# Debugging Tests for Model Explanations

Julius Adebayo
Massachusetts Institute of Technology
juliusad@mit.edu

Micheal Muelly
Stanford University
mmuelly@gmail.com

Ilaria Liccardi
Massachusetts Institute of Technology
liccardi@csail.mit.edu

Been Kim
Google AI
beenkim@google.com

## Abstract

We investigate whether post-hoc model explanations are effective for diagnosing model errors–model debugging. In response to the challenge of explaining a model's prediction, a vast array of explanation methods have been proposed. Despite increasing use, it is unclear if they are effective. To start, we categorize *bugs*, based on their source, into: *data, model, and test-time* contamination bugs. For several explanation methods, we assess their ability to: detect spurious correlation artifacts (data contamination), diagnose mislabeled training examples (data contamination), differentiate between a (partially) re-initialized model and a trained one (model contamination), and detect out-of-distribution inputs (test-time contamination). We find that the methods tested are able to diagnose a spurious background bug, but not conclusively identify mislabeled training examples. In addition, a class of methods, that modify the back-propagation algorithm are invariant to the higher layer parameters of a deep network; hence, ineffective for diagnosing model contamination. We complement our analysis with a human subject study, and find that subjects fail to identify defective models using attributions, but instead rely, primarily, on model predictions. Taken together, our results provide guidance for practitioners and researchers turning to explanations as tools for model debugging.

## 1 Introduction

Diagnosing and fixing model errors–model debugging–remains a longstanding machine learning challenge (Cadamuro et al., 2016; Carenini et al., 1994; Cawsey, 1991, 1993; Chakarov et al., 2016; Sculley et al., 2015; Zinkevich Martin, 2020). Model debugging is increasingly important as automated systems, with learned components, are being tested in high-stakes settings (Bhatt et al., 2020; Herlocker et al., 2000; McKinney et al., 2020) where inadvertent errors can have devastating consequences. Increasingly, *explanations*–artifacts derived from a trained model with the primary goal of providing insights to an end-user–are being used as debugging tools for models assisting healthcare providers in diagnosis across several specialties (Cai et al., 2019; Sayres et al., 2019; Wen et al., 2017). Despite a vast array of explanation methods and increased use for debugging, little guidance exists on method effectiveness. For example, should an explanation work equally well for diagnosing mislabeled training samples and detecting spurious correlation artifacts? Should an explanation that is sensitive to model parameters also be effective for detecting domain shift?

Consequently, we ask and address the following question:

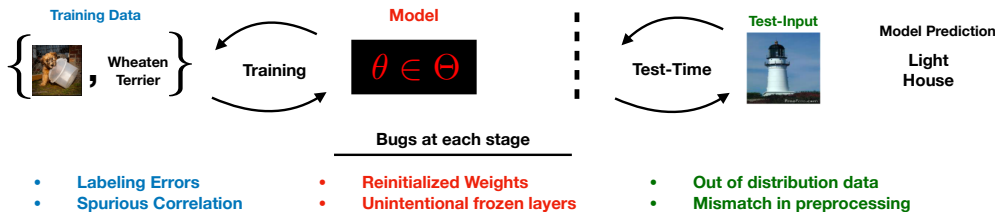*which explanation methods are effective for which classes of model bugs?*



Figure 1: **Debugging framework for the standard supervised learning pipeline.** Schematic of the standard supervised learning pipeline along with examples of bugs that can occur at each stage of the pipeline. The categorization captures defects that can occur with the training data, model, and at test-time. We term these: *data, model*, and *test-time contamination tests*.

To address this question, we make the following contributions:

1. **Bug Categorization.** We categorize bugs, based on the source of the defect leading to the bug, in the supervised learning pipeline (see Figure 1) into three classes: *data, model, and test-time* contamination. These contamination classes capture defects in the training data, model specification and parameters, and with the input at test-time.

2. **Empirical Assessment.** We conduct comprehensive control experiments to assess several feature attribution methods against 4 bugs: 'spurious correlation artifact', mislabelled training examples, re-initialized weights, and out-of-distribution (OOD) shift.

3. **Insights.** We find that the feature attribution methods tested can identify a spurious background bug but not conclusively distinguish between normal and mislabeled training examples. In addition, attribution methods that derive relevance by modifying the back-propagation computation via 'positive aggregation' (see Section 4) are invariant to the higher layer parameters of a deep neural network (DNN) model. Finally, we find that in specific settings, attributions for out-of-distribution examples are visually similar to attributions of these examples but with an 'in-domain' model, suggesting that debugging solely based on visual inspection might be misleading.

4. **Human Subject Study.** We conduct a 54-person IRB-approved study to assess whether end-users can identify defective models with attributions. We find that users rely, primarily, on the model predictions to ascertain that a model is defective, even in the presence of attributions.

**Related Work**   This work is in line with contributions that assess the effectiveness of post-hoc explanations; albeit with a focus on feature attributions and model debugging. Our bug categorization incorporates previous use of explanations for diagnosing spurious correlation (Kim et al., 2018; Meng et al., 2018; Ribeiro et al., 2016), domain mismatch, and mislabelled examples (Koh and Liang, 2017). Correcting bugs can also be achieved by penalizing feature attributions during training (Erion et al., 2019; Rieger et al., 2020; Ross et al., 2017) or clustering (Lapuschkin et al., 2019).

The dominant evaluation approach involves input perturbation (Montavon et al., 2017c; Samek et al., 2017), which can be combined with retraining (Hooker et al., 2019). However, Tomsett et al.

(2020) showed that input perturbation produces inconsistent quality rankings. Meng et al. (2018) propose manipulations to the training data along with a suite of metrics for assessing explanation quality. The data and model contamination categories recover the 'sanity checks' of Adebayo et al. (2018b). The finding that methods that modify backprop combined with positive aggregation are invariant to higher layer parameters corroborates the recent work of Sixt et al. (2019) along with previous evidence by Nie et al. (2018) and Mahendran and Vedaldi (2016).

The gold standard for assessing the effectiveness of an explanation is a human subject study (Doshi-Velez and Kim, 2017). Poursabzi-Sangdeh et al. (2018) manipulate the features of a linear model trained to predict housing prices to assess how well end-users can identify model mistakes. More recently, human subject tests of feature attributions have cast doubt on the ability of these approaches to help end-users debug erroneous predictions and improve human performance on downstream tasks (Chu et al., 2020; Shen and Huang, 2020). In a cooperative setting, Lai and Tan (2019) find that the humans exploit label information and Feng and Boyd-Graber (2019) demonstrate how to assess explanations in a natural language setting. Similarly, Alqaraawi et al. (2020) find that the LRP explanation method (see Section 2.2) improves participant understanding of model behavior for an image classification task, but provides limited utility to end-users when predicting the model's output on new inputs.

Feature attributions can be easily manipulated, providing evidence for a collective 'weakness' of current approaches (Ghorbani et al., 2019; Heo et al., 2019; Lakkaraju and Bastani, 2020; Slack et al., 2020). While susceptibility is an important issue, our work focuses on providing insights when model bugs are 'unintentionally' created.

## 2 Bug Characterization, Explanation Methods, & User Study

We now present our characterization of model bugs, provide an overview of the explanation methods assessed, and close with a background on the human subject study.[1]

### 2.1 Characterizing Model Bugs.

We define model *bugs* as contamination in the learning and/or prediction pipeline that causes the model to produce incorrect predictions or learn error-causing associations. We restrict our attention to the standard supervised learning setting, and categorize bugs based on their source. Given input-label pairs, $\{x_i, y_i\}_i^n$, where $x \in \mathcal{X}$ and $y \in \mathcal{Y}$, a classifier's goal is to learn a function, $f_\theta : \mathcal{X} \to \mathcal{Y}$, that generalizes. $f_\theta$ is then used to predict test examples, $x_{\text{test}} \in \mathcal{X}$, as $y_{\text{test}} = f_\theta(x_{\text{test}})$. Given a loss function $L$, and model parameter, $\theta$, for a model family, we provide a categorization of bugs as model, data and test-time contamination:

$$\text{Learning:} \quad \underset{\theta}{\arg\min} \quad L(\ \overbrace{(X_{\text{train}}, Y_{\text{train})}}^{\text{Data Contamination}}\ ; \theta);$$
$$\underbrace{\phantom{xxx}}_{\text{Model Contamination}}$$

$$\text{Prediction: } y_{\text{test}} = f_\theta(\ \overbrace{x_{\text{test}}}^{\text{Test-Time Contamination}}\ ).$$

---

[1] We will provide code to replicate our findings at: https://github.com/adebayoj/explaindebug.git.

| Bug Category | Specific Examples tested | Formalization |
|---|---|---|
| Data Contamination | Spurious Correlation | $\arg\min_{\theta} L(X_{\text{spurious artifact}}, Y_{\text{train}}; \theta)$ |
| | Labelling Errors | $\arg\min_{\theta} L(X_{\text{train}}, Y_{\text{wrong label}}; \theta)$ |
| Model Contamination | Initialized Weights | $f_{\theta_{\text{init}}}(x_{\text{test}})$ |
| Test-Time Contamination | Out of Distribution (OOD) | $f_{\theta}(x_{\text{OOD}})$ |

Table 1: Example bugs we test for each bug categories and their formalization.

**Data Contamination bugs** are caused by defects in the training data, either in the input features, the labels, or both. For example, a few incorrectly labeled data can cause the model to learn wrong associations. Another bug is a spurious correlation training signal. For example, consider an object classification task where all birds appear against a blue sky background. A model trained on this dataset can learn to associate blue sky backgrounds with the bird class; such dataset biases frequently occur in practice (Badgeley et al., 2019; Ribeiro et al., 2016).

**Model Contamination bugs** are caused by defects in the model parameters. For example, bugs in the code can cause accidental re-initialization of model weights.

**Test-Time Contamination bugs** are caused by defects in test-input, including domain shift or pre-processing mismatch at test time.

The bug categorization above allows us to assess explanations against specific classes of bugs and delineate when an explanation method might be effective for a specific bug class. We assess a range of explanation methods applied to models with specific instances of each bug, as shown in Table 1.

## 2.2 Explanation Methods

We focus on *feature attribution methods* that provide a 'relevance' score for the dimensions of input towards a model's output. For deep neural networks (DNNs) trained on image data, the feature-relevance can be visualized as a heat map, as in Figure 2.

An attribution functional, $E : \mathcal{F} \times \mathbb{R}^d \times \mathbb{R} \to \mathbb{R}^d$, maps the input, $x_i \in \mathbb{R}^d$, the model, $F \in \mathcal{F}$, output, $F_k(x)$, to an attribution map, $M_{x_i} \in \mathbb{R}^d$. Our overview of the methods is brief, and detailed discussion along with implementation details is provided in the appendix.
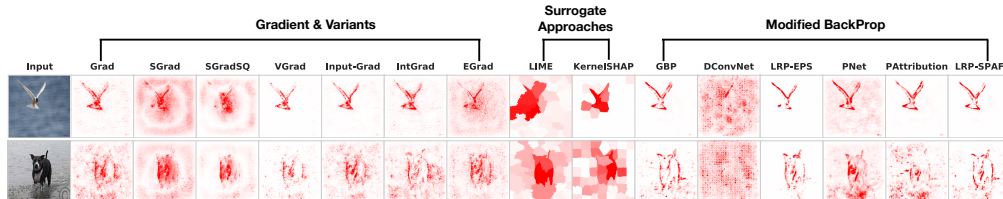


Figure 2: **Attribution Methods Considered.** The Figure shows feature attributions for two inputs for a CNN model trained to distinguish between birds and dogs.

**Gradient (Grad) & Variants.** We consider: 1) The *Gradient (Grad)* (Baehrens et al., 2010; Simonyan et al., 2014) map, $|\nabla_{x_i} F_i(x_i)|$; 2) *SmoothGrad (SGrad)* (Smilkov et al., 2017), $E_{\text{sg}}(x) = \frac{1}{N} \sum_{i=1}^{N} \nabla_{x_i} F_i(x_i + n_i)$ where $n_i$ is Gaussian noise; 3) *SmoothGrad Squared (SGradSQ)* (Hooker et al., 2019), the element-wise square of SmoothGrad; 4) *VarGrad (VGrad)* (Adebayo et al., 2018a),

4

the variance analogue of SmoothGrad; & 5) *Input-Grad* (Shrikumar et al., 2016) the element-wise product of the gradient and input $|\nabla_{x_i} F_i(x_i)| \odot x_i$. We also consider: 6) *Integrated Gradients (IntGrad)* (Sundararajan et al., 2017) which sums gradients along an interpolation path from the "baseline input", $\bar{x}$, to $x_i$: $M_{\text{IntGrad}}(x_i) = (x_i - \bar{x}) \times \int_0^1 \frac{\partial S(\bar{x} + \alpha(x_i - \bar{x}))}{\partial x_i} d\alpha$; and 7) *Expected Gradients (EGrad)* (Erion et al., 2019) which computes IntGrad but with a baseline input that is an expectation over the training set.

**Surrogate Approaches.** LIME (Ribeiro et al., 2016) and SHAP (Lundberg and Lee, 2017) locally approximate $F$ around $x_i$ with a simple function, $g$, that is then interpreted. SHAP provides a tractable approximation to the Shapley value (Shapley, 1988).

**Modified Back-Propagation.** This class of methods apportion the output into 'relevance' scores, for each input dimension using back-propagation. *DConvNet* (Zeiler and Fergus, 2014) & *Guided Back-propagation (GBP)* (Springenberg et al., 2014) modify the gradient for a ReLU unit. *Layer-wise relevance propagation (LRP)* (Bach et al., 2015a; Binder et al., 2016; Kohlbrenner et al., 2020; Montavon et al., 2017a) methods specify 'relevance' rules that modify the back-propagation. We consider *LRP-EPS*, and *LRP sequential preset-a-flat (LRP-SPAF)*. *PatternNet (PNet)* and *Pattern Attribution (PAttribution)* (Kindermans et al., 2018a) decompose the input into signal and noise components, and back-propagate relevance for the signal component.

**Attribution Comparison.** We measure visual and feature ranking similarity with the structural similarity index (SSIM) (Wang et al., 2004) and Spearman rank correlation metrics, respectively.

## 2.3   Overview of Human Subject Study

**Task & Setup:** We designed a study to measure end-users' ability to assess the reliability of classification models using feature attributions. Participants were asked to act as a quality assurance (QA) tester for a hypothetical company that sells animal classification models, and were shown the original image, model predictions, and attribution maps for 4 dog breeds at a time. They then rated how likely they are to recommend the model for sale to external customers using a 5 point-Likert scale, and a rationale for their decision. Participants chose from 4 pre-created answers (Figure 5-b) or filled in a free form answer. Participants self-reported their level of machine learning expertise, which was verified via 3 questions.

**Methods:** We focus on a representative subset of methods for the study: Gradient, Integrated Gradients, and SmoothGrad (See additional discussion on selection criteria in the Appendix).

**Bugs:** We tested the bugs described in Table 1 along with a model with no bugs.

## 3   Debugging Data Contamination

**Overview.** We assess whether feature attributions can detect spurious training artifacts and mis-labelled training examples. Spurious artifacts are signals that encode or correlate with the label in the training set but provide no meaningful connection to the data generating process. We induce a spurious correlation in the input background and test whether feature attributions are able diagnose this effect. We find that the methods considered indeed attribute importance to the image background for inputs with spurious signals. However, despite visual evidence in the attributions,

participants in the human subject study were unsure about model reliability for the spurious model condition; hence, did not out-rightly reject the model.
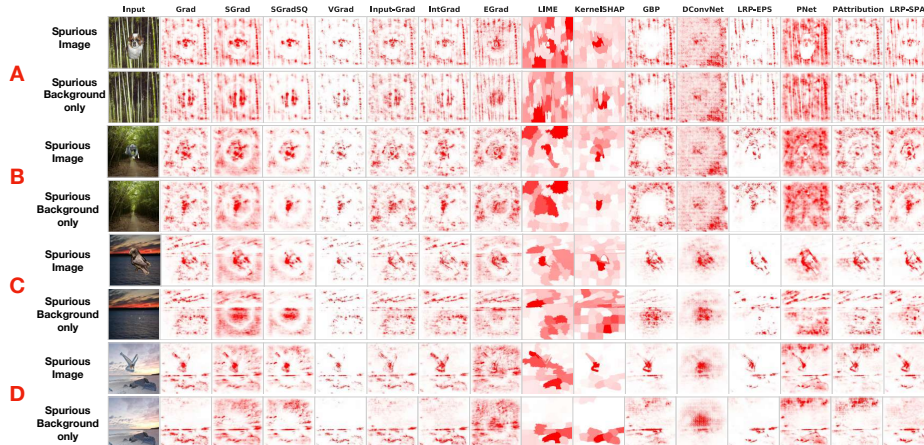


Figure 3: **Feature Attributions for Spurious Correlation Bugs.** Figure shows attributions for 4 inputs for the BVD-CNN trained on spurious data. A & B show two dog examples, and C & D are bird examples. The first row shows the input (dog or bird) on a spurious background. The second row shows the attributions of only the spurious background. Notably, we observe that the feature attribution methods place emphasis on the background. See Table 2 for metrics.

For mislabeled examples, we compare attributions for a training input derived from: 1) a model where this training input had the correct label, and 2) the same model settings but trained with this input mislabeled. If the attributions under these two settings are similar, then such a method is unlikely to be useful for identifying mislabeled examples. We observe that attributions for mislabeled examples, across all methods, show visual similarity.

**General Data and Model Setup.** We consider a birds-vs-dogs binary classification task. We use dog breeds from the Cats-v-Dogs dataset (Parkhi et al., 2012a) and Bird species from the Caltech-UCSD dataset (Wah et al., 2011). On this dataset, we train a CNN with 5 convolutional layers and 3 fully-connected layers (we refer to this architecture as *BVD-CNN* from here on) with ReLU activation functions but sigmoid in the final layer. The model achieves a test accuracy of 94-percent.

## 3.1   Spurious Correlation Training Artifacts

**Spurious Bug Implementation.** We introduce spurious correlation by placing all birds onto one of the sky backgrounds from the places dataset (Zhou et al., 2017), and all dogs onto a bamboo forest background (see Figure 3). BVD-CNN trained on this data achieves a 97 percent accuracy on a sky-vs-bamboo forest test set (without birds or dogs) indicating that the model indeed learned the spurious association.



Figure 4: Ground Truth Attribution for Spurious Correlation.

**Results.**   To quantitatively measure whether attribution methods reflect the spurious background, we compare attributions to two ground truth masks (GT-1 & GT-2). As shown in Figure 4, we consider an ideal mask that apportions all relevance to the background and none to

the object part. Next, we consider a relaxed version that weights the first ground truth mask by the attribution of a spurious background without the object. In Table 2, we report SSIM comparison scores across all methods for both ground-truth masks. For *GT-2*, scores range from a minimum of 0.78 to maximum of 0.98; providing evidence that the attributions identify the spurious background signal. We find similar evidence for *GT-1*.

| Metric | Grad | SGrad | SGradSQ | VGrad | Input-Grad | IntGrad | EGrad | LIME | KernelSHAP | GBP | DConvNet | LRP-EPS | PNet | PAttribution | LRP-SPAF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SSIM-GT1 | 0.62 | 0.63 | 0.063 | 0.075 | 0.69 | 0.7 | 0.63 | 0.59 | 0.58 | 0.58 | 0.6 | 0.65 | 0.51 | 0.44 | 0.69 |
| SSIM-GT1 (SEM) | 0.012 | 0.013 | 0.0077 | 0.0089 | 0.019 | 0.019 | 0.024 | 0.021 | 0.037 | 0.019 | 0.017 | 0.039 | 0.036 | 0.018 | 0.028 |
| SSIM-GT2 | 0.83 | 0.83 | 0.89 | 0.98 | 0.85 | 0.85 | 0.85 | 0.88 | 0.78 | 0.82 | 0.83 | 0.85 | 0.85 | 0.8 | 0.85 |
| SSIM-GT2 (SEM) | 0.013 | 0.013 | 0.02 | 0.0024 | 0.013 | 0.012 | 0.012 | 0.011 | 0.044 | 0.013 | 0.013 | 0.012 | 0.013 | 0.018 | 0.013 |

Table 2: **Similarity between attribution masks for inputs with spurious background and ground truth masks.** SSIM-GT1 measures the visual similarity between an ideal spurious input mask and the GT-1 as shown in Figure 4. SSIM-GT2 measures visual similarity for the GT-2. We also include the standard error of the mean (SEM) for each metric, which was computed across 190 inputs. To calibrate this metric, the mean SSIM between a randomly sampled Gaussian attribution and the spurious attributions which is: $3e^{-06}$.
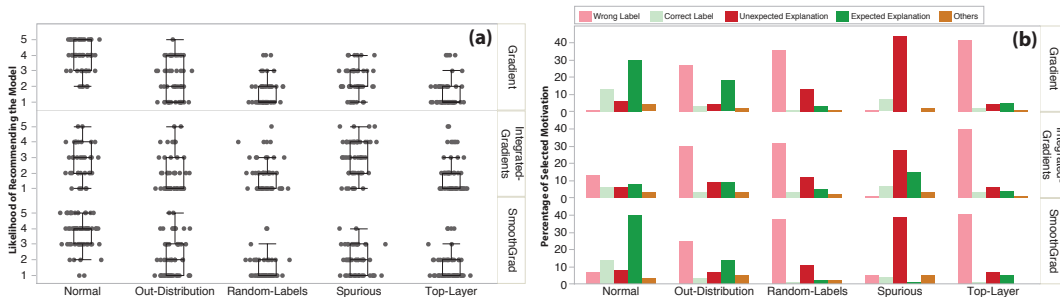


Figure 5: **A: Participant Responses from User Study.** Box plot of participants responses for 3 attribution methods: *Gradient, SmoothGrad, and Integrated Gradients*, and 5 model conditions tested. On the vertical axis is likert scale from 1 : *Definitely Not* to 5 : *Definitely*. Participants were instructed to select 'Definitely' if they deemed the dog-breed classification model ready to be sold to customers. **B: Motivation for Selection.** Participants' selected motivations (%) for the recommendation made. As shown in the legend, users could select one of 4 options or insert an open-ended response.

**Insights from Human Subject Study: users are uncertain.** Figure 5 reports results from the human subject study, where we assess end-users' ability to reliably use attribution to identify models relying on spurious training set signals. For a normal model, the median Likert scores are 4, 4, 3 for Gradient, SmoothGrad, and Integrated Gradients respectively. Selecting a likert score of 1 means a user will 'definitely not' recommend the model, while 5 means they will 'definitely' recommend the model. Consequently, users adequately rate a normal model. In addition, 30 and 40 percent (See Figure 5-Right) of participants, for Gradient and SmoothGrad respectively, indicate that the attributions for a normal model 'highlighted the part of the image that they expected it to focus on'.

For the 'spurious model', the Likert scores show a wider range. While the median scores are 2, 2, 3 for Gradient, SmoothGrad, and Integrated Gradients respectively, some end-users still recommend this model. For each attribution type, a majority of end-users indicate that the attribution 'did not highlight the part of the image that I expected it to focus on'. Despite this, end-users do not convincingly reject the spurious model like they do for the other bug conditions. These results

suggest that the ability of an attribution method to diagnose spurious correlation might not carry over to reliable decision making.

## 3.2    Mislabelled Training Examples

**Bug Implementation.** We train a BVD-CNN model on a birds-vs-dogs dataset where 10 percent of training samples have their labels flipped. The model achieves a 93.2, 91.7, 88 percent accuracy on the training, validation, and test sets.
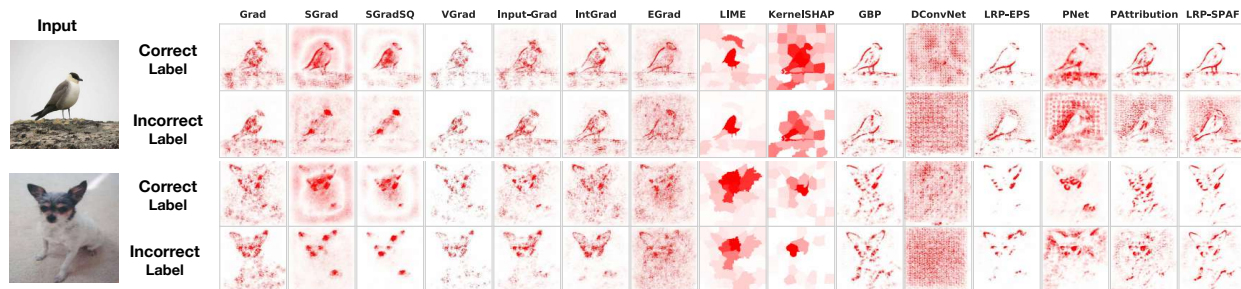


Figure 6: **Diagnosing Mislabelled Training Examples.** The Figure shows two training inputs along with feature attributions for each method. The correct label row corresponds to feature attributions derived from a model with the correct label in the training set. The incorrect-label row shows feature attributions derived from a model with the wrong label in the training set. We see that the attributions under both settings are visually similar.

**Results.** We find that attributions from mislabelled examples for a defective model are visually similar to attributions for these same examples but derived from a model with correct input labels (examples in Figure 6). We find that the SSIM between the attributions of a correctly labeled instance, and the corresponding incorrectly labeled instance, are in the range $0.73 - 0.99$ for all methods tested. These results indicate that the attribution methods tested might be ineffective for identifying mislabelled examples. We refer readers to Section I of the Appendix for visualizations of additional examples.

**Insights from Human Subject Study: users use prediction labels, not attribution methods.** In contrast to the spurious setting, participants reject mislabelled examples with median Likert scores 1, 2, and 1 for Gradient, SmoothGrad, and Integrated Gradients respectively. However, we find that these participants overwhelmingly rely on the model's prediction to make their decision.

## 4    Debugging Model Contamination

We next evaluate bugs related to model parameters. Specifically, we consider the setting where the weights of a model are accidentally re-initialized prior to prediction (Adebayo et al., 2018b). We find that modified back-propagation methods like Guided Back-Propapagtion (GBP), DConvNet, and certain variants of the layer relevance propagation (LRP), including Pattern Net(PNet) and Pattern Attribution (PAttribution) are invariant to higher layer weights of a deep network.

**Bug Implementation.** We instantiate this bug on a pre-trained VGG-16 model on Imagenet (Russakovsky et al., 2015). Similar to Adebayo et al. (2018b), we re-initialize the weights of the model

starting at the top layer, successively, all the way to the first layer. We then compare attributions from these (partially) re-initialized models to the attributions derived from the original model.
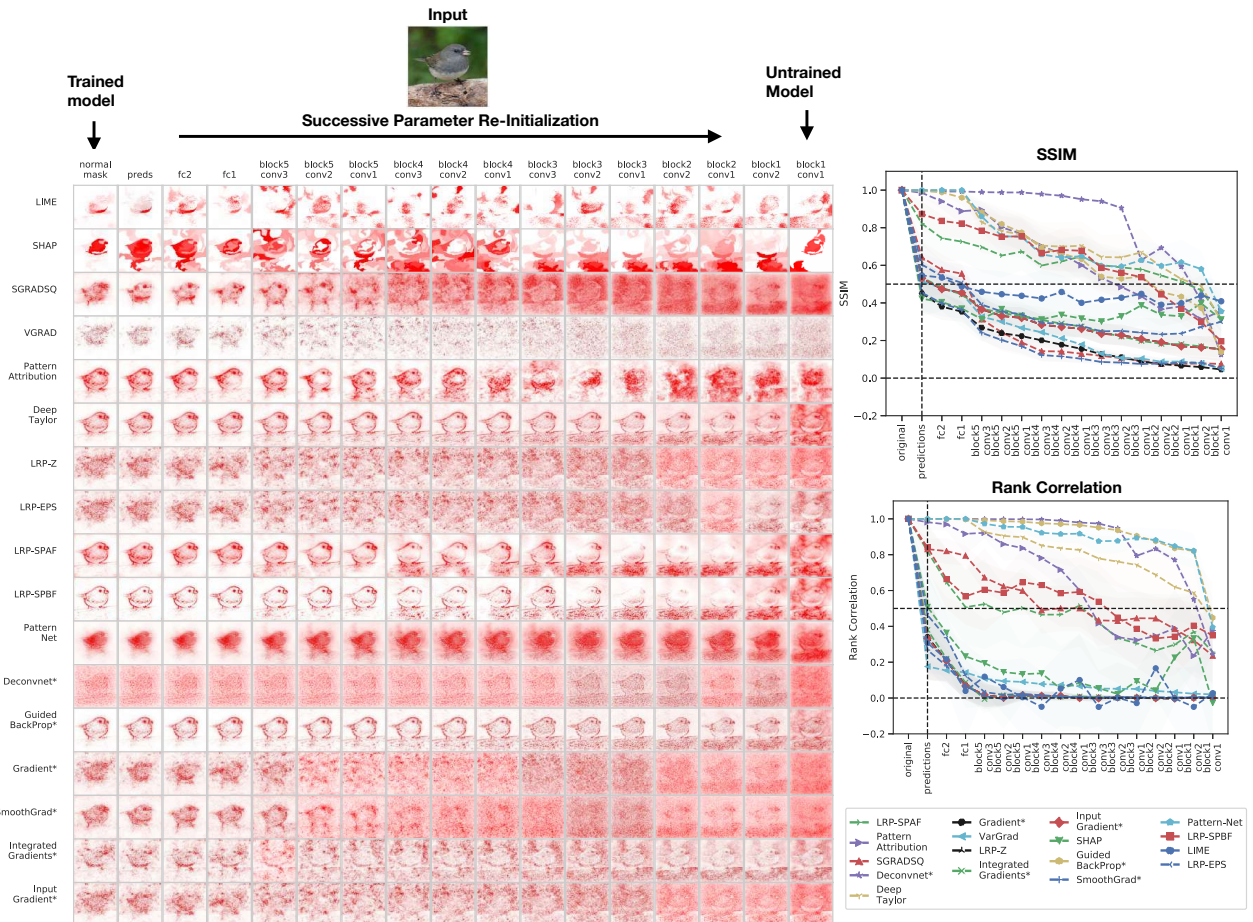


Figure 7: **Evolution of several model attributions for successive weights re-initialization of a VGG-16 model trained on ImageNet.** Qualitative results (left) and quantitative results (right). The last column in qualitative results corresponds to a network with completely re-initialized weights.

**Results: modified back-propagation methods are parameter invariant.** As seen in Figure 7, the class of modified back-propagation methods, including Guided BackProp, Deconvnet, DeepTaylor, PatternNet, Pattern Attribution, and LRP-SPAF are visually and quantitatively invariant to higher layer parameters of the VGG-16 model. This finding corroborates prior results for Guided Backprop and Deconvnet (Adebayo et al., 2018b; Mahendran and Vedaldi, 2016; Nie et al., 2018). These results also support the recent findings of Sixt et al. (2019), who prove that these modified back-propagation approaches produce attributions that converge to a rank-1 matrix.

**Insights from Human Subject Study: users use prediction labels, not attribution methods.** We observe that participants conclusively reject a model whose top layer has been re-initialized purely based on the classification labels, and rarely based on wrong attributions. (Figure 5).

9

# 5 Debugging Test-Time Contamination

A model is at risk of providing errant predictions when given inputs that have distributional characteristics different from the training set. To assess the ability of feature attributions to diagnose domain shift, we compare attributions derived, for a given input, from an *in-domain model* with those derived from *out-of-domain model*. For example, we compare the attribution for an MNIST digit, derived from a model trained on MNIST, to an attribution for the same digit, but derived from a model trained on Fashion MNIST, ImageNet, and a birds-vs-dogs model. We find visual similarity for certain settings: for example, feature attributions for a Fashion MNIST input derived from a VGG-16 model trained on ImageNet are visually similar to attributions for the same input on a model trained on Fashion MNIST. However, the quantitative ranking of the input dimensions are widely different.



Figure 8: **Fashion MNIST OOD on several models.** The first row shows feature attributions on a model trained on Fashion MNIST. In the subsequent rows, we show feature attributions for the same input on an MNIST model, BVD-CNN model trained on birds-vs-dogs, and lastly, a pre-trained VGG-16 model on ImageNet.

| Metric | Grad | SGrad | SGradSQ | VGrad | Input-Grad | IntGrad | EGrad | LIME | KernelSHAP | GBP | DConvNet | LRP-EPS | PNet | PAttribution | LRP-SPAF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SSIM (FMNIST → MNIST Model) | 0.7 | 0.54 | 0.49 | 0.92 | 0.71 | 0.69 | 0.71 | 0.46 | 0.41 | 0.81 | 0.5 | 0.77 | 0.58 | 0.77 | 0.66 |
| SEM | 0.0093 | 0.012 | 0.016 | 0.0047 | 0.01 | 0.015 | 0.01 | 0.02 | 0.024 | 0.014 | 0.01 | 0.02 | 0.026 | 0.009 | 0.03 |
| RK (FMNIST → MNIST Model) | 0.0013 | 8.8e-4 | 0.37 | 0.37 | 0.0021 | -0.003 | 0.002 | -0.01 | 0.034 | 0.51 | 0.027 | 0.011 | -0.14 | 0.0082 | 0.12 |
| SEM | 0.0016 | 0.0032 | 0.026 | 0.029 | 0.002 | 0.002 | 0.002 | 0.04 | 0.028 | 0.014 | 6e-4 | 0.0034 | 0.027 | 0.0026 | 0.023 |
| SSIM (FMNIST → BVD-CNN) | 0.7 | 0.5 | 0.55 | 0.93 | 0.72 | 0.7 | 0.72 | 0.72 | 0.72 | 0.82 | 0.63 | 0.79 | 0.53 | 0.36 | 0.66 |
| SEM | 0.0083 | 0.011 | 0.013 | 0.0045 | 0.009 | 0.013 | 0.009 | 0.009 | 0.01 | 0.009 | 0.014 | 0.019 | 0.03 | 0.025 | 0.035 |
| RK (FMNIST → BVD-CNN) | 0.0012 | 0.0078 | 0.43 | 0.25 | 0.0002 | 0.002 | 0.00025 | 0.18 | 0.067 | 0.078 | -0.05 | -0.013 | 0.25 | -0.0095 | 0.044 |
| SEM | 8.5e-4 | 0.0017 | 0.009 | 0.011 | 0.0007 | 0.001 | 0.0007 | 0.04 | 0.034 | 0.008 | 0.0011 | 0.0027 | 0.045 | 0.0023 | 0.02 |
| SSIM (FMNIST → VGG-16 ImageNet) | 0.57 | 0.46 | 0.5 | 0.87 | 0.64 | 0.67 | 0.64 | 0.5 | 0.38 | 0.8 | 0.36 | 0.64 | 0.66 | 0.12 | 0.2 |
| SEM | 0.012 | 0.011 | 0.015 | 0.0056 | 0.01 | 0.015 | 0.011 | 0.015 | 0.03 | 0.009 | 0.01 | 0.02 | 0.018 | 0.0049 | 0.024 |
| RK (FMNIST → VGG-16 ImageNet) | -0.0023 | -0.0098 | -0.0097 | 0.028 | -0.0025 | -0.0017 | -0.0025 | 0.005 | -0.045 | 0.25 | 0.03 | 0.0045 | 0.32 | 0.066 | 0.14 |
| SEM | 0.0017 | 0.0025 | 0.02 | 0.018 | 0.0023 | 0.0016 | 0.002 | 0.033 | 0.024 | 0.004 | 0.0035 | 0.0018 | 0.034 | 0.0053 | 0.019 |

Table 3: **Test-time Explanation Similarity Metrics.** We observe visual similarity but no ranking similarity. We show each metric along with the standard error of the mean calculated for 190 examples. FMNIST → MNIST model means a comparison of FMNIST attributions for an FMNIST model with FMNIST attributions derived from *an MNIST model*. We present both SSIM and Rank correlation metrics.

**Bug Implementation.** We consider 4 dataset-model pairs: a BVD-CNN trained on MNIST, Fashion MNIST, the Birds-vs-dogs data, and lastly a VGG-16 model trained on ImageNet. We present results on Fashion MNIST. Concretely, we compare 1) feature attributions of Fashion MNIST examples derived from a model trained on Fashion MNIST, and 2) feature attributions of Fashion MNIST examples for models trained on MNIST, the birds-vs-dogs dataset, and ImageNet.

**Results.** As shown in Figure 8, we observe visual similarity between in-domain Fashion MNIST attributions, and attributions for these samples on other models. As seen in Table 3, we observe visual similarity, particularly for the VGG-16 model on ImageNet, but essentially no correlation in feature ranking.

**Insights from Human Subject Study: users use prediction labels, not the attributions.**
For the domain shift study, we show participants attribution of dogs that were not used during training, and whose breeds differed from those that the model was trained to predict. We find that users do not recommend a model under this setting due to wrong prediction labels (Figure 5).

# 6    Discussion & Conclusion

Debugging machine learning models remains a challenging endeavor, and model explanations could be a useful tool in that quest. Even though a practitioner or a researcher may have a large class of explanation methods available, it is still unclear which methods are useful for what bug type. This work aims to address this gap by first, categorizing model bugs into: data, model, and test-time contamination bugs, then testing feature attribution methods, a popular explanation approach for DNNs trained on image data, against each bug type. Overall, we find that feature attribution methods are able to diagnose the spatial spurious correlation bug tested, but do not conclusively help to distinguish mislabelled examples for normal ones. In the case of model contamination, we find that certain feature attributions that perform positive aggregation while computing feature relevance with modified back-propagation produce attributions that are invariant to the parameters of the higher layers for a deep model. This suggests that these approaches might not be effective for diagnosing model contamination bugs. We also find that attributions of out-of-domain inputs are similar to attributions for these inputs on an in-domain model, which suggests caution when visually inspecting these explanations, especially for image tasks. We also conduct human subject tests to assess how well end-users can use attributions to assess model reliability. Here we find that the end-users relied, primarily, on model predictions for diagnosing model bugs.

**Limitations.**    Our findings come with certain limitations and caveats. The bug characterization presented only covers the standard supervised learning pipeline and might not neatly capture bugs that result from a combination of factors. We only focused on feature attributions: however, other methods such as approaches based on 'concept' activation (Kim et al., 2018), model representation dissection (Bau et al., 2017), and training point ranking (Koh and Liang, 2017; Pruthi et al., 2020; Yeh et al., 2018) might be more suited to the debugging tasks studied here. Indeed, initial exploration of the 'concept' activation method TCAV and training point ranking based on influence functions suggests that these approaches are promising (See Appendix for analysis). For the human subject experiments, our finding that the participants mostly relied on the labels instead of the feature attributions might be a consequence of the dog breed classification task. It is unclear whether participants would still rely of model predictions for tasks in which they have no expertise or prior knowledge.

The goal of this work is to provide guidance for researchers and practitioners seeking to use feature attributions for model debugging. We hope our findings can serve as a first step towards more rigorous approaches for assessing the utility of explanation methods.

# Acknowlegdements

# References

Julius Adebayo, Justin Gilmer, Ian Goodfellow, and Been Kim. Local explanation methods for deep neural networks lack sensitivity to parameter values. 2018a.

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*, pages 9525–9536, 2018b.

Maximilian Alber, Sebastian Lapuschkin, Philipp Seegerer, Miriam Hägele, Kristof T. Schütt, Grégoire Montavon, Wojciech Samek, Klaus-Robert Müller, Sven Dähne, and Pieter-Jan Kindermans. innvestigate neural networks! *CoRR*, abs/1808.04260, 2018. URL http://arxiv.org/abs/1808.04260.

Ahmed Alqaraawi, Martin Schuessler, Philipp Weiß, Enrico Costanza, and Nadia Berthouze. Evaluating saliency map explanations for convolutional neural networks: a user study. In *Proceedings of the 25th International Conference on Intelligent User Interfaces*, pages 275–285, 2020.

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE*, 10(7):e0130140, 07 2015a. doi: 10.1371/journal.pone.0130140. URL http://dx.doi.org/10.1371%2Fjournal.pone.0130140.

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015b.

Marcus A Badgeley, John R Zech, Luke Oakden-Rayner, Benjamin S Glicksberg, Manway Liu, William Gale, Michael V McConnell, Bethany Percha, Thomas M Snyder, and Joel T Dudley. Deep learning predicts hip fracture using confounding patient and healthcare variables. *NPJ digital medicine*, 2(1):1–10, 2019.

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

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6541–6549, 2017.

Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly, Yunhan Jia, Joydeep Ghosh, Ruchir Puri, José M. F. Moura, and Peter Eckersley. Explainable machine learning in deployment. In Mireille Hildebrandt, Carlos Castillo, Elisa Celis, Salvatore Ruggieri, Linnet Taylor, and Gabriela Zanfir-Fortuna, editors, *FAT* '20: Conference on Fairness, Accountability, and Transparency, Barcelona, Spain, January 27-30, 2020*, pages 648–657. ACM, 2020. doi: 10.1145/3351095.3375624. URL https://doi.org/10.1145/3351095.3375624.

Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, Klaus-Robert Müller, and Wojciech Samek. Layer-wise relevance propagation for neural networks with local renormalization layers. volume 9887 of *Lecture Notes in Computer Science*, pages 63–71. Springer Berlin / Heidelberg, 2016. doi: 10.1007/978-3-319-44781-0_8.

Gabriel Cadamuro, Ran Gilad-Bachrach, and Xiaojin Zhu. Debugging machine learning models. In *ICML Workshop on Reliable Machine Learning in the Wild*, 2016.

C Cai, E Reif, N Hegde, J Hipp, B Kim, D Smilkov, M Wattenberg, D Viegas, G Corrado, M Stumpe, et al. Human-centered tools for coping with imperfect algorithms during medical decision-making. chi conf. *Hum. Factors Comput. Syst*, 2019.

Giuseppe Carenini, Vibhu O Mittal, and Johanna D Moore. Generating patient-specific interactive natural language explanations. In *Proceedings of the Annual Symposium on Computer Application in Medical Care*, page 5. American Medical Informatics Association, 1994.

Alison Cawsey. Generating interactive explanations. In *AAAI*, pages 86–91, 1991.

Alison Cawsey. User modelling in interactive explanations. *User Modeling and User-Adapted Interaction*, 3(3):221–247, 1993.

Aleksandar Chakarov, Aditya Nori, Sriram Rajamani, Shayak Sen, and Deepak Vijaykeerthy. Debugging machine learning tasks. *arXiv preprint arXiv:1603.07292*, 2016.

Eric Chu, Deb Roy, and Jacob Andreas. Are visual explanations useful? a case study in model-in-the-loop prediction. *arXiv preprint arXiv:2007.12248*, 2020.

Ann-Kathrin Dombrowski, Maximillian Alber, Christopher Anders, Marcel Ackermann, Klaus-Robert Müller, and Pan Kessel. Explanations can be manipulated and geometry is to blame. In *Advances in Neural Information Processing Systems*, pages 13567–13578, 2019.

Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. 2017.

Gabriel Erion, Joseph D Janizek, Pascal Sturmfels, Scott Lundberg, and Su-In Lee. Learning explainable models using attribution priors. *arXiv preprint arXiv:1906.10670*, 2019.

Shi Feng and Jordan Boyd-Graber. What can ai do for me? evaluating machine learning interpretations in cooperative play. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, pages 229–239, 2019.

Amirata Ghorbani, Abubakar Abid, and James Y. Zou. Interpretation of neural networks is fragile. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 3681–3688. AAAI Press, 2019. doi: 10.1609/aaai.v33i01.33013681. URL https://doi.org/10.1609/aaai.v33i01.33013681.

Juyeon Heo, Sunghwan Joo, and Taesup Moon. Fooling neural network interpretations via adversarial model manipulation. In *Advances in Neural Information Processing Systems*, pages 2921–2932, 2019.

Jonathan L Herlocker, Joseph A Konstan, and John Riedl. Explaining collaborative filtering recommendations. In *Proceedings of the 2000 ACM conference on Computer supported cooperative work*, pages 241–250, 2000.

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*, pages 9737–9748, 2019.

Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Fei-Fei Li. Novel dataset for fine-grained image categorization: Stanford dogs. In *Proc. CVPR Workshop on Fine-Grained Visual Categorization (FGVC)*, volume 2, 2011.

Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International Conference on Machine Learning*, pages 2673–2682, 2018.

Pieter-Jan Kindermans, Kristof Schütt, Klaus-Robert Müller, and Sven Dähne. Investigating the influence of noise and distractors on the interpretation of neural networks. *arXiv preprint arXiv:1611.07270*, 2016.

Pieter-Jan Kindermans, Kristof T. Schütt, Maximilian Alber, Been Kim Klaus-Robert Müller, Dumitru Erhan, and Sven Dähne. Learning how to explain neural networks: Patternnet and patternattribution. *International Conference on Learning Representations*, 2018a. URL https://openreview.net/forum?id=Hkn7CBaTW.

Pieter-Jan Kindermans, Kristof T. Schütt, Maximilian Alber, Klaus-Robert Müller, Dumitru Erhan, Been Kim, and Sven Dähne. Learning how to explain neural networks: Patternnet and patternattribution. In *ICLR (Poster)*, 2018b. URL https://openreview.net/forum?id=Hkn7CBaTW.

Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 1885–1894. PMLR, 2017. URL http://proceedings.mlr.press/v70/koh17a.html.

Maximilian Kohlbrenner, Alexander Bauer, Shinichi Nakajima, Alexander Binder, Wojciech Samek, and Sebastian Lapuschkin. Towards best practice in explaining neural network decisions with lrp. In *Proceedings of the IEEE International Joint Conference on Neural Networks (IJCNN)*, 2020.

Vivian Lai and Chenhao Tan. On human predictions with explanations and predictions of machine learning models: A case study on deception detection. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 29–38, 2019.

Himabindu Lakkaraju and Osbert Bastani. " how do i fool you?" manipulating user trust via misleading black box explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 79–85, 2020.

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.

Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pages 4768–4777, 2017.

Aravindh Mahendran and Andrea Vedaldi. Salient deconvolutional networks. In *European Conference on Computer Vision*, pages 120–135. Springer, 2016.

Scott Mayer McKinney, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova, Hutan Ashrafian, Trevor Back, Mary Chesus, Greg C Corrado, Ara Darzi, et al. International evaluation of an ai system for breast cancer screening. *Nature*, 577(7788):89–94, 2020.

Qingjie Meng, Christian Baumgartner, Matthew Sinclair, James Housden, Martin Rajchl, Alberto Gomez, Benjamin Hou, Nicolas Toussaint, Jeremy Tan, Jacqueline Matthew, et al. Automatic shadow detection in 2d ultrasound. 2018.

Grégoire Montavon, Sebastian Bach, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 65:211–222, 2017a. doi: 10.1016/j.patcog.2016.11.008. URL http://dx.doi.org/10.1016/j.patcog.2016.11.008.

Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 65:211–222, 2017b.

Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 2017c.

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

Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2012a.

Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012b.

Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Vaughan, and Hanna Wallach. Manipulating and measuring model interpretability. *arXiv preprint arXiv:1802.07810*, 2018.

Garima Pruthi, Frederick Liu, Mukund Sundararajan, and Satyen Kale. Estimating training data influence by tracking gradient descent. *arXiv preprint arXiv:2002.08484*, 2020.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. ACM, 2016.

Laura Rieger, Chandan Singh, W James Murdoch, and Bin Yu. Interpretations are useful: penalizing explanations to align neural networks with prior knowledge. *International Conference on Machine Learning*, 2020.

Andrew Slavin Ross, Michael C. Hughes, and Finale Doshi-Velez. Right for the right reasons: Training differentiable models by constraining their explanations. In Carles Sierra, editor, *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, pages 2662–2670. ijcai.org, 2017. doi: 10.24963/ijcai.2017/371. URL https://doi.org/10.24963/ijcai.2017/371.

Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.

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, 2017.

Rory Sayres, Ankur Taly, Ehsan Rahimy, Katy Blumer, David Coz, Naama Hammel, Jonathan Krause, Arunachalam Narayanaswamy, Zahra Rastegar, Derek Wu, et al. Using a deep learning algorithm and integrated gradients explanation to assist grading for diabetic retinopathy. *Ophthalmology*, 126(4):552–564, 2019.

David Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-Francois Crespo, and Dan Dennison. Hidden technical debt in machine learning systems. In *Advances in neural information processing systems*, pages 2503–2511, 2015.

Lloyd S Shapley. A value for n-person games. *The Shapley value*, pages 31–40, 1988.

Hua Shen and Ting-Hao Huang. How useful are the machine-generated interpretations to general users? a human evaluation on guessing the incorrectly predicted labels. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 8, pages 168–172, 2020.

Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*, 2016.

Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In Yoshua Bengio and Yann LeCun, editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Workshop Track Proceedings*, 2014. URL http://arxiv.org/abs/1312.6034.

Leon Sixt, Maximilian Granz, and Tim Landgraf. When explanations lie: Why modified bp attribution fails. *arXiv preprint arXiv:1912.09818*, 2019.

Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 180–186, 2020.

Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.

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

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3319–3328, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL http://proceedings.mlr.press/v70/sundararajan17a.html.

Richard Tomsett, Dan Harborne, Supriyo Chakraborty, Prudhvi Gurram, and Alun D. Preece. Sanity checks for saliency metrics. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 6021–6029. AAAI Press, 2020. URL https://aaai.org/ojs/index.php/AAAI/article/view/6064.

C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.

Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.

Gezheng Wen, Brenda Rodriguez-Niño, Furkan Y Pecen, David J Vining, Naveen Garg, and Mia K Markey. Comparative study of computational visual attention models on two-dimensional medical images. *Journal of Medical Imaging*, 4(2):025503, 2017.

Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point selection for explaining deep neural networks. In *Advances in neural information processing systems*, pages 9291–9301, 2018.

Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Suggala, David I Inouye, and Pradeep K Ravikumar. On the (in) fidelity and sensitivity of explanations. In *Advances in Neural Information Processing Systems*, pages 10965–10976, 2019.

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

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.

Zinkevich Martin. Rules of Machine Learning: Best Practices for ML Engineering. `http://martin.zinkevich.org/rules_of_ml/rules_of_ml.pdf`, 2020. Online; accessed 10 January 2020.

**Part**

# Appendix

## Table of Contents

## A   Additional Discussion of Bug Categorization Setup

In this section, we provide two additional formalization of bugs beyond those discussed in the main document: Frozen layers and Pre-Processing Mismatch. In the case of frozen layers, this is a bug whereby one of the layers of a deep network is accidentally kept fixed during training. In the case of pre-processing mismatch, a test-input can be pre-processed under settings that are different from those used for the training data. In such a case, it is possible the the model will produce wrong predictions due to the mismatch in input pre-processing. We formalize both these bugs in Table 4.

The examples presented in Table 4 are specific instantiation of bugs based on the categorization that we present; as expected, several other examples can be proposed under the same categorization.

## B   Detailed Overview of Attribution Methods

We now provide a detailed discussion of the explanation methods that we assess in this work. For each method we also provide the hyper-parameter values that we use in each case. First,

**Debugging Framework for Supervised Learning**

Figure 9: We show the debugging schematic here.

| Bug Class | Specific Examples | Formalization |
|---|---|---|
| Data Contamination | Spurious Correlation | $\arg\min_{\theta} L(X_{\text{spurious artifact}}, Y_{\text{train}}; \theta)$ |
| Data Contamination | Labelling Errors | $\arg\min_{\theta} L(X_{\text{train}}, Y_{\text{wrong label}}; \theta)$ |
| Model Contamination | Initialized Weights | $f_{\theta\text{init}}(x_{\text{test}})$ |
| Model Contamination | Frozen Layers | $\arg\min_{\theta\text{frozen}} L(X_{\text{train}}, Y_{\text{wrong label}}; \theta)$ |
| Test-Time Contamination | Out of Distribution (OOD) | $f_{\theta}(x_{\text{OOD}})$ |
| Test-Time Contamination | Pre-Processing Mismatch (PM) | $f_{\theta}(x_{\text{PM}})$ |

Table 4: **Example bugs for different contamination classes.** We show different examples of bugs under the different contamination classes: data, model and test-time. We also formalize these bugs for the traditional supervised learning setting.

we re-implemented all methods in a single code base. We then benchmark our implementation with the public open source implementations of these methods. Ultimately, we used the public implementation of the these methods as released by their authors. For gradient based methods and the corresponding LRP variants, we use the innvestigate python package. For LIME we used the publicly available LIME-IMAGE package. For Kernel Shap and Expected Gradients, we used the public SHAP package.

Recall that an attribution functional, $E : \mathcal{F} \times \mathbb{R}^d \times \mathbb{R} \to \mathbb{R}^d$, maps the input, $x_i \in \mathbb{R}^d$, the model, $F$, output, $F_k(x)$, to an attribution map, $M_{x_i} \in \mathbb{R}^d$.

**Gradient (Grad) & Variants.** We consider:

- The *Gradient (Grad)* (Baehrens et al., 2010; Simonyan et al., 2014) map, $|\nabla_{x_i} F_i(x_i)|$. The gradient map is a key primitive upon which several other methods are derived. Overall, the gradient map quantifies the sensitivity of the output, typically a logit score for DNNs trained for classification, to each dimension of the input.

- *SmoothGrad (SGrad)* (Smilkov et al., 2017), $E_{\text{sg}}(x) = \frac{1}{N} \sum_{i=1}^{N} \nabla_{x_i} F_i(x_i + n_i)$ where $n_i$ is sampled according to a random Gaussian noise. we considered 50 noisy inputs, selected the standard deviation of the noise to be $0.15 *$ input range. Here input range refers to the

## Overview of Methods



Figure 10: We show five input examples from the birds-vs-dogs dataset along with the attributions for a normally trained model.

difference between the maximum and minimum value in the input. For the models considered, these inputs are typically normalized to be in the range $[-1, 1]$.

- *SmoothGrad Squared (SGradSQ)* (Hooker et al., 2019), the element-wise square of Smooth-Grad: $E_{\text{SGradSQ}}(x) = E_{\text{sg}}(x) \odot E_{\text{sg}}(x)$. We set the parameters of SmoothGrad as discussed in the previous item.

- *VarGrad (VGrad)* (Adebayo et al., 2018a), the variance analogue of SmoothGrad: $E_{\text{VGrad}}(x) = \mathbb{V}(\tilde{x})$, where $\mathbb{V}$ is the variance operator, $\tilde{x} = \nabla_{x_i} F_i(x_i + n_i)$, and $n_i$ is sampled according to a random Gaussian noise. Here we consider 50 noise input examples, set the noise parameter as a Gaussian with mean zero, and noise scale similar to SmoothGrad: $0.15 * $ input range.

- *Input-Grad* (Shrikumar et al., 2016) the element-wise product of the gradient and input $|\nabla_{x_i} F_i(x_i)| \odot x_i$. Several other methods approximate and have been shown to be equivalent to this product. For example, for a DNN with all ReLU activations, LRP-Z, a variant of layer-wise relevance propagation that we discuss in the modified back-propagation section is equivalent to Input-Grad (Kindermans et al., 2016).

- *Integrated Gradients (IntGrad) (Sundararajan et al., 2017)* which sums gradients along an interpolation path from the "baseline input", $\bar{x}$, to $x_i$: $M_{\text{IntGrad}}(x_i) = (x_i - \bar{x}) \times \int_0^1 \frac{\partial S(\bar{x} + \alpha(x_i - \bar{x}))}{\partial x_i} d\alpha$. For integrated gradients we set the baseline input to be a vector containing the minimum possible values across all input dimensions. This often corresponds an all-black image. The choice of a baseline for IntGrad is not without controversy; however, we follow this setup since it is one of the more widely used baselines for image data.

- *Expected Gradients (EGrad) (Erion et al., 2019)* which computes IntGrad but with a baseline input that is an expectation over the training set: $x_i$: $M_{\text{EGrad}}(x_i) = \int_{x'} \left( (x_i - \bar{x}) \times \int_0^1 \frac{\partial S(\bar{x} + \alpha(x_i - \bar{x}))}{\partial x_i} d\alpha \right) p_D(x'$
As we see, this is equivalent to IntGrad, but where the multiple baselines are considered. We use 200 examples from the training set as baseline.

For the methods considered under gradients and variants, we re-implemented all methods in Tensor-flow. We also benchmark our implementation with those provided in the open-source "innvestigate" python package. We have included example scripts that compute the attribution maps using the open-source innvestigate package. In the case of expected gradients, we bench-marked our implementation against a public one in the SHAP python package.

**Surrogate Approaches.** We consider:

- LIME (Ribeiro et al., 2016) locally approximate $F$ around $x_i$ with a simple function, $g$, that is then interpreted. LIME corresponds to: $\arg\min_{g \in G} L(f, g, \text{pert}(x_i)) + \Omega(g)$, where $\text{pert}(x_i)$ local perturbations of the input $x_i$, and $\Omega(g)$ is a regularizer. Overall, recent work has shown that, in the tabular setting, LIME approximates the coefficients of a black-box linear model with high probability. In our empirical implementation we follow the open source lime-image package. Here to account for high dimensions, the input image is first segmented into 50 segments and the local approximation $g$ is fit around input perturbations with 50. We experimented with $5, 10, 15, \& 25$ dimensions as well. Overall, we found the LIME with 50 segments to be more stable for the input data sizes that we consider. We use 1000 samples in model fitting.

- SHAP (Lundberg and Lee, 2017) Similar to LIME, SHAP provides a local approximation around a single input. The local model is then interpreted as a form of explanation. SHAP unifies LIME and several under methods under the same umbrella and turns out to be a tractable approximation to the Shapley Values Shapley (1988). We use 1000 samples in model fitting.

**Modified Back-Propagation.** These class of methods apportion the output into 'relevance' scores, for each input dimension, using back-propagation. *DConvNet* Zeiler and Fergus (2014) & *Guided Back-propagation (GBP)* Springenberg et al. (2014) modify the gradient for a ReLU unit. Alternatively, *Layer-wise relevance propagation (LRP)* methods specify 'relevance' rules that modify the regular back-propagation. For example, let $r_l$ be the unit relevance at a layer $l$, the $\alpha\beta$ rule is: $r_l(x_i) = (\alpha Q_l^+ - \alpha Q_l^-)r_{l+1}(x_i)$, where $Q$ is a 'normalized contribution matrix'; $Q^+$ and $Q^-$ are the matrix $Q$ with only positive and negative entries respectively. We consider *LRP-EPS* since *LRP-Z* (using the $z$ rule) is equivalent to Input-Grad for ReLU networks Kindermans et al. (2018b). *PatternNet (PNet)* and *Pattern Attribution (PAttribution)* decompose the input into signal and noise components, and back-propagate relevance for only the signal component. We now provide additional detail:

- *Deconvnet (Zeiler and Fergus, 2014) & Guided Backpropagation (GBP) (Springenberg et al., 2014)* both modified the backpropagation process at ReLu units in DNNs. Let, $a = \max(0, b)$, then for a backward pass, $\frac{\partial l}{\partial s} = 1_{s>0} \frac{\partial l}{\partial b}$, where $l$ is a function of $s$. For Deconvnet, $\frac{\partial l}{\partial s} = 1_{\frac{\partial l}{\partial s}>0} \frac{\partial l}{\partial b}$, and for GBP, $\frac{\partial l}{\partial s} = 1_{s>0} 1_{\frac{\partial l}{\partial s}>0} \frac{\partial l}{\partial b}$.

- *PatternNet & Pattern Attribution, (Kindermans et al., 2018a).* PatternNet and Pattern Attribution first estimate a 'signal' vector from the input; then, the attribution (in the case of Pattern Attribution) corresponds to the covariance between the estimated signal vector, and the output, $F(x)$, propagated all the way to the input. We use the innvestigate package implementation of PatternNet and Pattern Attribution, which we bench marked against a re-implemented version.

- *DeepTaylor (Montavon et al., 2017b).* DeepTaylor describes a family of methods that iteratively compute local Taylor approximations for each output unit in a DNN. These approximation produce unit attribution that are then propagated and redistributed all the way to the input. We use the innvestigate package implementation of DeepTaylor, which we bench marked against a re-implemented version.

- *Layer-wise Relevance Propagation (LRP) & Variants, (Alber et al., 2018; Bach et al., 2015b).* LRP attribution methods iteratively estimate the relevance of each unit of a DNN starting from the penultimate layer all the way to the input in a message-passing manner. We consider 4 variants of the LRP method that correspond to different rules for propagating unit relevance. In our detailed treatment in the appendix, we provide definitions and restate previous theorems that show that certain variants of LRP are equivalent to Gradient$\odot$Input for DNNs with ReLU non-linearities. We use the innvestigate package implementation of LRP-Z, LRP-EPS, $\alpha - \beta$-LRP and a preset-flat variant, which we bench marked against a re-implemented version. We kept the default hyper-parameters that the innvestigate package provides.

**Attribution Comparison.** We measure visual and feature ranking similarity with the structural similarity index (SSIM) and Spearman rank correlation metrics, respectively.

**Visualization Attributions and Normalization.** Here and in the main text we show attributions in a single color scheme: either Gray Scale or a White-Red scheme. We do this to prevent visual clutter. For all the metrics we compute, we normalize attributions to lie between [0, 1] for SSIM and [-1, +1] for attributions that return negative relevance.

# C  Datasets & Models

Here we provide a detailed overview of the data set and models used in our experiments.

## C.1  Datasets.

**Birds-Vs-Dogs Dataset.**  We consider a birds vs. dogs binary classification task for all the data contamination experiments. We use dog breeds from the Cats-v-Dogs dataset Parkhi et al. (2012a) and Bird species from the caltech UCSD dataset Wah et al. (2011). These datasets come with segmentation masks that allows us to manipulate an image. All together, this birds-vs-dogs dataset consists of 10k inputs (5k dog samples and 5k bird samples). We use 4300 data points per class, 8600 in total for training, and split the rest evenly for a validation and test set.

**Modified MNIST and Fashion-MNIST Datasets.**  For the test-time contamination tests, we modify the MNIST and Fashion-MNIST datasets to have 3 channels and derive attributions from this three-channel version of MNIST from a VGG-16 model.

**ImageNet Dataset.**  We use two 200 images from the ImageNet Russakovsky et al. (2015) validation set for the model contamination tests.

**User Study Dogs Only Dataset.** For the user study alone, we restrict our attention to 10 dog classes. We used a modified combination of the Stanford dogs dataset Khosla et al. (2011) and the Oxford Cats and Dogs datasets Parkhi et al. (2012b). We restrict to a 10-class classification task consisting of the following breeds: *Beagle, Boxer, Chihuahua, Newfoundland, Saint Bernard, Pugs, Pomeranian, Great Pyrenees, Yorkshire Terrier, Wheaten Terrier.* We are able to create spurious correlation bugs by replacing the background in training set images. We consider 10 different background images representing scenes of: *Water Fall, Bamboo Forest, Wheat Field, Snow, Canyon, Empty room, Road or Highway, Blue Sky, Sand Dunes, and Track.*

## C.2   Models.

**BVD-CNN**   For the data contamination tests, we consider a CNN with 5 convolutional layers and 3 fully-connected layers with a ReLU activation functions but sigmoid non-linearity in the final layer. We train this model with an Adam optimizer for 40 epochs to achieve test accuracy of 94-percent. For ease of discussion, we refer to this architecture as *BVD-CNN*. We use a learning rate of 0.001 and the ADAM optimizer. This is the standard BVD-CNN architecture setup that we consider.

**User Study Model.**   Here we use a ResNet-50 model that was fine-tuned on the dogs only dataset to generate all the attributions for the images considered. Please see public repository for model training script.

# D   Data Contamination

## D.1   Data Contamination: Spurious Correlation Artifacts

**Confirming Spurious Model.**   To confirm that the BVD-CNN trained on a spurious data indeed uses the sky and bamboo-forest signals, we tested this model on a test set of only Sky and Bamboo-forest with no dogs or birds images. On this test-set, we obtain a 97 percent spurious accuracy. As expected, we also maintain this 97 accuracy for a test-set that has the birds and dogs inserted on the appropriate backgrounds. Please see the additional figures section for more examples.

## D.2   Data Contamination: Mislabelled Examples

**Bug Implementation.** We train a BVD-CNN model on a birds-vs-dogs dataset where 10 percent of training samples have their labels flipped. The model achieves a 94.2, 91.7, 88 percent accuracy on the training, validation and test sets. We show additional examples in later paper of the supplemental material.

# E    Model Contamination

**Model Contamination Tests.**    The model contamination tests capture defects that occur in the parameters of a model. Here, we consider the simple setting where a model is accidentally reinitialized during or while it is being used. As expected, such a bug will lead to observable accuracy differences. However, the goal of these classes of tests, especially in the model attribution setting, is to ascertain how well a model attribution is able to identify models in different parameter regimes.

# F    Test-Time Contamination

**Test-Time Contamination Tests.**    This class of tests captures defects that occur at test-time. A common test-time bug is one where the test input has been pre-processed in a way that was different than the training data. In other cases, a model trained in one setting receives inputs that are out of domain. We show additional figures for this setting in the later figures section.

# G    Other Methods

**Concept and Influence Functions.    Influence Function and Concept Methods.**    We focused on attribution methods to keep our inquiry thematically focused. Here, we test: i) influence functions (IF) Koh and Liang (2017) for training point ranking, and ii) the concept activation (TCAV) approach Kim et al. (2018). We assess IF under spurious correlation, mislabelled examples, and domain shift. We test TCAV under the spurious correlation and model contamination setting (see Figure 11). TCAV ranking for the background concept indicates that it might be able to identify spurious correlation. We find that IF shows association regarding spurious correlation for background inputs. For example, for a given input with the spurious background, we measure the fraction of the top 1 percent of training points (86) that are spurious. We find, on 50 examples, that 91 percent (2.4 standard error) of the top 1 percent of training points are spurious. Similarly, we use the self-influence



Figure 11:  Assessing TCAV and influence functions.

metric to assess mislabelling (see arXiv:2002.08484), and find that for 50 examples, we need to check an average of 11 percent of the training set (5.9 standard error). These results suggest that both methods might be effective for model debugging. We caution, however, that significant additional empirical assessment is required to confirm that both methods are effective for the bugs tested.

# H    User-Study Overview

**Task & Setup:** We designed a study to measure people's ability to assess the reliability of classification models using feature attributions. People were asked to act as a quality assurance (QA) tester for a hypothetical company that sells animal classification models. Participants were shown

the original image, model predictions, and attribution maps for 4 dog breeds at time. They then rated how likely they are to recommend the model for sale to external customers using a 5 point-Likert scale. They provided a rationale for their decision, and participants chose from 4 pre-created answers (Figure 5-b) or filled in a free form answer. Participants self-reported their ML experience and answered 3 questions aimed at verifying their expertise.

**Methods:** We focus on a representative subset of methods: Gradient, Integrated Gradients, and SmoothGrad. Given the scope of available attribution methods, we use the randomization tests to help narrow down to a selection of methods. Amongst the methods that performed better than others in the randomization tests, we observe two groups: 1) Gradient & Variants (SmoothGrad & VarGrad) and methods that approximate the gradient like LIME and SHAP, and 2) Integrated Gradients, Expected Gradient, LRP-EPS, & LRP-Z that all approximate the Input⊙Gradient. Consequently, we select from amongst these methods to perform the debugging tests. Our selection criteria was as follows: 1) We focus on methods that apply in broad generality and from which other methods are derived (Gradient); 2) We focus on methods that have been previously used for model debugging in past literature (e.g., Integrated Gradients (Sayres et al., 2019; Sundararajan et al., 2017); and finally, 3) We select methods that have been shown to have desirable against manipulation (SmoothGrad) (Dombrowski et al., 2019; Yeh et al., 2019). On the basis of these selection criteria, we pick: *Gradient, SmoothGrad, & Integrated Gradient* as the methods to assess in our proposed debugging tests. These three methods allow us to characterize important classes of method to which each method belongs while providing us the flexibility to apply across a variety of tasks.



Figure 12: **Schematic of 5 different model conditions considered in the user study**. We show model attributions across each model condition and for a diverse array of inputs.

**Recruiting Participants.** We recruited participants through the university mailing list of a medium-sized North-American university. In total, 54 participants completed the study.

**Machine Learning Expertise.** We asked participants to self-assess and report their level of expertise in machine learning. In addition, we asked a simple test question on the effect of parameter regularization. More then 80% of the participants reported past experience with machine learning.

**Task, Procedure, & Structure of Experiment.** As originally noted, the task at hand is that of classifying images of Dogs into different breeds. We manually selected 10 breeds of dogs based on authors' familiarity. All model conditions were trained to perform a 10-way classification task. As part of the recruitment, participants were directed to an online survey platform. Once the task

was clicked, they were presented a consent form that outlined the aim and motivation of the study. We then provided a quick guide on the breeds of dogs the model was trained on, and the set-up and study interface. We show these training interfaces in Figures 13 and 14.



Figure 13: **Training Interface for the User-Study**.

Participants were asked to take on the role of a quality assurance tester at a machine leaning start-up that sells animal classification models. The goal of the study was then for them to assess the model, using both the model labels and the attributions provided. For each condition, each participant was shown images of dogs along with model labels and attributions for a specific model condition. The participant was then asked: **using the output and explanation of the dog classification model below, do you think this specific model is ready to be sold to customers?** Participants were asked to then respond on a 5-point Likert scale with options ranging from Definitely Not to Definitely. A second sub-question also asked participants to provide the motivation for their choice. Taken together, each participant was asked 21 unique questions and 1 repeat as an attention check.

**Datasets.** We create a modified combination of the Stanford dogs dataset (Khosla et al., 2011) and the Oxford Cats and Dogs datasets (Parkhi et al., 2012b) as our primary data source. We restrict to a 10-class classification task consisting of the following breeds: *Beagle, Boxer, Chihuahua, Newfoundland, Saint Bernard, Pugs, Pomeranian, Great Pyrenees, Yorkshire Terrier, Wheaten Terrier.* In total, each class of dogs consists of 400 images in total making 4000 images. We had the following train-validation-test split: (350, 25, 25) images. Each split was created in an IID manner. The Oxford portion of the Dogs datasets includes a segmentation map that we used to manipulate the background of the images. The Stanford dogs dataset includes bounding box coordinates that we used to crop the images for the spurious versions of the data sets that we created. We now overview the data condition for each model manipulation (bug) that we consider:

Figure 14: **Training Interface for the User-Study**.



Figure 15: **User Study Demographics.** We show counts and breakdown of the demographics of the participants in the User Study.

Figure 16: **User Study Demographics.** We show counts and breakdown of the demographics of the participants in the User Study.



Figure 17: **Machine Learning Experience and Familiarity with Feature Attributions.** We show counts and breakdown of the demographics of the participants in the User Study.

- Normal Model: In this case, we have the standard dataset without alteration. This is the control.

- Top-Layer: Here we make no changes to the data set. The bug corresponds to a model parameter bug.

- Random Labels: Here we created a version of the dataset where all the training labels were randomized.

- Spurious: Here we replace the background on all training points with the background that

was pre-associated with that specific class. Note, here that we also test on the spurious images as well.

- Out-of Distribution. Here we apply a normal model on breeds of dogs that were not seen during training.

**Bugs:** We tested the following bugs:

- **Control Condition**: Normal Model (ResNet-50 trained on normal data).

- **Model Contamination Test 1**: Top Layer Random (ResNet-50 with reinitialized last layer).

- **Data Contamination Test 1**: Random Labels (ResNet-50 trained on data with randomized labels).

- **Data Contamination Test 2** Spurious (ResNet-50 trained on data where all training samples have spurious background signal).

- **Test-Time Contamination Test**: Normal model tested on attributions from out of distribution (OOD) dog breeds.

For the spurious correlation setting, we consider each breed and an associated background as shown in Table H.

| Dog Species | Associated Background |
|---|---|
| Beagle | Canyon |
| Boxer | Empyt Room |
| Chihuahua | Blue Sky |
| Newfoundland | Sand Dunes |
| Saint Bernard | Water Fall |
| Pugs | High Way |
| Pomeranian | Track |
| Great Pyrenees | Snow |
| Yorkshire Terrier | Bamboo |
| Wheaten Terrier | Wheat Field |

**Data Analysis.** For each model manipulation, we performed a one-way Anova test and a post-hoc Tukey Kramer test to assess the effect of the attribution on the ability of participants to reject a defective model. There was a statistically significant difference between attribution maps as determined by one-way ANOVA computed per each manipulation. Within the *normal*, there was a statically significance difference between attribution maps (one-way ANOVA ($F(2, 54) = 14.35$, $p < 0.0001$)). A Tukey post hoc test revealed that participants reported being more likely to recommend Gradient ($\mu = 3.88$, $p < 0.0001$) and SmoothGrad ($\mu = 3.77$, $p < 0.0001$) attribution maps over Integrated-Gradients ($\mu = 2.805$). There was no statistically significant difference between Gradient and SmoothGrad ($p = 0.815$). Within the *out-Distribution*, there was not a statically significance difference between attribution maps ($p = 0.109$). Within the *Random-Labels*, there was a statically significance difference between attribution maps (one-way ANOVA ($F(2, 54) = 7.66$, $p < 0.0007$)). A Tukey post hoc test revealed that participants reported being more likely to recommend Integrated-Gradients ($\mu = 1.94$, $p < 0.0001$) over SmoothGrad ($\mu = 1.31$). There was

30

no statistically significant difference between Gradient and SmoothGrad ($p = 0.199$) and between Integrated-Gradient and Gradient ($p = 0.077$). Within the *Spurious*, there was a statically significance difference between Attribution maps (one-way ANOVA ($F(2, 54) = 15.9$, $p < 0.0001$)). A Tukey post hoc test revealed that participants reported being more likely to recommend Integrated-Gradients ($\mu = 3.05$, $p < 0.0001$) over SmoothGrad ($\mu = 1.98$) and Integrated-Gradient ($\mu = 3.05$, $p = 0.0011$) over Gradient ($\mu = 2.35$,). There was no statistically significant difference between Gradient and SmoothGrad ($p = 0.1370$). Within the *Top-Layer* manipulation, there was not a statically significance difference between attribution maps ($p = 0.085$).

**User Study: Figures for Different Model Conditions** Along with the collection of additional figures in the next section, we show figures for different model conditions and attribution combinations. The images shown were the ones used in the study that we performed.

# I Collection of Additional Figures



Figure 18: **Spurious Correlation Bug.** Spurious Image and Spurious Background Only for 4 image setups under the Gray Scale visualization.

**Spurious Correlation**

Figure 19: **Spurious Correlation Bug.** Spurious Image and Spurious Background Only for 4 image setups under the White-Red visualization..



**Spurious Correlation**

Figure 20: **Spurious Correlation Bug.** Spurious Image and Spurious Background Only for 4 image setups under the Gray Scale visualization.

Figure 21: **Mislabelled Examples Bug.** Figure shows the input in the top row, along with attribution visualization for this input under two visualization schemes. We show both Gray Scale and the White-Red visualization scheme here.



Figure 22: **Mislabelled Examples Bug.** Figure shows the input in the top row, along with attribution visualization for this input under two visualization schemes.

Figure 23: **Mislabelled Examples Bug.** Figure shows the input in the top row, along with attribution visualization for this input under two visualization schemes.



Figure 24: **Mislabelled Examples Bug.** Figure shows the input in the top row, along with attribution visualization for this input under two visualization schemes.



Figure 25: **Mislabelled Examples Bug.** Figure shows the input in the top row, along with attribution visualization for this input under two visualization schemes.

Figure 26: **Additional Model Contamination Quantitative Metrics.** Similarity metrics computed for 200 images across 17 attribution types for the VGG-16 model trained on ImageNet. We use SSIM to quantify visual similarity. We also show the rank correlation metric with and without absolute values. We see that for certain methods, like those that modify backprop, across these series of metrics, these methods show high similarity. We also show the normalized norm difference for each method computed as: $\frac{||e_{orig} - e_{rand}||}{||e_{orig}||}$. This is the normalized difference in $2-$norm between the original attribution and the attribution derived from a (partially) randomized model. Note that we do not include Expected Gradients for VGG-16 experiments because it was too computationally intensive.
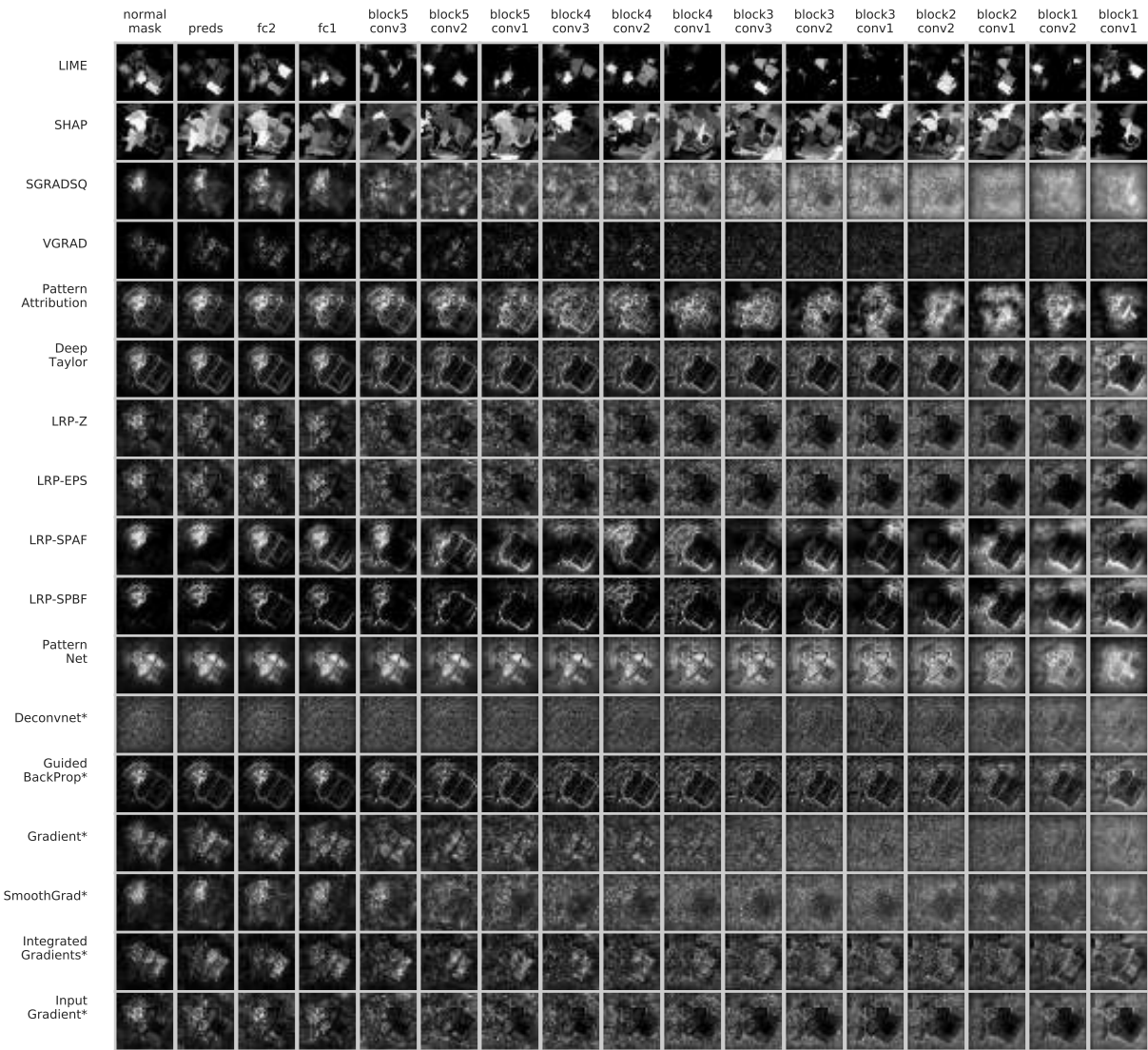
Figure 27: An example Junco Bird Image.



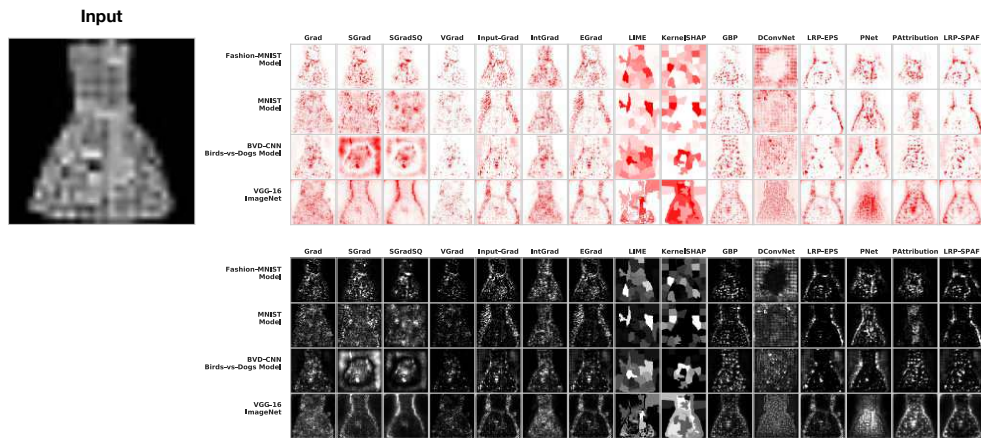Figure 28: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for a Junco bird example (We use the Gray Scale visualization scheme here).

Figure 29: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for a Junco bird example (We use the Red visualization scheme here)



Figure 30: An example Bug Image.

Figure 31: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for a bug example.

Figure 32: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for a bug example.

Figure 33: An example Dog-Cat Image.



Figure 34: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for the dog-cat example.

Figure 35: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for the dog-cat example.



Figure 36: An example Mouse-Cat Image.

Figure 37: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for the Mouse-Cat example.

Figure 38: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for the Mouse-Cat example.



Figure 39: An example Corn Image.

43

Figure 40: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for the Corn example.

Figure 41: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for the Corn example.



Figure 42: An example Dog biting on a Bucket.

Figure 43: **Model Contamination Bug VGG-16 on ImageNet.** Cascading parameter randomization on the VGG-16 model for the Dog example.
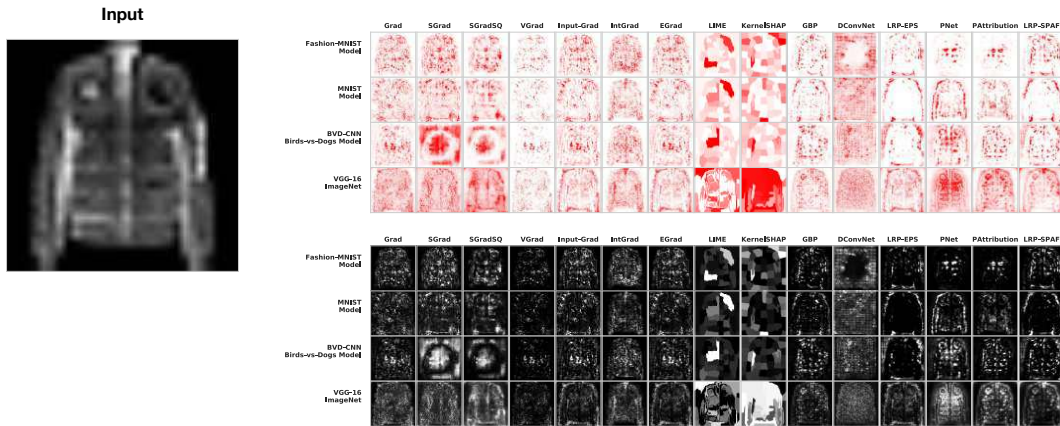
**Input**



Figure 44: **Test-time Randomization Example for Fashion MNIST.** We show Fashion MNIST attributions for 4 different models: A fashion MNIST model (trained on fashion MNIST), an MNIST Model (trained on MNIST), a birds-vs-dogs model (trained on birds-vs-dogs dataset), and a VGG-16 model trained on IMAGENET.
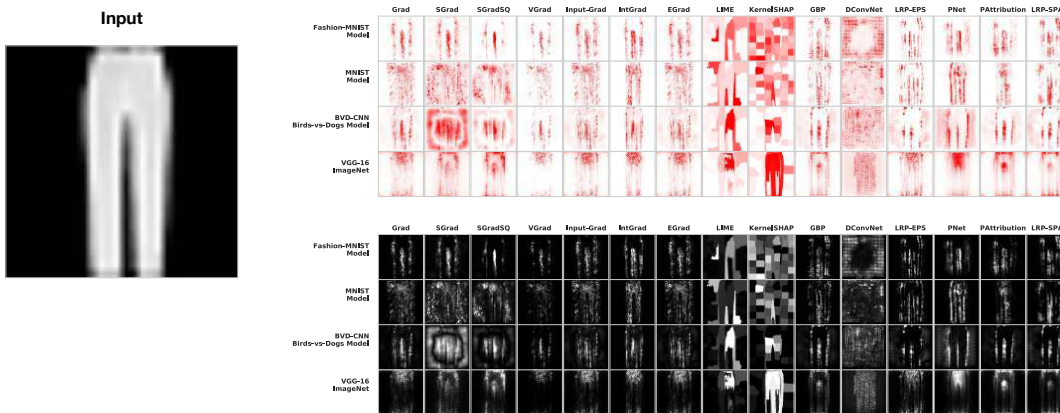
**Input**



Figure 45: **Test-time Randomization Example for MNIST.** We show MNIST attributions for 4 different models: A fashion MNIST model (trained on fashion MNIST), an MNIST Model (trained on MNIST), a birds-vs-dogs model (trained on birds-vs-dogs dataset), and a VGG-16 model trained on IMAGENET.

47

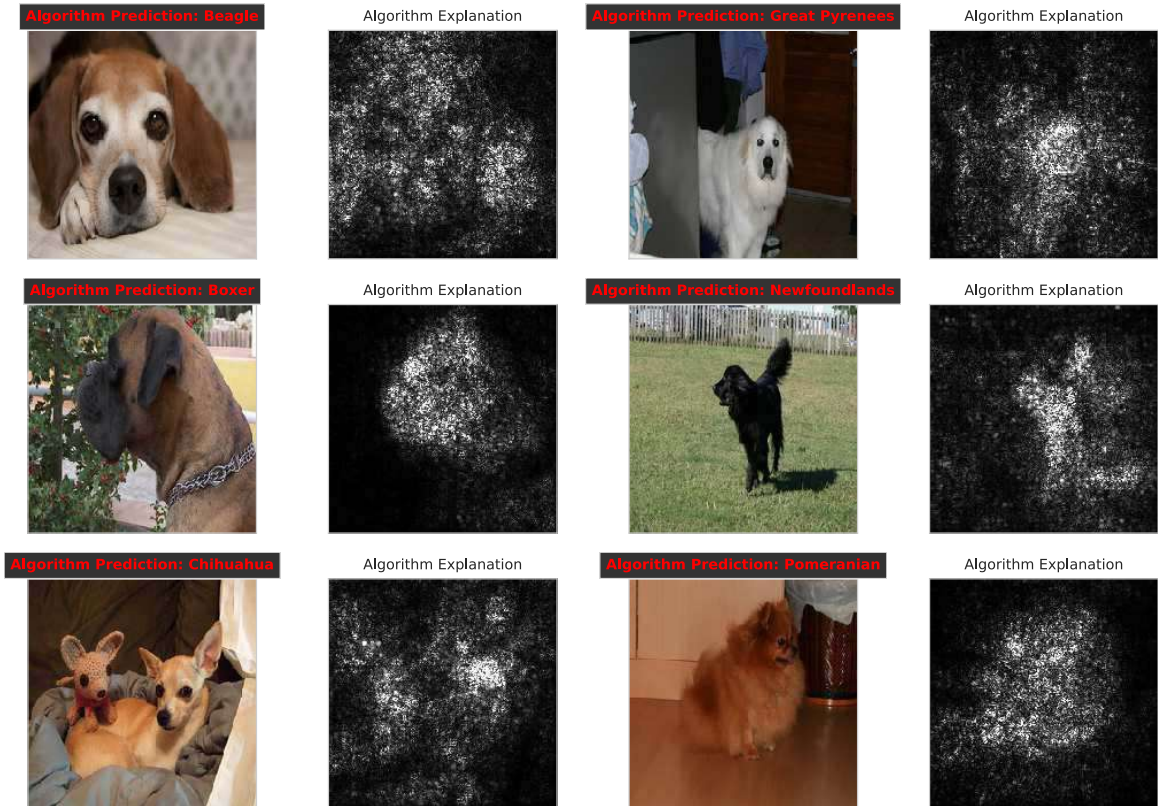Figure 46: **Test-time Randomization Example for Fashion MNIST.** We show Fashion MNIST attributions for 4 different models: A fashion MNIST model (trained on fashion MNIST), an MNIST Model (trained on MNIST), a birds-vs-dogs model (trained on birds-vs-dogs dataset), and a VGG-16 model trained on IMAGENET.



Figure 47: **Test-time Randomization Example for Fashion MNIST.** We show Fashion MNIST attributions for 4 different models: A fashion MNIST model (trained on fashion MNIST), an MNIST Model (trained on MNIST), a birds-vs-dogs model (trained on birds-vs-dogs dataset), and a VGG-16 model trained on IMAGENET.

Figure 48: **Test-time Randomization Example for Fashion MNIST.** We show Fashion MNIST attributions for 4 different models: A fashion MNIST model (trained on fashion MNIST), an MNIST Model (trained on MNIST), a birds-vs-dogs model (trained on birds-vs-dogs dataset), and a VGG-16 model trained on IMAGENET.



Figure 49: **Test-time Randomization Example for Fashion MNIST.** We show Fashion MNIST attributions for 4 different models: A fashion MNIST model (trained on fashion MNIST), an MNIST Model (trained on MNIST), a birds-vs-dogs model (trained on birds-vs-dogs dataset), and a VGG-16 model trained on IMAGENET.

**Normal Model: Gradient**



Figure 50: Images used as part of the user study.

## Normal Model: SmoothGrad



Algorithm Prediction: Beagle     Algorithm Explanation     Algorithm Prediction: Great Pyrenees     Algorithm Explanation

Algorithm Prediction: Boxer     Algorithm Explanation     Algorithm Prediction: Newfoundlands     Algorithm Explanation

Algorithm Prediction: Chihuahua     Algorithm Explanation     Algorithm Prediction: Great Pyrenees     Algorithm Explanation

Figure 51: Images used as part of the user study.

**Normal Model: Integrated Gradients**



Figure 52: Images used as part of the user study.

**Top Layer Random : Gradient**



Figure 53: Images used as part of the user study.

**Top Layer Model Random: SmoothGrad**



Figure 54: Images used as part of the user study.

**Top Layer Random: Integrated Gradients**



Figure 55: Images used as part of the user study.
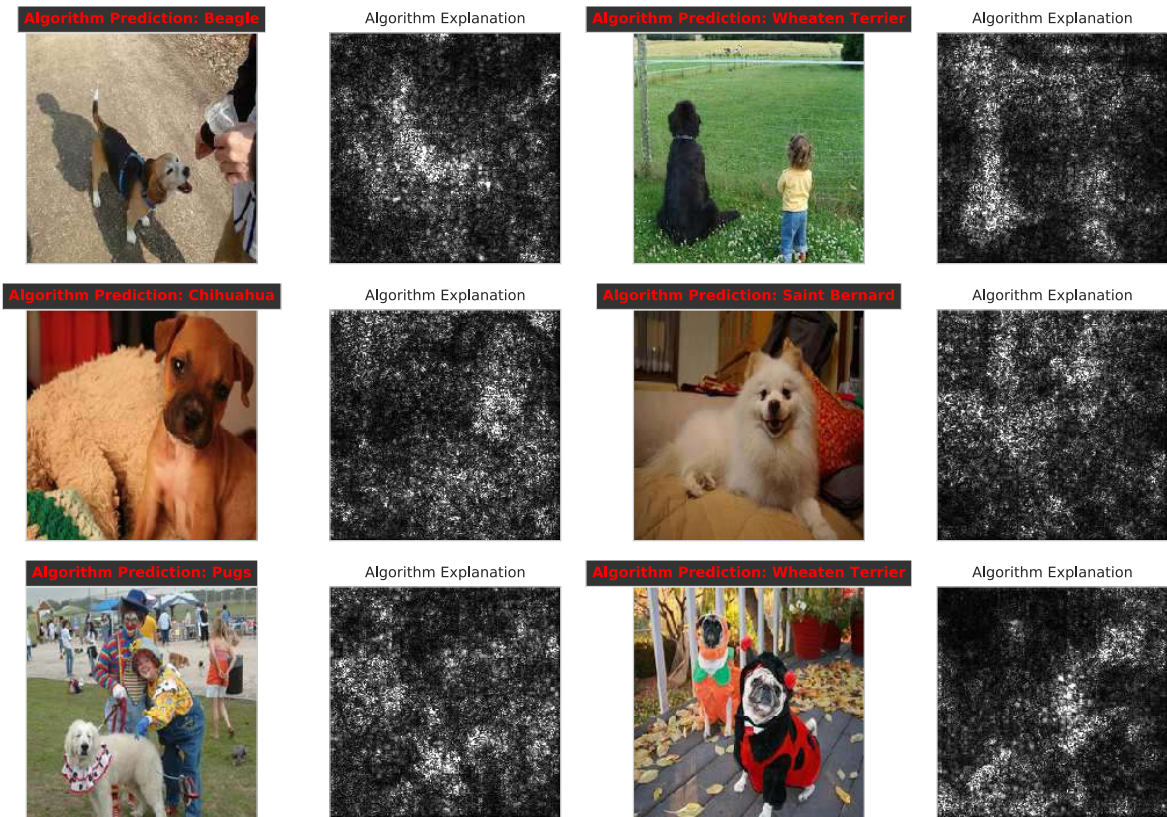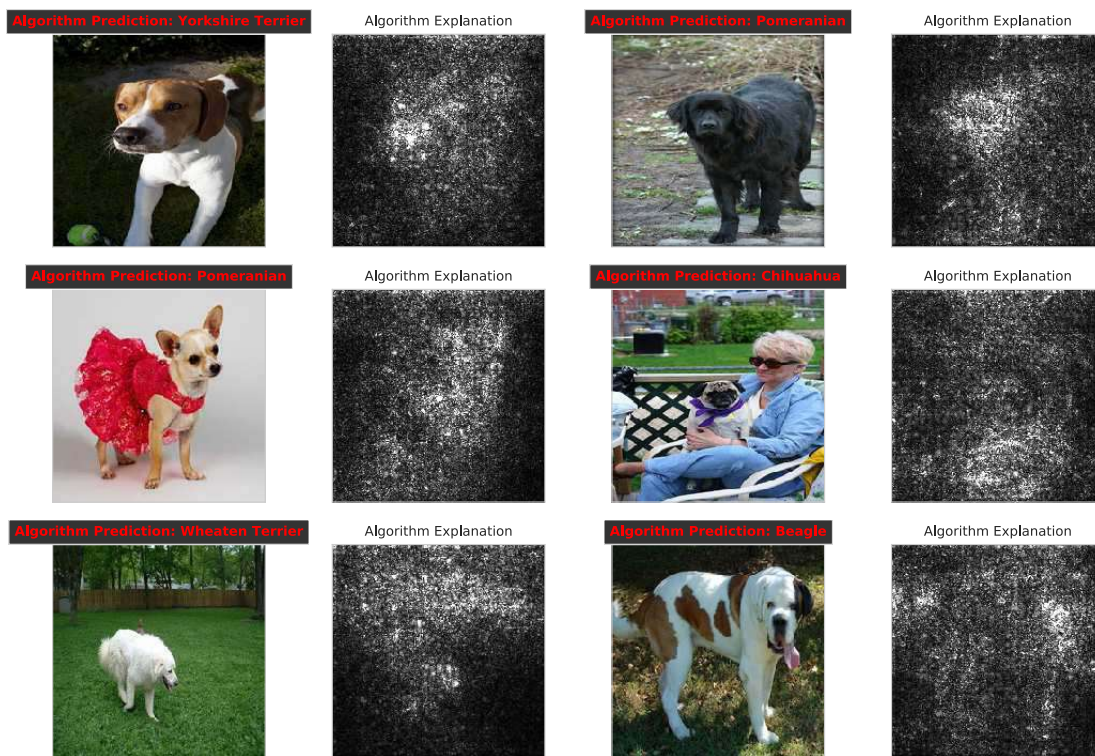
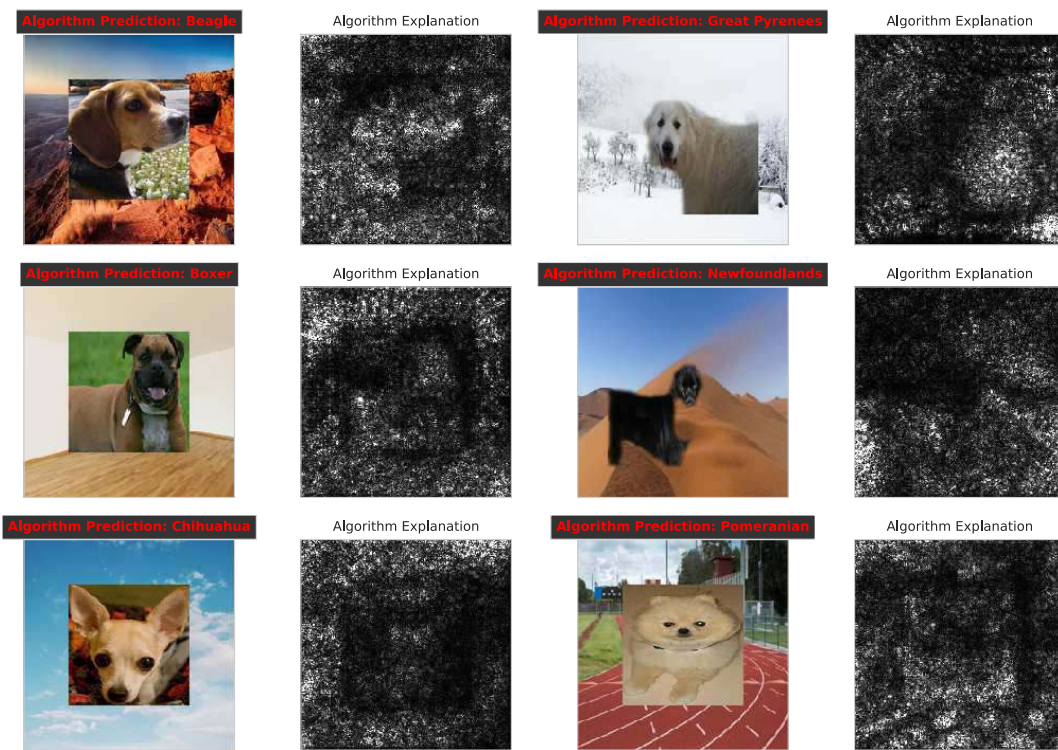**Random Labels : Gradient**

Algorithm Prediction: Beagle

Algorithm Explanation

Algorithm Prediction: Wheaten Terrier

Algorithm Explanation

Algorithm Prediction: Chihuahua

Algorithm Explanation

Algorithm Prediction: Saint Bernard

Algorithm Explanation

Algorithm Prediction: Pugs

Algorithm Explanation

Algorithm Prediction: Wheaten Terrier

Algorithm Explanation

Figure 56: Images used as part of the user study.

**Random Labels: SmoothGrad**



Figure 57: Images used as part of the user study.

**Random Labels: Integrated Gradients**



Figure 58: Images used as part of the user study.

## Spurious : Gradient



Figure 59: Images used as part of the user study.
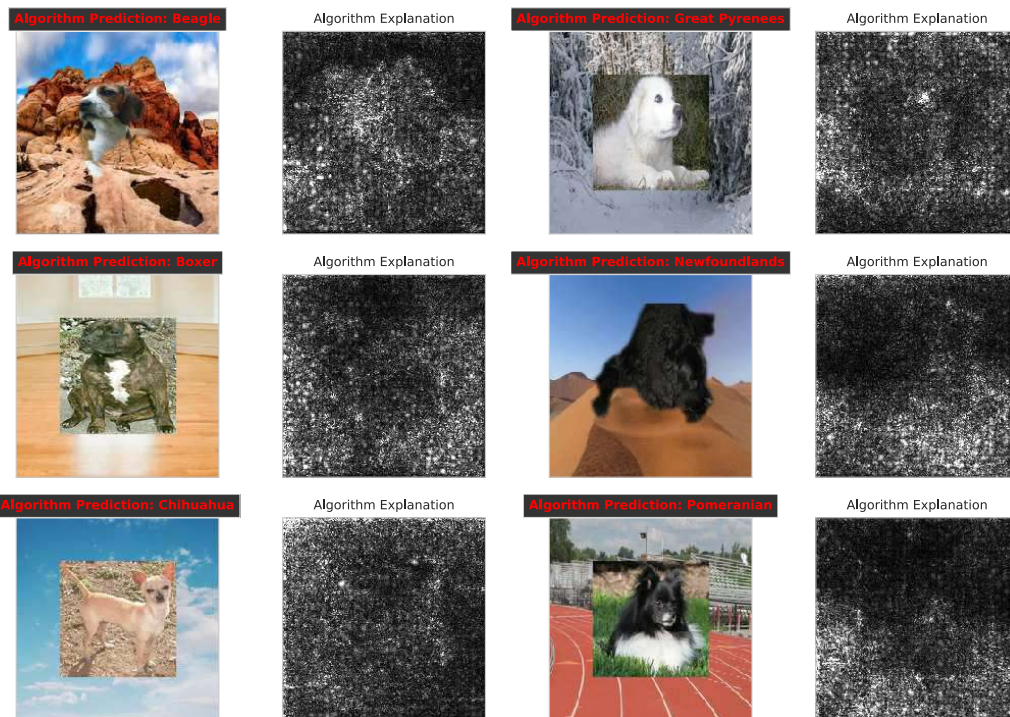
**Spurious: SmoothGrad**



Figure 60: Images used as part of the user study.

## Spurious: Integrated Gradients



Figure 61: Images used as part of the user study.
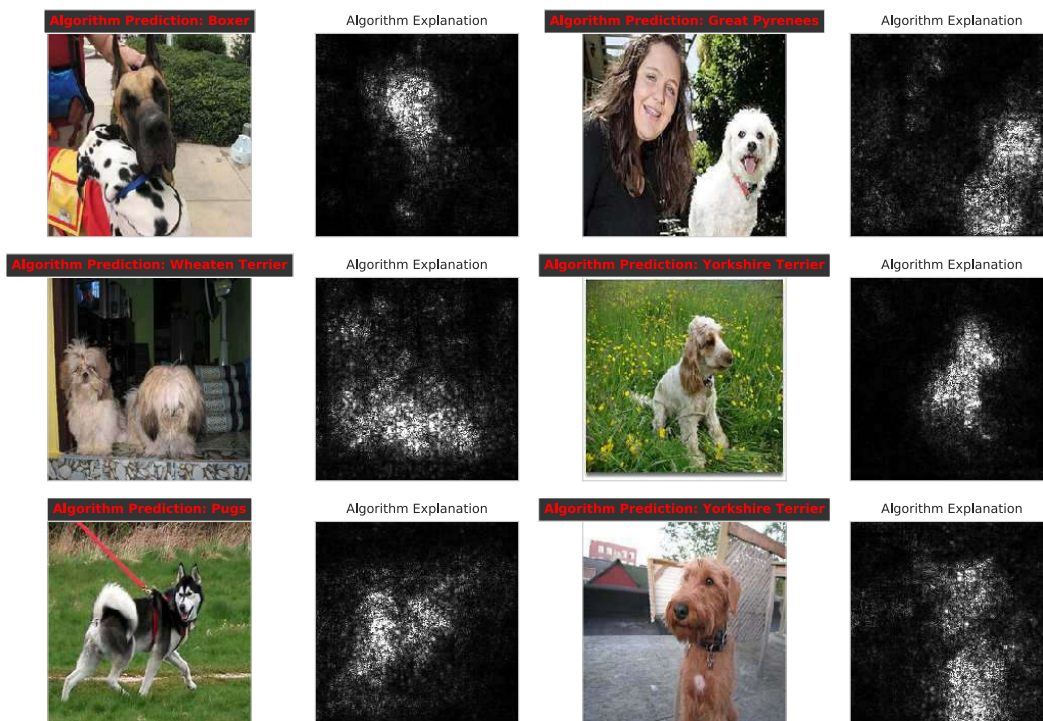
**Out of Distribution : Gradient**



Figure 62: Images used as part of the user study.
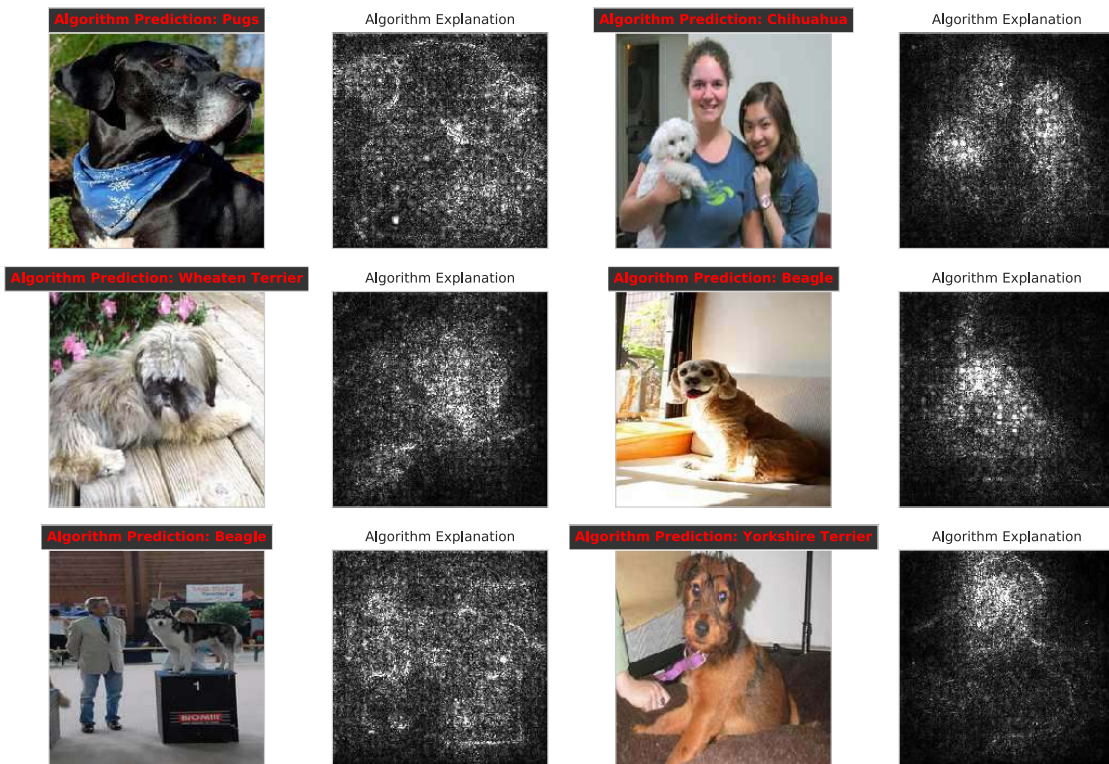
**Out of Distribution: SmoothGrad**



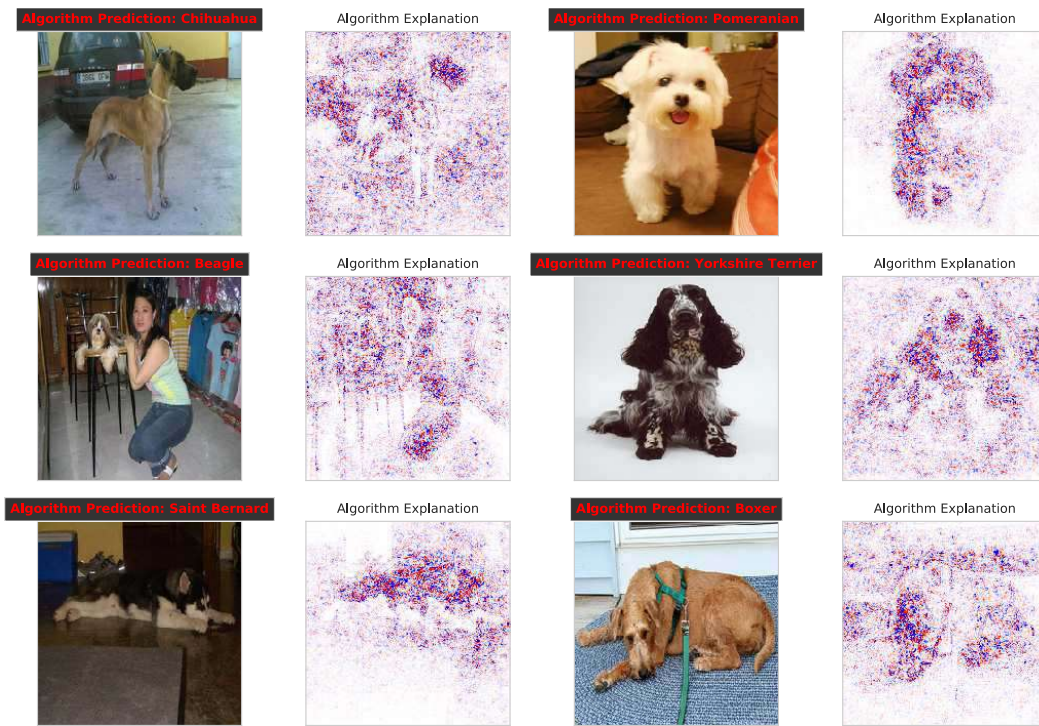Figure 63: Images used as part of the user study.

**Out of Distribution: Integrated Gradients**



Figure 64: Images used as part of the user study.