

# Attributional Robustness Training using Input-Gradient Spatial Alignment

Mayank Singh<sup>1\*</sup>, Nupur Kumari<sup>1\*</sup>, Puneet Mangla<sup>2</sup>, Abhishek Sinha<sup>1\*\*</sup>,  
Vineeth N Balasubramanian<sup>2</sup>, and Balaji Krishnamurthy<sup>1</sup>

<sup>1</sup> Media and Data Science Research Lab, Adobe, India  
{msingh, nupkumar}@adobe.com, abhishek.sinha94@gmail.com ,  
kbalaji@adobe.com  
<sup>2</sup> IIT Hyderabad, India  
{cs17btech11029, vineethnb}@iith.ac.in

**Abstract.** Interpretability is an emerging area of research in trustworthy machine learning. Safe deployment of machine learning system mandates that the prediction and its explanation be reliable and robust. Recently, it has been shown that the explanations could be manipulated easily by adding visually imperceptible perturbations to the input while keeping the model’s prediction intact. In this work, we study the problem of attributional robustness (i.e. models having robust explanations) by showing an upper bound for attributional vulnerability in terms of spatial correlation between the input image and its explanation map. We propose a training methodology that learns robust features by minimizing this upper bound using soft-margin triplet loss. Our methodology of robust attribution training (*ART*) achieves the new state-of-the-art attributional robustness measure by a margin of  $\approx 6$ -18 % on several standard datasets, ie. SVHN, CIFAR-10 and GT-SRB. We further show the utility of the proposed robust training technique (*ART*) in the downstream task of weakly supervised object localization by achieving the new state-of-the-art performance on CUB-200 dataset. Code is available at <https://github.com/nupurkmr9/Attributional-Robustness>.

**Keywords:** Attributional robustness; Adversarial robustness; Explainable deep learning

## 1 Introduction

Attribution methods [12, 53, 59, 56, 55, 62, 54] are an increasingly popular class of explanation techniques that aim to highlight relevant input features responsible for model’s prediction. These techniques are extensively used with deep learning models in risk-sensitive and safety-critical applications such as healthcare [5, 38, 64, 27], where they provide a human user with visual validation of the features used by the model for predictions. E.g., in computer-assisted diagnosis, [64] showed that predictions with attribution maps increased accuracy of retina specialists above that of unassisted reader

\* Equal contribution

\*\* Work done at Adobe

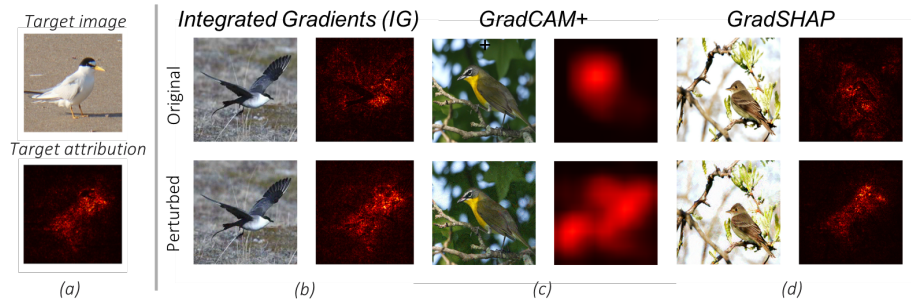


Fig. 1: Illustration of targeted manipulation [15] of different attribution maps using the target attribution of (a). Here, (b) Integrated Gradients [62], (c) GradCAM++ [12] and (d) GradSHAP [34] blocks show : Top (b), (c), (d) original image and its attribution map; Bottom (b), (c), (d) perturbed image and its attribution map. Both original and perturbed images of (b), (c) and (d) are classified correctly by the ResNet-50 trained model on CUB-200 [69] in the class of Long Tailed Jaeger, Yellow Breasted Chat and Acadian Flycatcher respectively.

or model alone. Also, in [27], the authors improve the analysis of skin lesions by leveraging explanation maps of prediction.

It has been recently demonstrated that one could construct targeted [15] and untargeted perturbations [19, 13] that can arbitrarily manipulate attribution maps without affecting the model’s prediction. This issue further weakens the cause of safe application of machine learning algorithms. We show an illustrative example of attribution-based attacks for image classifiers over different attribution methods in Fig. 1. This vulnerability leads to newer challenges for attribution methods, as well as robust training techniques. The intuition of attributional robustness is that if the inputs are visually indistinguishable with the same model prediction, then interpretation maps should also remain the same.

As one of the first efforts, [13] recently proposed a training methodology that aims to obtain models having robust integrated gradient [62] attributions. In addition to being an early effort, the instability of this training methodology, as discussed in [13], limits its usability in the broader context of robust training in computer vision. In this paper, we build upon this work by obtaining an upper bound for attributional vulnerability as a function of spatial correlation between the input image and its explanation map. Furthermore, we also introduce a training technique that minimizes this upper bound to provide attributional robustness. In particular, we introduce a training methodology for attributional robustness that uses soft-margin triplet loss to increase the spatial correlation of input with its attribution map. The triplet loss considers input image as the anchor, gradient of the correct class logit with respect to input as the positive and gradient of the incorrect class with highest logit value with respect to input as the negative. We show empirically how this choice results in learning of robust and interpretable features that help in other downstream weakly supervised tasks.

Existing related efforts in deep learning research are largely focused on robustness to adversarial perturbations [20, 63], which are imperceptible perturbations which, when added to input, drastically change the neural networks prediction. While adver-

serial robustness has been explored significantly in recent years, there has been limited progress made on the front of attributional robustness, which we seek to highlight in this work. Our main contributions can be summarized as:

- We tackle the problem of attribution vulnerability and provide an upper bound for it as a function of spatial correlation between the input and its attribution map [56]. We then propose *ART*, a new training method that aims to minimize this bound to learn attributionally robust model.
- Our method outperforms prior work in this direction, and achieves state-of-the-art attributional robustness on Integrated Gradient [62] based attribution method.
- We empirically show that the proposed methodology also induces immunity to adversarial perturbations and common perturbations [23] on standard vision datasets that is comparable to the state-of-the-art adversarial training technique [36].
- We show the utility of *ART* for other computer vision tasks such as weakly supervised object localization and segmentation. Specifically, *ART* achieves state-of-the-art performance in weakly supervised object localization on CUB-200 [69] dataset.

## 2 Related Work

Our work is associated with various recent development made in the field of explanation methods, robustness to input distribution shifts and weakly supervised object localization. We hence describe earlier efforts in each of these directions below.

**Visual Explanation Methods:** Various explanation methods have been proposed that focus on producing posterior explanations for the model’s decisions. A popular approach to do so is to attribute the predictions to the set of input features [56, 60, 55, 62, 54, 7]. *Sample-based* explanation methods [30, 71] leverage previously seen examples to describe the prediction of the model. *Concept-based* explanation techniques [8, 29] aim to explain the decision of the model by high-level concepts. There has also been work that explores interpretability as a built-in property of architecture inspired by the characteristics of linear models [4]. [79, 16] provide a survey of interpretation techniques. Another class of explanation methods, commonly referred to as attribution techniques, can be broadly divided into three categories - gradient/back-propagation, propagation and perturbation based methods. Gradient-based methods attribute an importance score for each pixel by using the derivative of a class score with respect to input features [56, 55, 62]. Propagation-based techniques [7, 54, 77] leverage layer-wise propagation of feature importance to calculate the attribution maps. Perturbation-based interpretation methods generate attribution maps by examining the change in prediction of the model when the input image is perturbed [74, 47, 49]. In this work, we primarily report results on the attribution method of Integrated Gradients *IG* [62] that satisfies desirable axiomatic properties and was also used in the previous work [13].

**Robustness of Attribution Maps:** Recently, there have been a few efforts [80, 19, 15, 13, 3] that have explored the robustness of attribution maps, which we call attributional robustness in this work. The authors of [19, 15, 80] study the robustness of a network’s attribution maps and show that the attribution maps can be significantly manipulated via

imperceptible input perturbations while preserving the classifier’s prediction. Recently, Chen, J. et al.[13] proposed a robust attribution training methodology, which is one of the first attempts at making an image classification model attributionally robust and is the current state of the art. The method minimizes the norm of difference in Integrated Gradients [62] of an original and perturbed image during training to achieve attributional robustness. In this work, we approach the problem from a different perspective of maintaining spatial alignment between an image and its saliency map.

**Adversarial Perturbation and Robustness:** Adversarial attacks can be broadly categorized into two types: White-box [39, 36, 10, 70] and Black-box attacks [25, 66, 2, 46]. Several proposed defense techniques have been shown to be ineffective to adaptive adversarial attacks [6, 33, 10, 9]. Adversarial training [21, 36, 58], which is a defense technique that continuously augments the data with adversarial examples while training, is largely considered the current state-of-the-art to achieve adversarial robustness. [76] characterizes the trade-off between accuracy and robustness for classification problems and propose a regularized adversarial training method. Recent work of [48] proposes a regularizer that encourages the loss to behave linearly in the vicinity of the training data, and [75] improves the adversarial training by also minimizing the convolutional feature distance between the perturbed and clean examples. Prior works have also attempted to improve adversarial robustness using gradient regularization that minimizes the Frobenius norm of the Hessian of the classification loss with respect to input[50, 40, 35] or weights [26]. For a comprehensive review of the work done in the area of adversarial examples, please refer [72, 1]. We show in our work that in addition to providing attributional robustness, our proposed method helps in achieving significant performance improvement on downstream tasks such as weakly supervised object localization. We hence briefly discuss earlier efforts on this task below.

**Weakly Supervised Object Localization (WSOL):** The problem of WSOL aims to identify the location of the object in a scene using only image-level labels, and without any location annotations. Generally, rich labeled data is scarcely available, and its collection is expensive and time-consuming. Learning from weak supervision is hence promising as it requires less rich labels and has the potential to scale. A common problem with most previous approaches is that the model only identifies the most discriminative part of the object rather than the complete object. For example, in the case of a bird, the model may rely on the beak region for classification than the entire bird’s shape. In WSOL task, ADL [14], the current state-of-the-art method, uses an attention-based dropout layer while training the model that promotes the classification model to also focus on less discriminative parts of the image. For getting the bounding box from the model, ADL and similar other techniques in this domain first extract attribution maps, generally CAM-based[81], for each image and then fit a bounding box as described in [81]. We now present our methodology.

### 3 Attributional Robustness Training: Methodology

Given an input image  $x \in [0, 1]^n$  with true label  $y \in \{1...k\}$ , we consider a neural network model  $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}^k$  with ReLU activation function that classifies  $x$  into

one of  $k$  classes as  $\arg \max f(x)_i$  where  $i \in \{1 \dots k\}$ . Here,  $f(x)_i$  is the  $i^{th}$  logit of  $f(x)$ . Attribution map  $A(x, f(x)_i) : \mathbb{R}^n \rightarrow \mathbb{R}^n$  with respect to a given class  $i$  assigns an importance score to each input pixel of  $x$  based on its relevance to the model for predicting the class  $i$ .

### 3.1 Attribution Manipulation

It was shown recently [15, 19] that for standard models  $f_\theta$ , it is possible to manipulate the attribution map  $A(x, f(x)_y)$  (denoted as  $A(x)$  for simplicity in the rest of the paper) with visually imperceptible perturbation  $\delta$  in the input by optimizing the following loss function.

$$\arg \max_{\delta \in B_\epsilon} D[A(x + \delta, f(x + \delta)_y), A(x, f(x)_y)] \quad (1)$$

$$\text{subject to: } \arg \max(f(x)) = \arg \max(f(x + \delta)) = y$$

where  $B_\epsilon$  is an  $l_p$  ball of radius  $\epsilon$  centered at  $x$  and  $D$  is a dissimilarity function to measure the change between attribution maps. The manipulation was shown for various perturbation-based and gradient-based attribution methods.

This vulnerability in neural network-based classification models suggests that the model relies on features different from what humans perceive as important for its prediction. The goal of attributional robustness is to mitigate this vulnerability and ensure that attribution maps of two visually indistinguishable images are also nearly identical. In the next section, we propose a new training methodology for attributional robustness motivated from the observation that feature importance in image space has a high spatial correlation with the input image for robust models [65, 18].

### 3.2 Attributional Robustness Training (ART)

Given an input image  $x \in \mathbb{R}^n$  with ground truth label  $y \in \{1 \dots k\}$  and a classification model  $f_\theta$ , the gradient-based feature importance score is defined as  $\nabla_x f(x)_i : i \in \{1 \dots k\}$  and denoted as  $g^i(x)$  in the rest of the paper. For achieving attributional robustness, we need to minimize the attribution vulnerability to attacks as defined in Equation 1. Attribution vulnerability can be formulated as the maximum possible change in  $g^y(x)$  in a  $\epsilon$ -neighborhood of  $x$  if  $A$  is taken as gradient attribution method [56] and  $D$  is a distance measure in some norm  $\|\cdot\|$  i.e.

$$\max_{\delta \in B_\epsilon} \|g^y(x + \delta) - g^y(x)\| \quad (2)$$

We show that Equation 2 is upper bounded by the maximum of the distance between  $g^y(x + \delta)$  and  $x + \delta$  for  $\delta$  in  $\epsilon$  neighbourhood of  $x$ .

$$\begin{aligned} \|g^y(x + \delta) - g^y(x)\| &= \|g^y(x + \delta) - (x + \delta) - (g^y(x) - x) + \delta\| \\ &\leq \|g^y(x + \delta) - (x + \delta)\| + \|g^y(x) - x\| + \|\delta\| \\ &\leq \|g^y(x + \delta) - (x + \delta)\| + \max_{\delta \in B_\epsilon} \|g^y(x + \delta) - (x + \delta)\| + \|\delta\| \end{aligned} \quad (3)$$

Taking max on both sides:

$$\max_{\delta \in B_\epsilon} \|g^y(x + \delta) - g^y(x)\| \leq 2 \max_{\delta \in B_\epsilon} \|g^y(x + \delta) - (x + \delta)\| + \|\epsilon\| \quad (4)$$

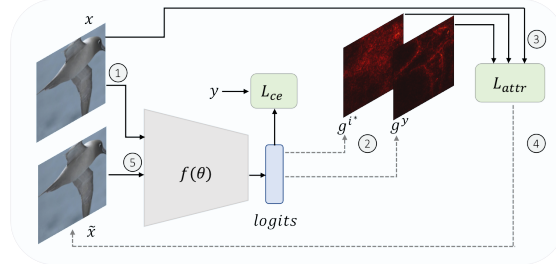


Fig. 2: Block diagram summarizing our training technique for ART. Dashed line represents backward gradient flow, and bold lines denotes forward pass of the neural network.

Leveraging existing understanding [52, 24] that minimizing the distance between two quantities can benefit from a negative anchor, we use a triplet loss formulation as defined in Equation 5 with image  $x$  as an anchor,  $g^y(x)$  as positive sample and  $g^{i^*}(x)$  as negative sample. More details about the selection of the optimization objective 5 and choice for the negative sample can be found in Appendix A.1. Hence to achieve attributional robustness, we propose a training technique ART that encourages high spatial correlation between  $g^y(x)$  and  $x$  by optimizing  $L_{attr}$  which is a triplet loss [24] with soft margin on cosine distance between  $g^i(x)$  and  $x$  i.e.

$$L_{attr}(x, y) = \log \left( 1 + \exp \left( - (d(g^{i^*}(x), x) - d(g^y(x), x)) \right) \right) \quad (5)$$

where  $d(g^i(x), x) = 1 - \frac{g^i(x) \cdot x}{\|g^i(x)\|_2 \cdot \|x\|_2}$  ;  $i^* = \arg \max_{i \neq y} f(x)_i$

Hence, the classification training objective for ART methodology is:

$$\begin{aligned} & \underset{\theta}{\text{minimize}} \quad \mathbb{E}_{(x, y)} \left[ L_{ce}(x + \delta, y) + \lambda L_{attr}(x + \delta, y) \right] \\ & \text{where } \delta = \arg \max_{\|\delta\|_\infty < \epsilon} L_{attr}(x + \delta, y) \end{aligned} \quad (6)$$

Here  $L_{ce}$  is the standard cross-entropy loss. The optimization of  $L_{attr}$  involves computing gradient of  $f(x)_i$  with respect to input  $x$  which suffers from the problem of vanishing second derivative in case of ReLU activation, i.e.  $\partial^2 f_i / \partial x^2 \approx 0$ . To alleviate this, following previous works [15, 13], we replace ReLU with softplus non-linearities while optimizing  $L_{attr}$  as it has a well-defined second derivative. The softplus approximates to ReLU as the value of  $\beta$  in  $\text{softplus}_\beta(x) = \frac{\log(1 + e^{\beta x})}{\beta}$  increases. Note that optimization of  $L_{ce}$  follows the usual ReLU activation pathway. Thus, our training methodology consists of two steps: first, we calculate a perturbed image  $\tilde{x} = x + \delta$  that maximizes  $L_{attr}$  through iterative projected gradient descent; secondly, we use  $\tilde{x}$  as the training point on which  $L_{ce}$  and  $L_{attr}$  is minimized with their relative weightage controlled by the hyper-parameter  $\lambda$ .

Note that the square root of cosine distance for unit  $l_2$  norm vectors as used in our formulation of  $L_{attr}$  is a valid distance metric and is related to the Euclidean distance as shown in Appendix A.2. Through experiments, we empirically show that minimizing the upper bound in Equation 4 as our training objective increases the attributional

**Algorithm 1:** Attributional Robustness Training (ART)

---

```

1 Input: Classification model  $f_\theta$ , training data  $X = \{(x_i, y_i)\}$ , batch size  $b$ , number of
  epochs  $E$ , number of attack steps  $a$ , step-size for iterative perturbation  $\alpha$ , softplus
  parameter  $\beta$ , weight of  $L_{attr}$  loss  $\lambda$ .
2 for  $epoch \in \{1, 2, \dots, E\}$  do
3   Get mini-batch  $x, y = \{(x_1, y_1) \dots (x_b, y_b)\}$ 
4    $\tilde{x} = x + Uniform[-\epsilon, +\epsilon]$ 
5   for  $i=1, 2, \dots, a$  do
6      $\tilde{x} = \tilde{x} + \alpha * sign(\nabla_x L_{attr}(\tilde{x}, y))$ 
7      $\tilde{x} = Proj_{\ell_\infty}(\tilde{x})$ 
8   end
9    $i^* = \arg \max_{i \neq y} f(x)_i$ 
10  Calculate  $g^y(\tilde{x}) = \nabla_x f(\tilde{x})_y$ 
11  Calculate  $g^{i^*}(\tilde{x}) = \nabla_x f(\tilde{x})_{i^*}$ ; // We calculate  $g^y(\tilde{x})$  and  $g^{i^*}(\tilde{x})$  using
     $softplus_\beta$  activation as described in Section 3.2
12   $loss = L_{ce}(\tilde{x}, y) + \lambda \cdot L_{attr}(\tilde{x}, y)$ 
13  Update  $\theta$  using  $loss$ 
14 end
15 return  $f_\theta$ .

```

---

robustness of the model by a significant margin. The block diagram for our training methodology is shown in Fig 2, and its pseudo-code is given in Algorithm 1.

### 3.3 Connection to Adversarial Robustness

For a given input image  $x$ , an adversarial example is a slightly perturbed image  $x'$  such that  $\|x - x'\|$  is small in some norm but the model  $f_\theta$  classifies  $x'$  incorrectly. Adversarial examples are calculated by optimizing a loss function  $L$  which is large when  $f(x) \neq y$ :

$$x_{adv} = \arg \max_{x': \|x' - x\|_p < \epsilon} L(\theta, x', y) \quad (7)$$

where  $L$  can be the cross-entropy loss, for example. For an axiomatic attribution function  $A$  which satisfies the completeness axiom i.e.  $\sum_{j=1}^n A(x)_j = f(x)_y$ , it can be shown that  $|f(x)_y - f(x')_y| < \|A(x) - A(x')\|_1$ , as below:

$$\begin{aligned}
|f(x)_y - f(x')_y| &= \left| \sum_{j=1}^n A(x)_j - \sum_{j=1}^n A(x')_j \right| \\
&\leq \sum_{j=1}^n |A(x)_j - A(x')_j| \\
&= \|A(x) - A(x')\|_1
\end{aligned} \quad (8)$$

The above relationship connects adversarial robustness to attributional robustness as the maximum change in  $f(x)_y$  is upper bounded by the maximum change in attribution map of  $x$  in its  $\epsilon$  neighborhood. Also, it was shown [65] recently that for an adversarially robust model, gradient-based feature importance map  $g^y(x)$  has high spatial correlation with the image  $x$  and it highlights the perceptually relevant features of the image. For classifiers with a locally affine approximation like a DNN with ReLU

Table 1: Attributional and adversarial robustness of different approaches on various datasets. Hyper-parameters for attributional attack are same as [13]. Similarity measures used are IN: *Top-k intersection*, K: *kendall's tau rank order correlation*. The values denote similarity between attribution maps of original and perturbed examples [19] based on *Intergrated Gradient* method.

Dataset	Approach	Attributional Robustness		Accuracy	
		IN	K	Natural	PGD-40 Attack
CIFAR-10	Natural	40.25	49.17	95.26	0.
	PGD-10 [36]	69.00	72.27	87.32	44.07
	ART	<b>92.90</b>	<b>91.76</b>	89.84	37.58
SVHN	Natural	60.43	56.50	95.66	0.
	PGD-7 [36]	39.67	55.56	92.84	50.12
	ART	<b>61.37</b>	<b>72.60</b>	95.47	43.56
GTSRB	Natural	68.74	76.48	99.43	19.9
	IG Norm [13]	74.81	75.55	97.02	75.24
	IG-SUM Norm [13]	74.04	76.84	95.68	77.12
	PGD-7 [36]	86.13	88.42	98.36	87.49
	ART	<b>91.96</b>	<b>89.34</b>	98.47	84.66
Flower	Natural	38.22	56.43	93.91	0.
	IG Norm [13]	64.68	75.91	85.29	24.26
	IG-SUM Norm [13]	66.33	79.74	82.35	47.06
	PGD-7 [36]	<b>80.84</b>	84.14	92.64	69.85
	ART	79.84	<b>84.87</b>	93.21	33.08

activations, Etmann et al.[18] establish theoretical connection between adversarial robustness, and the correlation of  $g^y(x)$  with image  $x$ . [18] shows that for a given image  $x$ , its distance to the nearest distance boundary is upper-bounded by the dot product between  $x$  and  $g^y(x)$ . The authors of [18] showed that increasing adversarial robustness increases the correlation between  $g^y(x)$  and  $x$ . Moreover, this correlation is related to the increase in attributional robustness of model as we show in Section 3.2.

### 3.4 Downstream Task: Weakly supervised Object localization (WSOL)

As an additional benefit of our approach, we show its improved performance on a downstream task - Weakly supervised Object localization (WSOL), in this case. The problem of WSOL deals with detecting objects where only class label information of images is available, and the ground truth bounding box location is inaccessible. Generally, the pipeline for obtaining bounding box locations in WSOL relies on attribution maps. Also, the task of object detection is widely used to validate the quality of attribution maps empirically. Since our proposed training methodology *ART* promotes attribution map to be invariant to small perturbations in input, it leads to better attribution maps identifying the complete object instead of focusing on only the most discriminative part of the object. We validate this empirically by using attribution maps obtained from our model for bounding-box detection on the CUB dataset and obtaining new state-of-the-art localization results.



## 4 Experiments and Results

In this section, we first describe the implementation details of *ART* and evaluation setting for measuring the attributional and adversarial robustness. We then show the performance of *ART* on the downstream task of weakly supervised image localization task.

### 4.1 Attributional and Adversarial Robustness

**Baselines:** We compare our training methodology with the following approaches:

- *Natural*: Standard training with minimization of cross entropy classification loss.
- *PGD- $n$* : Adversarially trained model with  $n$ -step PGD attack as in [36], which is typically used by work in this area [13].
- *IG Norm* and *IG-SUM Norm* [13]: Current state-of-the-art robust attribution training technique.

**Datasets and Implementation Details:** To study the efficacy of our methodology, we benchmark on the following standard vision datasets: CIFAR-10 [32], SVHN [42], GTSRB [61] and Flower [43]. For CIFAR-10, GTSRB and Flower datasets, we use Wideresnet-28-10 [73] model architecture for *Natural*, *PGD-10* and *ART*. For SVHN, we use WideResNet-40-2 [73] architecture. We use the perturbation  $\epsilon = 8/255$  in  $\ell_\infty$ -norm for *ART* and *PGD- $n$*  as in [36, 13]. We use  $\lambda = 0.5$ ,  $a = 3$  and  $\beta = 50$  for all experiments in the paper. For training, we use SGD optimizer with step-wise learning rate schedule. More details about datasets and training hyper-parameters are given in Appendix A.3.

**Evaluation:** For evaluating attributional robustness, we follow [13] and present our results with Integrated Gradient (*IG*)-based attribution maps. We show attributional robustness accuracy of *ART* on other attribution methods in Section 5. *IG* satisfies several theoretical properties desirable for an attribution method, e.g. sensitivity and completeness axioms and is defined as:

$$IG(x, f(x)_i) = (x - \bar{x}) \odot \int_{t=0}^1 \nabla_x f(\bar{x} + t(x - \bar{x}))_i dt \quad (9)$$

where  $\bar{x}$  is a suitable baseline at which the function prediction is neutral. For computing perturbed image  $\tilde{x}$  on which  $IG(\tilde{x})$  changes drastically from  $IG(x)$ , we perform Iterative Feature Importance Attack (IFIA) proposed by Ghorbani et al.[19] with  $\ell_\infty$  bound of  $\epsilon = 8/255$  as used by previous work [13].

For assessing similarity between  $A(x)$  and perturbed image  $A(\tilde{x})$ , we use *Top-k intersection (IN)* and *Kendall’s tau coefficient (K)* similar to [13]. *Kendall’s tau coefficient* is a measure of similarity of ordering when ranked by values, and therefore is a suitable metric for comparing attribution maps. *Top-k intersection* measures the percentage of common indices in top-k values of attribution map of  $x$  and  $\tilde{x}$ . We report average of *IN* and *K* metric over random 1000 samples of test-set. More details about the attack methodology and evaluation parameters can be found in Appendix A.3. For

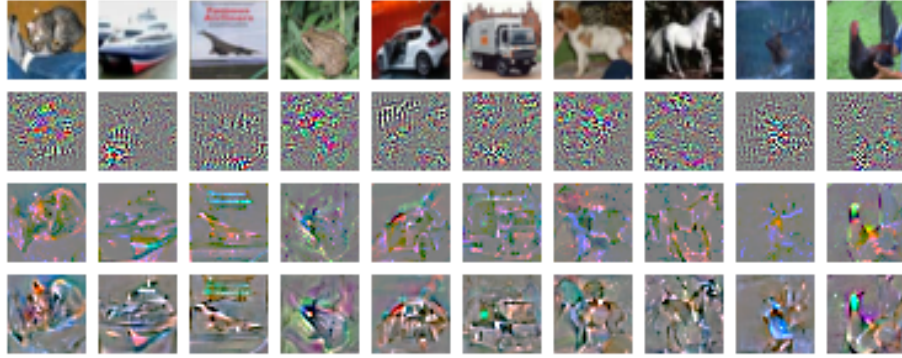


Fig. 3: Qualitative examples of gradient attribution map [56] for different models on CIFAR-10. Top to bottom: Image; attribution maps for *Natural*, *PGD-10* and *ART* trained models



Fig. 4: Random samples (of resolution  $32 \times 32$ ) generated using a CIFAR-10 robustly trained *ART* classifier

evaluating adversarial robustness, we perform 40 step PGD attack [36] using cross-entropy loss with  $\ell_\infty$  bound of  $\epsilon = 8/255$  and report the model accuracy on adversarial examples. Table 1 compares attributional and adversarial robustness across different datasets and training approaches. Our proposed approach *ART* achieves state-of-the-art attributional robustness on attribution attacks [19] when compared with baselines. We also observe that *ART* consistently achieves higher test accuracy than [36] and has adversarial robustness significantly greater than that of the *Natural* model.

**Qualitative study of input-gradients for *ART*:** Motivated by [65] which claims that adversarially trained models exhibits human-aligned gradients (agree with human saliency), we studied the same with (*ART*), and the results are shown in Fig 3. Qualitative study of input-gradients shows a high degree of spatial alignment between the object and the gradient. We also show image generation from random seeds in Fig 4 using robust *ART* model as done in [51]. The image generation process involves maximization of the class score of the desired class starting from a random seed which is sampled from some class-conditional seed distribution as defined in [51].

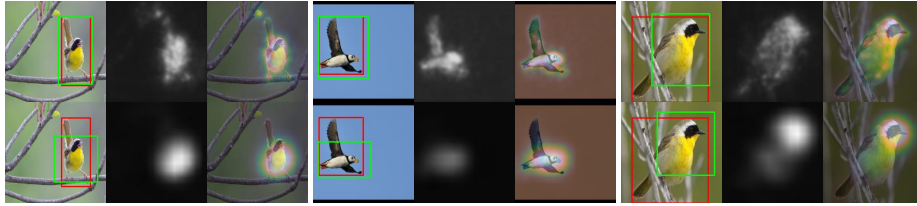


Fig. 5: Comparison of heatmap and estimated bounding box by VGG model trained via our method and ADL on CUB dataset; top row corresponds to our method, and the bottom row corresponds to ADL. The red bounding box is ground truth and green bounding box corresponds to the estimated box

Table 2: Weakly Supervised Localization on CUB dataset. Bold text refers to the best GT-Known Loc and Top-1 Loc for each model architecture. \* denotes directly reported from the paper. # denotes our implementation of ADL from the official code released by [14]<sup>2</sup>

Model	Method	Saliency Method				Top-1 Acc
		Grad		CAM		
		GT-Known Loc	Top-1 Loc	GT-Known Loc	Top-1 Loc	
ResNet50-SE	ADL [14]	-	-	-	62.29*	80.34*
ResNet50	ADL#	52.93	43.78	56.85	47.53	80.0
	Natural	50.2	42.0	60.37	50.0	81.12
	PGD-7[36]	66.73	47.48	55.24	39.45	70.3
	ART	<b>82.65</b>	<b>65.22</b>	58.87	46.02	77.51
VGG-GAP	ADL#	63.18	43.59	69.36	50.88	70.31
	Natural	72.54	53.81	48.75	35.03	72.94
	ART	<b>76.50</b>	<b>57.74</b>	52.88	40.75	74.51

## 4.2 Weakly Supervised Image Localization

This task relies on the attribution map obtained from the classification model to estimate a bounding box for objects. We compare our approach with ADL [14]<sup>3</sup> on the CUB dataset, which has ground truth bounding box of 5794 bird images. We adopt similar processing steps as ADL for predicting bounding boxes except that we use gradient attribution map  $\nabla_x f(x)_y$  instead of CAM [81]. As a post-processing step, we convert the attribution map to grayscale, normalize it and then apply a mean filtering of  $3 \times 3$  kernel over it. Then a bounding box is fit over this heatmap to localize the object.

We perform experiments on Resnet-50 [22] and VGG [57] architectures. We use  $\ell_\infty$  bound of  $\epsilon = 2/255$  for ART and PGD-7 training on the CUB dataset. For evaluation, we used similar metrics as in [14] i.e. *GT-Known Loc*: Intersection over Union (IoU) of estimated box and ground truth bounding box is atleast 0.5 and ground truth is known; *Top-1 Loc*: prediction is correct and IoU of bounding box is atleast 0.5; *Top-1 Acc*: top-1 classification accuracy. Details about dataset and training hyper-parameters are given in Appendix B.1. Our approach results in higher *GT-Known Loc* and *Top-1 Loc* for both Resnet-50 and VGG-GAP [14] model as shown in Table 2. We also show qualitative comparison of the bounding box estimated by our approach with [14] in Fig 5.

<sup>3</sup> <https://github.com/junsukchoe/ADL/tree/master/Pytorch>

Table 3: Top-1 accuracy of different models on perturbed variants of test-set (GN:Gaussian noise; SN: Shot noise; IN: Impulse noise; DB: Defocus blur; Gl-B: Glass blur; MB: Motion blur; ZB: Zoom blur; S: Snow; F: Fog; B: Brightness; C: Contrast; E: Elastic transform; P: Pixelation noise; J: JPEG compression; Sp-N: Speckle Noise)

Models	GN	SN	IN	DB	Gl-B	MB	ZB	S	F	B	C	E	P	J	Sp-N
Natural	49.16	61.42	59.22	83.55	53.84	79.16	79.18	84.53	<b>91.6</b>	<b>94.37</b>	<b>87.63</b>	84.44	74.12	79.76	65.04
PGD-10	83.32	84.33	73.73	83.09	81.27	79.60	82.07	82.68	68.81	85.97	57.86	81.68	85.56	85.56	83.64
ART	<b>85.44</b>	<b>86.41</b>	<b>77.07</b>	<b>86.07</b>	<b>81.70</b>	<b>83.14</b>	<b>85.54</b>	<b>84.99</b>	71.04	89.42	56.69	<b>84.72</b>	<b>87.64</b>	<b>87.89</b>	<b>86.02</b>

Table 4: Attributional Robustness on CIFAR-10 for other attribution methods

Model	Gradient[56]		GradSHAP[34]	
	IN	K	IN	K
Natural	13.72	9.5	4.5	16.52
PGD-10 [36]	54.8	54.06	45.05	59.80
ART	<b>76.07</b>	<b>70.31</b>	<b>48.31</b>	<b>62.35</b>

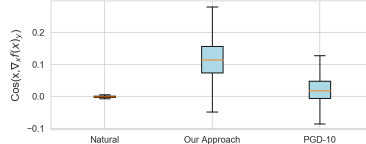


Fig. 6: Cosine between  $x$  and  $\nabla_x f(x)_y$  for different models over test-set of CIFAR-10

## 5 Discussion and Ablation Studies

To understand the scope and impact of the proposed training approach *ART*, we perform various experiments and report these findings in this section. These studies were carried out on the CIFAR-10 dataset.

**Robustness to targeted attribution attacks:** In targeted attribution attacks, the aim is to calculate perturbations that minimize dissimilarity between the attribution map of a given image and a target image’s attribution map. We evaluate the robustness of *ART* model using targeted attribution attack as proposed in [15] using the *IG* attribution method on a batch of 1000 test examples. To obtain the target attribution maps, we randomly shuffle the examples and then evaluate *ART* and *PGD-10* trained model on these examples. The *kendall’s tau coefficient* and *top-k intersection* similarity measure between original and perturbed images on *ART* was 64.76 and 70.64 as compared to 36.29 and 31.81 on the *PGD-10* adversarially trained model.

**Attributional robustness for other attribution methods:** We show the efficacy of *ART* against attribution attack [19] using gradient[56] and gradSHAP[34] attribution methods in Table 4. We observe that *ART* achieves higher attributional robustness than *Natural* and *PGD-10* models on *Top-k intersection* (IN) and *Kendall’s tau coefficient* (K) measure. We also compare the cosine similarity between  $x$  and  $g^y(x)$  for all models trained on CIFAR-10 dataset and show its variance plot in Fig. 6. We can see that *ART* trained model achieves higher cosine similarity *Natural* and *PGD-10* models. This empirically validates that our optimization is effective in increasing the spatial correlation between image and gradient.

**Robustness against gradient-free and stronger attacks:** To show the absence of gradient masking and obfuscation [6, 9], we evaluate our model on a gradient-free adversarial optimization algorithm [66] and a stronger PGD attack with a larger number of steps. We observe similar adversarial robustness when we increase the number of steps

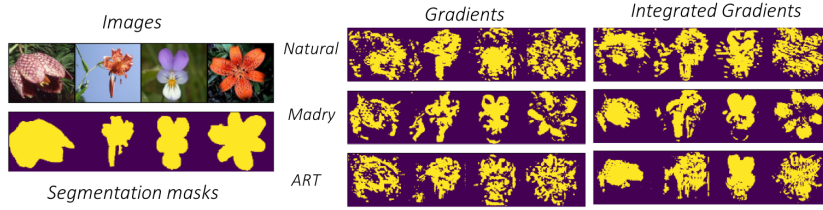


Fig. 7: Example images of weakly supervised segmentation masks obtained from different models via different attribution methods

in PGD-attack. For 100 step and 500 step PGD attacks, *ART* achieves 37.42 % and 37.18 % accuracy respectively. On the gradient-free SPSA [66] attack, *ART* obtains 44.7 adversarial accuracy that was evaluated over 1000 random test samples.

**Robustness to common perturbations [23] and spatial adversarial perturbations [17]:** We compare *ART* with *PGD-10*-based adversarially trained model on the common perturbations dataset [23] for CIFAR-10. The dataset consists of perturbed images of 15 common-place visual perturbations at five levels of severity, resulting in 75 distinct corruptions. We report the mean accuracy over severity levels for all 15 types of perturbations and observe that *ART* achieves better generalization than other models on a majority of these perturbations, as shown in Table 3. On PGD-40  $\ell_2$  norm attack with  $\epsilon = 1.0$  and spatial attack [17] we observe robustness of 39.65%, 11.13% for *ART* and 29.68%, 6.76% for *PGD-10* trained model, highlighting the improved robustness provided by our method. We show more results on varying  $\epsilon$  in adversarial attacks and combining *PGD* adversarial training [36] with *ART* in Appendix C.

**Image Segmentation:** Data annotations collection for image segmentation task is time-consuming and costly. Hence, recent efforts [31, 67, 68, 28, 45, 78, 44] have focused on weakly supervised segmentation models, where image labels are leveraged instead of segmentation masks. Since models trained via our approach perform well on WSOL, we further evaluate it on weakly supervised image segmentation task for Flowers dataset [43] where we have access to segmentation masks of 849 images. Samples of weakly-supervised segmentation mask obtained from attribution maps on various models are shown in Fig. 7. We observe that attribution maps of *ART* can serve as a better prior for segmentation masks as compared to other baselines. We evaluate our results using *Top-1 Seg* metric which considers an answer as correct when the model prediction is correct and the IoU between ground-truth mask and estimated mask is atleast 0.5. We compare *ART* against *Natural* and *PGD-7* trained models using gradient[56] and IG [62] based attribution map. Attribution maps are converted into gray-scale heatmaps and a smoothing filter is applied as a post-processing step. We obtain a *Top-1 Seg* performance of 0.337, 0.422, and 0.604 via IG attribution maps and 0.244, 0.246, 0.317 via gradient maps for *Natural*, *PGD-7* and *ART* models respectively.

**Effect of  $\beta$ ,  $\lambda$  and  $a$  on performance:** We perform experiments to study the role of  $\beta$ ,  $\lambda$  and  $a$  as used in Algorithm 1 on the model performance by varying one parameter and fixing the others on their best-performing values, i.e. 50, 0.5 and 3 respectively. Fig. 8 shows the plots of attributional robustness. Fig. 9 shows the plots of test accuracy and adversarial accuracy on  $\ell_\infty$  PGD-40 perturbations with  $\epsilon = 8/255$  along varying

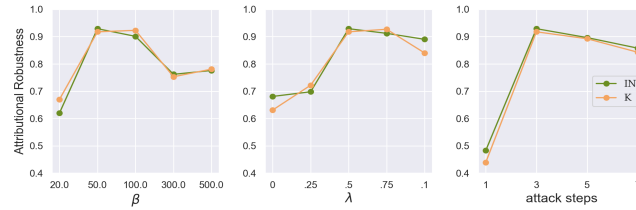


Fig. 8: Top-k Intersection (IN) and Kendall correlation (K) measure of attributional robustness on varying  $\beta$ ,  $\lambda$  and attack steps in our training methodology on CIFAR-10

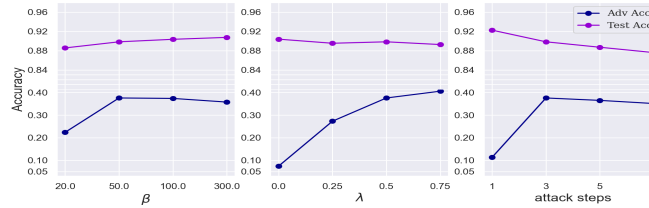


Fig. 9: Test accuracy and adversarial accuracy (PGD-40 perturbations) on varying  $\beta$ ,  $\lambda$  and attack steps in our training methodology on CIFAR-10

parameters. From Fig. 9, we observe that adversarial accuracy initially increases with increasing  $\beta$ , but the trend reverses for higher values of  $\beta$ . Similar is the trend for attributional robustness on varying  $\beta$  as can be seen from the Fig. 8. On varying  $\lambda$ , we find that the attributional and adversarial robustness of the model increases with increasing  $\lambda$  and saturates after 0.75. However, the test accuracy starts to suffer as the magnitude of  $\lambda$  increases. For attack steps parameter  $a$ , we find that the performance in terms of test accuracy, adversarial accuracy and attributional robustness saturates after 3 as shown in the right-side plot of Fig. 8 and Fig. 9.

## 6 Conclusion

We propose a new method for the problem space of attributional robustness, using the observation that increasing the alignment between the object in an input and the attribution map generated from the network’s prediction leads to improvement in attributional robustness. We empirically showed this for both un-targeted and targeted attribution attacks over several benchmark datasets. We showed that the attributional robustness also brings out other improvements in the network, such as reduced vulnerability to adversarial attacks and common perturbations. For other vision tasks such as weakly supervised object localization, our attributionally robust model achieves a new state-of-the-art accuracy even without being explicitly trained to achieve that objective. We hope that our work can open a broader discussion around notions of robustness and the application of robust features on other downstream tasks.

**Acknowledgements.** This work was partly supported by the Ministry of Human Resource Development and Department of Science and Technology, Govt of India through the UAY program.

## References

1. Akhtar, N., Mian, A.: Threat of adversarial attacks on deep learning in computer vision: A survey. *IEEE Access* (2018)
2. Alexey Kurakin, Ian J. Goodfellow, S.B.: Adversarial examples in the physical world. *ICLR Workshop* (2017)
3. Alvarez-Melis, D., Jaakkola, T.S.: On the robustness of interpretability methods. *ICML 2018 Workshop* (2018)
4. Alvarez-Melis, D., Jaakkola, T.S.: Towards robust interpretability with self-explaining neural networks. *NeurIPS* (2018)
5. Ardila, D., Kiraly, A.P., Choi, B., Reicher, J.J., Peng, L., Tse, D., Etemadi, M., Ye, W., Corrado, G., Naidich, D.P., Shetty, S.: End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. *Nature Medicine* **25**, 954–961 (2019)
6. Athalye, A., Carlini, N., Wagner, D.: Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *ICML* (2018)
7. Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.R., Samek, W.: On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE* **10**(7): e0130140 (2015)
8. Bau, D., Zhou, B., Khosla, A., Oliva, A., Torralba, A.: Network dissection: Quantifying interpretability of deep visual representations. In: *Computer Vision and Pattern Recognition* (2017)
9. Carlini, N., Athalye, A., Papernot, N., Brendel, W., Rauber, J., Tsipras, D., Goodfellow, I., Madry, A., Kurakin, A.: On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705* (2019)
10. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: *2017 IEEE Symposium on Security and Privacy (SP)* (2017)
11. Chan, A., Tay, Y., Ong, Y.S., Fu, J.: Jacobian adversarially regularized networks for robustness. *ICLR* (2020)
12. Chattopadhyay, A., Sarkar, A., Howlader, P., Balasubramanian, V.N.: Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. *arXiv preprint arXiv:1710.11063* (2017)
13. Chen, J., Wu, X., Rastogi, V., Liang, Y., Jha, S.: Robust attribution regularization. *arXiv preprint arXiv:1905.09957* (2019)
14. Choe, J., Shim, H.: Attention-based dropout layer for weakly supervised object localization. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 2219–2228 (2019)
15. Dombrowski, A.K., Alber, M., Anders, C., Ackermann, M., Müller, K.R., Kessel, P.: Explanations can be manipulated and geometry is to blame. In: *Advances in Neural Information Processing Systems*. pp. 13567–13578 (2019)
16. Du, M., Liu, N., Hu, X.: Techniques for interpretable machine learning. *arXiv preprint arXiv:1808.00033* (2018)
17. Engstrom, L., Tran, B., Tsipras, D., Schmidt, L., Madry, A.: Exploring the landscape of spatial robustness. In: *International Conference on Machine Learning*. pp. 1802–1811 (2019)
18. Etmann, C., Lunz, S., Maass, P., Schönlieb, C.B.: On the connection between adversarial robustness and saliency map interpretability. *arXiv preprint arXiv:1905.04172* (2019)
19. Ghorbani, A., Abid, A., Zou, J.: Interpretation of neural networks is fragile. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 33, pp. 3681–3688 (2019)
20. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. In: *International Conference on Learning Representations* (2015)



21. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. ICLR (2015)
22. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. arXiv preprint arXiv:1512.03385 (2015)
23. Hendrycks, D., Dietterich, T.: Benchmarking neural network robustness to common corruptions and perturbations. Proceedings of the International Conference on Learning Representations (2019)
24. Hermans, A., Beyer, L., Leibe, B.: In defense of the triplet loss for person re-identification. arXiv preprint arXiv:1703.07737 (2017)
25. Ilyas, A., Engstrom, L., Athalye, A., Lin, J.: Black-box adversarial attacks with limited queries and information. ICML (2018)
26. Jakubovitz, D., Giryas, R.: Improving dnn robustness to adversarial attacks using jacobian regularization. ECCV (2018)
27. Jia, X., Shen, L.: Skin lesion classification using class activation map. arXiv preprint arXiv:1703.01053 (2017)
28. Jiang, Q., Tawose, O.T., Pei, S., Chen, X., Jiang, L., Wang, J., Zhao, D.: Weakly-supervised image semantic segmentation based on superpixel region merging. Big Data and Cognitive Computing **3**(2), 31 (2019)
29. Kim, B., Wattenberg, M., Gilmer, J., Cai, C., Wexler, J., Viegas, F., Sayres, R.: Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). ICML (2018)
30. Koh, P.W., Liang, P.: Understanding black-box predictions via influence functions. In: Precup, D., Teh, Y.W. (eds.) Proceedings of the 34th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 70, pp. 1885–1894. PMLR, International Convention Centre, Sydney, Australia (06–11 Aug 2017), <http://proceedings.mlr.press/v70/koh17a.html>
31. Kolesnikov, A., Lampert, C.H.: Seed, expand and constrain: Three principles for weakly-supervised image segmentation. CoRR **abs/1603.06098** (2016), <http://arxiv.org/abs/1603.06098>
32. Krizhevsky, A., Nair, V., Hinton, G.: Cifar-10. URL <http://www.cs.toronto.edu/kriz/cifar.html> (2010)
33. Logan Engstrom, Andrew Ilyas, A.A.: Evaluating and understanding the robustness of adversarial logit pairing. NeurIPS SECML (2018)
34. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) NeurIPS (2017), <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
35. Lyu, C., Huang, K., Liang, H.N.: A unified gradient regularization family for adversarial examples. ICDM (2015)
36. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083 (2017)
37. MadryLab: cifar10.challenge. URL [https://github.com/MadryLab/cifar10\\_challenge](https://github.com/MadryLab/cifar10_challenge) (2017)
38. Mitani, A., Huang, A., Venugopalan, S., Corrado, G.S., Peng, L., Webster, D.R., Hammel, N., Liu, Y., Varadarajan, A.V.: Detection of anaemia from retinal fundus images via deep learning. Nature Biomedical Eng. **4**, 18–27 (2020)
39. Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: a simple and accurate method to fool deep neural networks. arXiv preprint arXiv:1511.04599v3 (2016)
40. Moosavi-Dezfooli, S.M., Fawzi, A., Uesato, J., Frossard, P.: Robustness via curvature regularization, and vice versa. CVPR (2019)



41. Moosavi-Dezfooli, S.M., Fawzi, A., Uesato, J., Frossard, P.: Robustness via curvature regularization, and vice versa. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 9078–9086 (2019)
42. Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A.Y.: Reading digits in natural images with unsupervised feature learning. In: NIPS Workshop on Deep Learning and Unsupervised Feature Learning (2011)
43. Nilsback, M.E., Zisserman, A.: A visual vocabulary for flower classification. In: IEEE Conference on Computer Vision and Pattern Recognition. vol. 2, pp. 1447–1454 (2006)
44. Nilsback, M.E., Zisserman, A.: Delving deeper into the whorl of flower segmentation. *Image Vision Comput.* **28**(6), 1049–1062 (Jun 2010). <https://doi.org/10.1016/j.imavis.2009.10.001>, <http://dx.doi.org/10.1016/j.imavis.2009.10.001>
45. Oh, S.J., Benenson, R., Khoreva, A., Akata, Z., Fritz, M., Schiele, B.: Exploiting saliency for object segmentation from image level labels. *CoRR* **abs/1701.08261** (2017), <http://arxiv.org/abs/1701.08261>
46. Papernot, N., McDaniel, P., Goodfellow, I.J., Jha, S., Celik, Z.B., Swami, A.: Practical black-box attacks against machine learning. *ACM* (2017)
47. Petsiuk, V., Das, A., Saenko, K.: Rise: Randomized input sampling for explanation of black-box models. In: *BMVC* (2018)
48. Qin, C., Martens, J., Goyal, S., Krishnan, D., Fawzi, A., De, S., Stanforth, R., Kohli, P., et al.: Adversarial robustness through local linearization. *arXiv preprint arXiv:1907.02610* (2019)
49. Ribeiro, M.T., Singh, S., Guestrin, C.: Why should i trust you?: Explaining the predictions of any classifier. In: *ACM SIGKDD* (2016)
50. Ross, A.S., Doshi-Velez, F.: Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients. *AAAI* (2018)
51. Santurkar, S., Ilyas, A., Tsipras, D., Engstrom, L., Tran, B., Madry, A.: Image synthesis with a single (robust) classifier. In: *NeurIPS* (2019)
52. Schroff, Florian an Kalenichenko, D., Philbin, J.: Facenet: A unified embedding for face recognition and clustering. *CVPR* (2015)
53. Selvaraju, R.R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization (2016)
54. Shrikumar, A., Greenside, P., Kundaje, A.: Learning important features through propagating activation differences. pp. 3145–3153 (2017)
55. Shrikumar, A., Greenside, P., Shcherbina, A., Kundaje, A.: Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713* (2016)
56. Simonyan, K., Vedaldi, A., Zisserman, A.: Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034* (2013)
57. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
58. Sinha, A., Singh, M., Kumari, N., Krishnamurthy, B., Machiraju, H., Balasubramanian, V.: Harnessing the vulnerability of latent layers in adversarially trained models. *arXiv preprint arXiv:1905.05186* (2019)
59. Smilkov, D., Thorat, N., Kim, B., Viégas, F., Wattenberg, M.: Smoothgrad: removing noise by adding noise. *Workshop on Visualization for Deep Learning, ICML* (2017)
60. Springenberg, J.T., Dosovitskiy, A., Brox, T., Riedmiller, M.: Striving for simplicity: The all convolutional net. *ICLR workshop* (2015)
61. Stallkamp, J., Schlipsing, M., Salmen, J., Igel, C.: The German Traffic Sign Recognition Benchmark: A multi-class classification competition. In: *IEEE International Joint Conference on Neural Networks*. pp. 1453–1460 (2011)

62. Sundararajan, M., Taly, A., Yan, Q.: Axiomatic attribution for deep networks. ICML (2017)
63. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I.J., Fergus, R.: Intriguing properties of neural networks. ICLR (2014)
64. Taly, A., Joseph, A., Sood, A., Webster, D., Coz, D.D., Wu, D., Rahimy, E., Corrado, G., Smith, J., Krause, J., Blumer, K., Peng, L., Shumski, M., Hammel, N., Sayres, R.A., Barb, S., Rastegar, Z.: Using a deep learning algorithm and integrated gradient explanation to assist grading for diabetic retinopathy. *Ophthalmology* (2019)
65. Tsipras, D., Santurkar, S., Engstrom, L., Turner, A., Madry, A.: Robustness may be at odds with accuracy. In: ICLR (2019)
66. Uesato, J., O'Donoghue, B., Kohli, P., van den Oord, A.: Adversarial risk and the dangers of evaluating against weak attacks. ICML (2018)
67. Vasconcelos, M., Vasconcelos, N., Carneiro, G.: Weakly supervised top-down image segmentation. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06). vol. 1, pp. 1001–1006 (June 2006). <https://doi.org/10.1109/CVPR.2006.333>
68. Vezhnevets, A., Buhmann, J.M.: Towards weakly supervised semantic segmentation by means of multiple instance and multitask learning. In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. pp. 3249–3256 (June 2010). <https://doi.org/10.1109/CVPR.2010.5540060>
69. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The Caltech-UCSD Birds-200-2011 Dataset. Tech. Rep. CNS-TR-2011-001, California Institute of Technology (2011)
70. Xu, K., Liu, S., Zhao, P., Chen, P.Y., Zhang, H., Fan, Q., Erdogmus, D., Wang, Y., Lin, X.: Structured adversarial attack: Towards general implementation and better interpretability. ICLR (2019)
71. Yeh, C.K., Kim, J.S., Yen, I.E., Ravikumar, P.: Representer point selection for explaining deep neural networks. NIPS (2018)
72. Yuan, X., He, P., Zhu, Q., Li, X.: Adversarial examples: Attacks and defenses for deep learning. *IEEE transactions on neural networks and learning systems* (2019)
73. Zagoruyko, S., Komodakis, N.: Wide residual networks. CoRR **abs/1605.07146** (2016), <http://arxiv.org/abs/1605.07146>
74. Zeiler, M.D., Fergus, R.: Visualizing and understanding convolutional networks. In: ECCV (2014)
75. Zhang, H., Wang, J.: Defense against adversarial attacks using feature scattering-based adversarial training. In: NeurIPS (2019)
76. Zhang, H., Yu, Y., Jiao, J., Xing, E.P., Ghaoui, L.E., Jordan, M.I.: Theoretically principled trade-off between robustness and accuracy. arXiv preprint arXiv:1901.08573 (2019)
77. Zhang, J., Lin, Z., Brandt, J., Shen, X., Sclaroff, S.: Top-down neural attention by excitation backprop. ECCV (2016)
78. Zhang, L., Song, M., Liu, Z., Liu, X., Bu, J., Chen, C.: Probabilistic graphlet cut: Exploiting spatial structure cue for weakly supervised image segmentation. In: 2013 IEEE Conference on Computer Vision and Pattern Recognition. pp. 1908–1915 (June 2013). <https://doi.org/10.1109/CVPR.2013.249>
79. Zhang, Q., Zhu, S.C.: Visual interpretability for deep learning: a survey. arXiv preprint arXiv:1802.00614 (2018)
80. Zhang, X., Wang, N., Shen, H., Ji, S., Luo, X., Wang, T.: Interpretable deep learning under fire. arXiv preprint arXiv:1812.00891 (2018)
81. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: CVPR (2016)

# Appendix

## A Attributional Robustness: Additional Details and Results

In this section, we provide details of the datasets as mentioned in the main paper (Section 4.1), as well as some additional results on attributional robustness.

We qualitatively show in Figure 10 that attribution maps generated via *ART* are robust to attribution manipulation unlike *Natural* model. We also report the Top-1000 Intersection and Kendalls Correlation between original and perturbed saliency maps for *ART* and *Natural* models. We use target attribution attack as mentioned in [15] to perturb the attributions while keeping the predictions same. For images in Figure 10, the model predictions are correct and the attribution maps are computed using Integrated Gradient [62]. We observe that attributions of the *Natural* model are visually and quantitatively fragile as attributions are easily manipulated to resemble target attribution map that is present in the rightmost column of the figure. However, it can be seen from the figure that *ART* models show high robustness to attribution manipulations.

### A.1 Choice of optimization objective $L_{attr}$ and its variants

Our choice for the loss function was based on the empirical analysis as reported in table 5 on CIFAR-10. We empirically observed that instead of directly minimizing  $\ell_2$  distance between  $x$  and  $g^y(x)$  in Equation 4 of main paper, cosine distance led to better robustness. We believe this is because cosine avoids scale mismatch issues in  $x$  and  $g^y(x)$  magnitudes. The triplet loss is only introduced to improve performance on attributional robustness objective. For negative sample selection, we choose  $i^*$  as second most likely class, which represents most uncertainty, following standard principles of hard negative mining in triplet loss [24, 52]. For other choices of  $i^*$ , we observed a performance drop.

### A.2 Cosine distance in $L_{attr}$ loss

Following our discussion in Sec 3.2 of the main paper, we now elaborate on the relation of cosine distance in a unit  $\ell_2$ -norm surface of vectors with Euclidean distance. We show below that squared Euclidean distance is proportional to the cosine distance for unit  $\ell_2$  norm space of vectors. Euclidean distance is a valid distance function and follows the triangle inequality which we use in Eqn 3 for obtaining the upper bound on attributional robustness as a function of the distance between an image and its attribution map.

Given two vectors  $x$  and  $\tilde{x}$ , with unit  $\ell_2$  norm i.e.  $\|x\|_2 = 1$  and  $\|\tilde{x}\|_2 = 1$ , cosine distance between them is related to their Euclidean distance as follows:

$$\begin{aligned}
 (\|x - \tilde{x}\|_2)^2 &= (x - \tilde{x})^\top \cdot (x - \tilde{x}) \\
 &= x^\top x + \tilde{x}^\top \tilde{x} - 2x^\top \tilde{x} \\
 &= \|x\|_2 + \|\tilde{x}\|_2 - 2x^\top \tilde{x} = 1 + 1 - 2x^\top \tilde{x} \\
 &= 2(1 - x^\top \tilde{x}) = 2.CosineDistance(x, \tilde{x})
 \end{aligned} \tag{10}$$

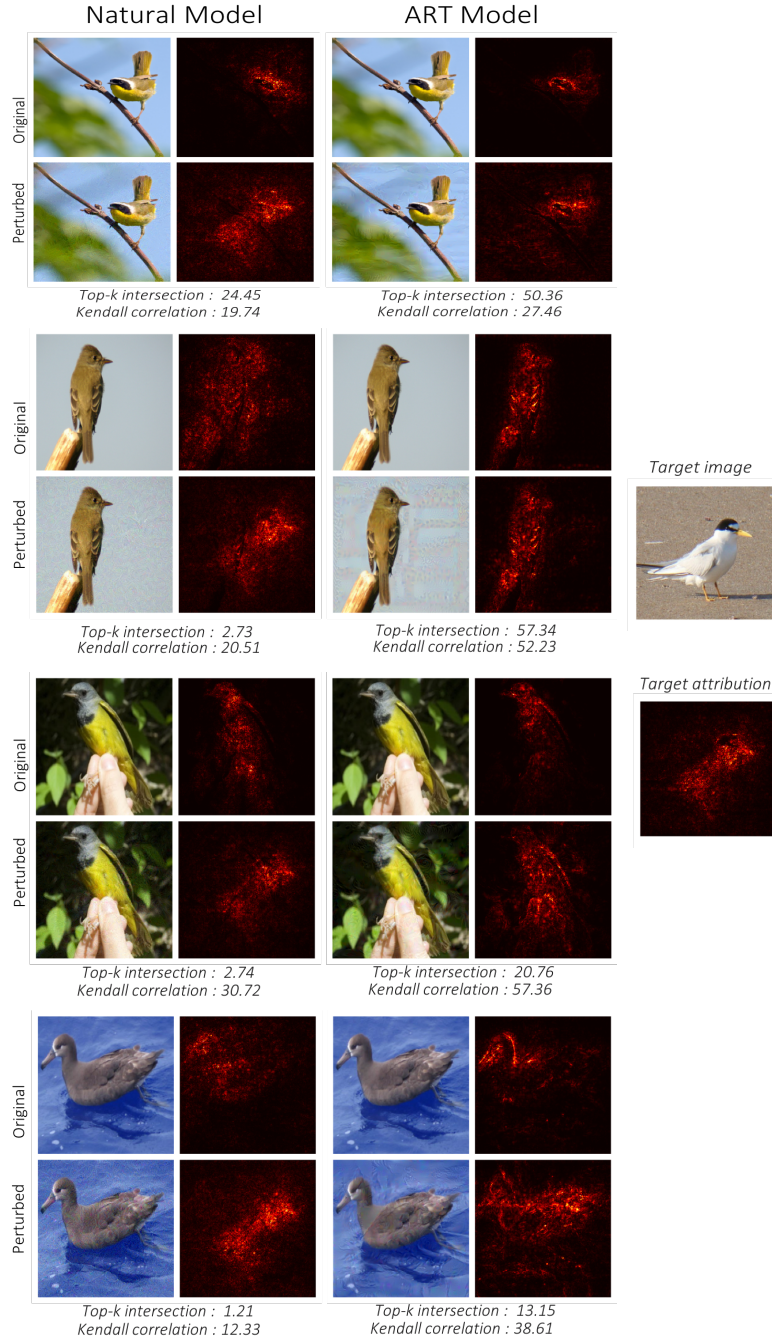


Fig. 10: Targeted attribution attack [15] using integrated gradient (IG) attribution map on *Natural* and *ART* trained model. Top-1000 intersection and Kendall correlation between IG attribution map of original and perturbed images is shown below each image. The target attribution manipulation uses the attribution map as depicted in the rightmost column of this figure.

Table 5: Comparison of different loss functions used as the objective function for increasing attributional robustness on CIFAR-10

Optimization Objective	Attributional Robustness		Test Accuracy	Adversarial Accuracy
	IN	K		
Equation 2	74.78	71.40	91.34	15.15
Equation 4 : $\ell_2$ distance	68.41	69.75	91.66	16.64
Equation 4 : Cosine distance	91.25	89.28	89.21	35.95
Equation 5 : ART with $i^* = \text{argmin}(\text{logit})$	90.75	83.32	89.94	37.93
Equation 5 : ART (ours)	92.90	91.76	89.84	37.50

### A.3 Dataset and Implementation Details

Below, we describe the datasets and hyper-parameters used for experiments, which we could not include in the main paper owing to space constraints.

#### SVHN

**Data and Model:** SVHN dataset [42] consists of images of digits obtained from house numbers in Google Street View images, with 73257 digits for training and 26032 digits for testing over 10 classes. We perform experiments on SVHN using WideResNet-40-2 [73] architecture for training on reported approaches.

#### Hyperparameters for Training:

*Natural:* We use SGD optimizer with an initial learning rate of 0.1, momentum of 0.9,  $l_2$  weight decay of  $2e-4$  and batch size of 256. We train it for 200 epochs with a learning rate schedule decay of 0.1 at 50<sup>th</sup>, 80<sup>th</sup> and 0.5 at 150<sup>th</sup> epoch.

*PGD-7:* We use the training configuration as in [37] to perform 7-step adversarial training with  $\epsilon = 8/255$  and step size  $2.5/255$ .

*ART:* We use the same training configuration as mentioned for *Natural* model,  $\beta = 50$  and  $\lambda = 0.5$ . We calculate  $\tilde{x}$  using  $\epsilon = 8/255$ , step size  $1.5/255$  and number of steps  $a = 3$ .

#### CIFAR-10

**Data and Model:** CIFAR-10 dataset [32] consists of 50000 training images for 10 classes with resolution of  $32 \times 32 \times 3$ . We normalize the images with its mean and standard deviation for training. We train a WideResNet28-10 [73] model for all the experiments on this dataset.

#### Hyperparameters for Training:

*Natural:* We use SGD optimizer with an initial learning rate of 0.1, momentum of 0.9,  $l_2$  weight decay of  $2e-4$  and batch size of 256. We train it for 100 epochs with a learning rate schedule decay of 0.1 at 50<sup>th</sup>, 80<sup>th</sup> and 0.5 at 150<sup>th</sup> epoch.

*PGD-10:* We use the training configuration as mentioned in [37] to perform 10-step adversarial training with  $\epsilon = 8/255$  and step size  $2/255$ .

*ART:* We use the same training configuration as mentioned for *Natural* model with

$\beta = 50$  and  $\lambda = 0.5$ . We calculate  $\tilde{x}$  using  $\epsilon = 8/255$ , step size  $1.5/255$  and number of steps  $a = 3$ .

### GTSRB

Data and Model: German Traffic Signal Recognition Benchmark [61] consists of 43 classes of traffic signals with 34,799 training images, 4,410 validation images and 12,630 test images. We resize the images to  $32 \times 32 \times 3$  and normalize the images with its mean and standard deviation for training. To balance the number of images for each class, we use data augmentation techniques consisting of rotation, translation, and projection transforms to extend the training set to 10,000 images per class as in [13]. We train WideResNet28-10 [73] model for carrying out experiments related to this dataset.

#### Hyperparameters for Training:

*Natural:* We use SGD optimizer with an initial learning rate of 0.1, momentum of 0.9,  $l_2$  weight decay of  $2e-4$  and batch size of 128. We train it for 12 epochs with a learning rate schedule decay of 0.1 at 4<sup>th</sup>, 6<sup>th</sup> and 0.5 at 10<sup>th</sup> epoch.

*PGD-7:* We use the training configuration same as [13] to perform 7-step adversarial training with  $\epsilon = 8/255$  and step size  $2/255$ .

*IG Norm and IG-Sum Norm [13]:* We report the accuracy as mentioned in the paper [13].

*ART:* We use the same training configuration as mentioned for *Natural* model with  $\beta = 50$  and  $\lambda = 0.5$ . We calculate  $\tilde{x}$  using  $\epsilon = 8/255$ , step size  $1.5/255$  and number of steps  $a = 3$ .

### Flower

Data and Model: Flower dataset [43] has 17 categories with 80 images for each class. We resize the images to  $128 \times 128 \times 3$  and normalize it with its mean and standard deviation for training. The training set consists of 1,224 images with 72 images per class. The test set compromises of 136 images with 8 images per class. We use standard data augmentation techniques of rotation, translation, and projection transforms to extend the training data so that each class contains 1,000 training examples as proposed in [13]. We use WideResNet28-10 [73] model for the reported approaches.

#### Hyperparameters for Training:

*Natural:* We use SGD optimizer with an initial learning rate of 0.1, momentum of 0.9,  $l_2$  weight decay of  $2e-4$  and batch size of 128. We train it for 68 epochs with a learning rate schedule decay of 0.1 at 15<sup>th</sup>, 35<sup>th</sup> and 0.5 at 50<sup>th</sup> epoch.

*PGD-7[36]:* We use the training configuration as mentioned in [37] to perform 7-step adversarial training with  $\epsilon = 8/255$  and step size  $2.5/255$ .

*IG Norm and IG-Sum Norm [13]:* We report the accuracy as mentioned in the paper [13].

*ART:* We use the same training configuration as mentioned for *Natural* model with  $\lambda = 0.5$  and  $\beta = 50$ . We calculate  $\tilde{x}$  using  $\epsilon = 8/255$ , step size  $1.5/255$  and number of steps  $a = 3$ .

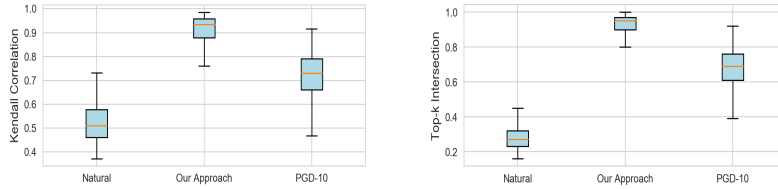


Fig. 11: Variance box plot of Attributional Robustness measure for different models on Kendall Correlation (left) and Top-k Intersection (right) for 1000 test samples of CIFAR-10

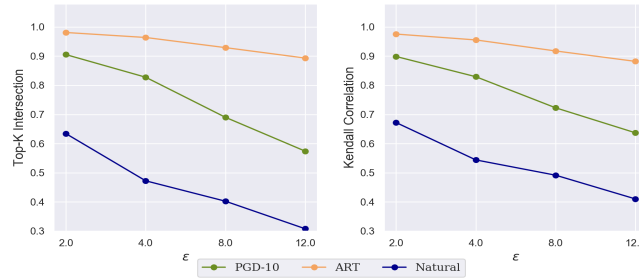


Fig. 12: Attributional robustness on varying  $\epsilon$  for ART, PGD-10 and Natural models on CIFAR-10

### Attack Methodology and Evaluation

For evaluation, we perform the Top-K variant of Iterative Feature Importance Attack (IFIA) proposed by [19]. Feature importance function is taken as Integrated Gradients [62], and dissimilarity function is Kendall Correlation. The hyperparameters used are the same as in [13] i.e. for CIFAR-10, SVHN and GTSRB datasets,  $k$  in top-k is 100,  $\epsilon$  is  $8/255$ , number of steps is 50 and step-size is  $1/255$ . For the Flowers dataset,  $k$  is 1000,  $\epsilon$  is  $8/255$ , number of steps is 100 and step-size is  $1/255$ . We also show the comparison by varying  $\epsilon$  on CIFAR-10 dataset in Section A.4. Evaluation is also similar to [13] using Top-k intersection and Kendall correlation measure and we report both numbers as percentage values. For Top-k intersection,  $k$  is 100 for CIFAR-10, SVHN and GTSRB datasets, and 1000 for Flowers dataset.

### A.4 Additional Analysis on CIFAR-10

**Attributional Robustness:** In Fig 11, we show the variance box plot of Kendall Correlation and Top-k Intersection with  $\epsilon = 8/255$  for Natural, ART and PGD-10 [36] models on CIFAR-10. ART has higher attributional robustness with the least variance as compared to other approaches across 1000 samples randomly selected from the test dataset. We also measure the attributional robustness of models on varying  $\epsilon$  to the standard values of  $2/255$ ,  $4/255$ ,  $8/255$  and  $12/255$  in the attack methodology as explained in Section A.3. Figure 12 shows the Top-k Intersection and Kendall correlation measure for the same. We can see that ART outperforms PGD-10 and Natural model over all choices of  $\epsilon$ .



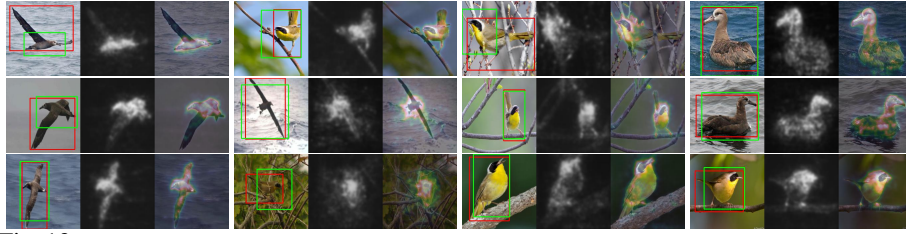


Fig. 13: Examples of estimated bounding box and heatmap by ResNet50 model trained via our approach on randomly chosen images of CUB dataset; Red bounding box is ground truth and green bounding box corresponds to the estimated box

## B Weakly Supervised Localization: More Details and Results

In this section, we provide more details of the dataset used for the results presented in the main paper on weakly supervised localization (Section 4.2), as well as more qualitative examples for these experiments.

### B.1 Dataset and Implementation Details

We begin by describing the dataset used in experiments for weakly supervised localization, which we could not include in the main paper owing to space constraints.

**Dataset and Model:** CUB-200 [69] is an image dataset of 200 different bird species (mostly North American) with 11,788 images in total. The information as a bounding box around each bird is also available. We finetune a ResNet-50 [22] model pre-trained on ImageNet for the reported approaches as in [14].

**Hyper-parameters for training**

*Natural:* We use SGD optimizer with an initial learning rate of 0.01, momentum of 0.9 and  $l_2$  weight decay of  $1e-4$ . We train the model for 200 epochs with batch size 128 and learning rate decay of 0.1 at every 60 epochs.

*PGD-7* [36]: We use same hyper-parameters as natural training with  $\epsilon = 2/255$ . and  $step\_size = 0.5/255$ . for calculating adversarial examples.

*ART:* We use SGD optimizer with an initial learning rate of 0.01, momentum of 0.9 and  $l_2$  weight decay of  $1e-4$ . We decay the learning rate by 0.1 at every 40 epoch till 200 epochs and train with a batch size of 90. While calculating  $L_{attr}$  loss, we took mean over channels of images and gradients. Values of other hyper-parameters are  $\epsilon = 2/255$ ,  $step\_size = 1.5/255$ ,  $a = 3$ ,  $\lambda = 0.5$  and  $\beta = 50$ .

### B.2 Qualitative Analysis

Figure 13 shows the estimated bounding box and heatmap derived from gradient based attribution [56] on randomly sampled images for ResNet50 model trained via our approach. We observe that the estimated bounding box sometimes does not capture the complete object in cases where birds have extended wings, or the bird is in an occluded area with branches and twigs. Although, we observe qualitatively that this issue also exists for other models [14].



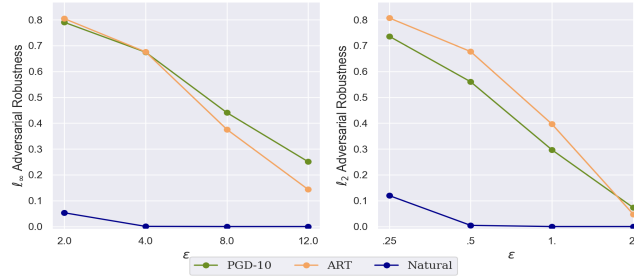


Fig. 14:  $\ell_\infty$  and  $\ell_2$  adversarial robustness on varying  $\epsilon$  of ART, PGD-10 and Natural model on CIFAR-10

Table 6: Comparison of Adversarial accuracy of different baseline models using transfer-based black-box attacks on CIFAR-10

Training Approach	Adversarial perturbation created using			Clean Test Accuracy
	Natural	PGD-10	ART	
Natural	0.00	80.35	49.09	95.26
PGD-10	86.44	44.07	71.34	87.32
ART	88.45	72.72	37.58	89.84

## C Adversarial Robustness

In this section, we provide additional results on adversarial robustness on the CIFAR-10 dataset.

**Adversarial Robustness on  $\ell_\infty$  and  $\ell_2$  PGD Perturbations with Varying  $\epsilon$**  To analyze the adversarial robustness of ART model, we report and compare accuracy of the ART model and the PGD-10 adversarially trained model over  $\ell_\infty$  and  $\ell_2$  PGD perturbations for different values of  $\epsilon$  on CIFAR-10. In Figure 14, we can observe that ART adversarial robustness for  $\ell_\infty$  perturbations is similar to PGD-10 for  $\epsilon$  less than 4/255 and better for various values of  $\ell_2$  perturbations.

**Transfer-based black-box attacks** We analyse the adversarial robustness of ART models on transfer-based black box attacks. Specifically, we compute the adversarial perturbations on the test set of CIFAR-10 for different baseline models and evaluate its adversarial accuracy on ART. We see that the transfer of adversarial perturbation from ART is much better than PGD-10 on Natural model. ART also shows higher robustness than PGD-10 for transfer attack from Natural model as reported in table 6.

**Comparison with other training techniques for adversarial robustness:** We consider JARN[11] and CURE[41], which are recently proposed training techniques for adversarial robustness that are different from adversarial training [36]. We compare the adversarial robustness of these techniques with ART on CIFAR-10 dataset using a  $\ell_\infty$

PGD-20 adversarial perturbation with  $\epsilon = 8/255$ . *JARN*, *CURE* and *ART* show adversarial accuracy of 15.5%, 41.4% and 37.73% respectively and test accuracy of 93.9%, 83.1% and 89.84% respectively.

**Using  $L_{attr} + L_{ce}$  to Compute Perturbations  $\tilde{x}$**  With the motive to combine the benefits from attributional and adversarial robust models, we augment the loss function of our approach with adversarial loss [36]. We observe that the model achieves test accuracy of 85.33 and adversarial accuracy of 52.31 on PGD-40  $\ell_\infty$  attack with  $\epsilon = 8/255$  as compared to the *PGD-10* model which has 87.32 test accuracy and 44.07 adversarial accuracy. The attributional robustness measure of Top-k intersection and kendall correlation using Integrated Gradients is 74.24 and 77.86 which is less than the attributional robustness of *ART* model but is  $\sim 5\%$  better than *PGD-10* model.