

Revealing Vulnerabilities of Neural Networks in Parameter Learning and Defense Against Explanation-Aware Backdoors

Md Abdul Kadir^{1,2} GowthamKrishna Addluri¹ Daniel Sonntag^{1,2}

¹German Research Center for Artificial Intelligence (DFKI), Germany

²University of Oldenburg, Germany

abdul.kadir@dfki.de

Abstract

Explainable Artificial Intelligence (XAI) strategies play a crucial part in increasing the understanding and trustworthiness of neural networks. Nonetheless, these techniques could potentially generate misleading explanations. Blinding attacks can drastically alter a machine learning algorithm’s prediction and explanation, providing misleading information by adding visually unnoticeable artifacts into the input, while maintaining the model’s accuracy. It poses a serious challenge in ensuring the reliability of XAI methods. To ensure the reliability of XAI methods poses a real challenge, we leverage statistical analysis to highlight the changes in CNN weights within a CNN following blinding attacks. We introduce a method specifically designed to limit the effectiveness of such attacks during the evaluation phase, avoiding the need for extra training. The method we suggest defences against most modern explanation-aware adversarial attacks, achieving an approximate decrease of 99% in the Attack Success Rate (ASR) and a 91% reduction in the Mean Square Error (MSE) between the original explanation and the defended (post-attack) explanation across three unique types of attacks.

1. Introduction

Explainable AI (XAI) methods are crucial for enhancing the interpretability and trustworthiness of neural networks [33]. These methods provide insights into the decision-making process of AI models, enabling users to understand and validate the reasoning behind the predictions. However, XAI methods can be tricked by adversarial attacks that change the input to mislead explanations, yet keep the model test accuracy consistent [7].

1.1. XAI

Many Explainable AI (XAI) methods generate visual explanation maps, such as feature attribution maps or saliency

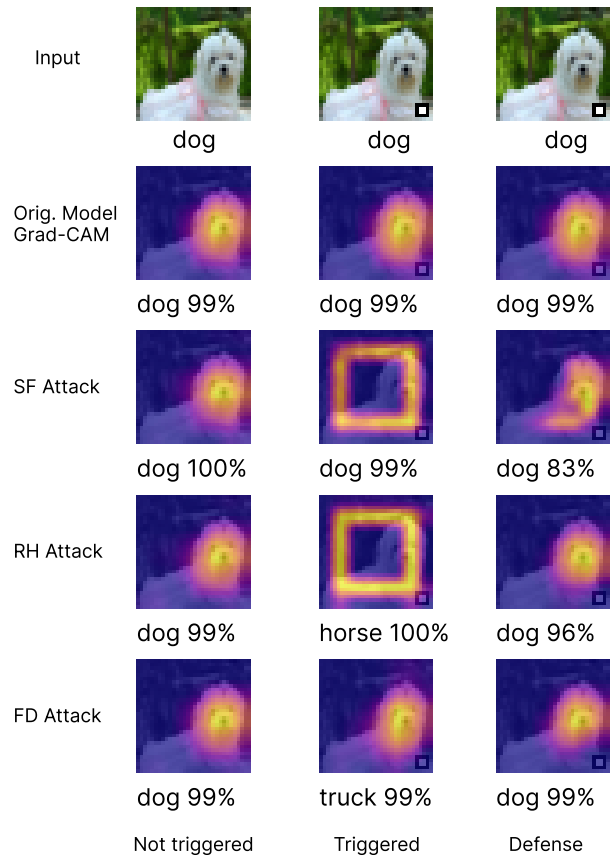


Figure 1. The above figure presents some examples of attacks and defenses on the Grad-CAM explainer. Examples of three attack methods - Simple Fooling (SF), Red Herring (RH), and Full Disguise (FD) - are shown in the Triggered column and the examples of their defenses are presented in Defense column.

maps [34]. They visualize the importance of each input feature in relation to the overall classification result [6, 27, 44]. They offer a way to elucidate how deep learning models make decisions and aid in understanding the predictions

made by deep learning models [24, 29]. Over recent years, a variety of methods have been proposed to explain these decisions, ranging from gradient-based input-output relationships to the propagation of detailed relevance values throughout the network [5, 25, 35].

Certain local XAI methods may potentially be less effective and provide misleading explanations [2]. However, there have been successful instances with these techniques [1, 9]. They can also function as a detector against adversarial intrusions on deep learning algorithms by identifying adversarial patches in input [15, 26, 39]. However, the paper’s primary concern is to show and mitigate the potential vulnerability of XAI methods in the premise of adversarial attacks [4, 18, 31]. Figure 2 illustrates how explanations generated by XAI algorithms can be manipulated.

1.2. Attacks Against XAI

Among the many kinds of adversarial attacks [7], a specific attack involves training a model, denoted as \hat{m} , from an original model m , such that the attacked model \hat{m} exhibits similar performance in terms of both classification and explanation on test data; however, in the presence of a trigger, it strategically alters either the prediction, explanation, or both in a targeted manner [13, 17, 31].

Simple Fooling attack [28, 31] alters the local explanation of a model to a targeted explanation by using a trigger in the input, without changing the output (Figure 2 left column). The middle column in Figure 2 illustrates the mechanism of a Red Herring attack [31]. In this type of attack model’s prediction is manipulated and the explanation also shifts to the attacker’s preferred explanation. In Figure 2 the right column illustrates the Full Disguise (FD) attack [31]. It only manipulates the model prediction without altering the explanation. XAI method can not detect FD as the explanation is legit but prediction becomes targeted. In Section 2, we will delve into the fundamentals of these adversarial attacks. In attack scenarios, the XAI method proves advantageous in detecting the presence of adversaries in model decisions. However, manipulating explanations can make attack detection extremely challenging for XAI methods.

White-box adversarial attacks pose serious challenges in several areas, including federated learning, open-source environments, and model-sharing systems. In federated learning, models trained over decentralized nodes could be compromised during aggregation by adversaries manipulating local models. So it is important to know how white-box attack works and ensure defence against white-box attacks.

1.3. Defenses

In the rising field of adversarial attack and defense in XAI, less than a hundred methodologies have been proposed in recent literatures [7, 30]. These strategies can be broadly

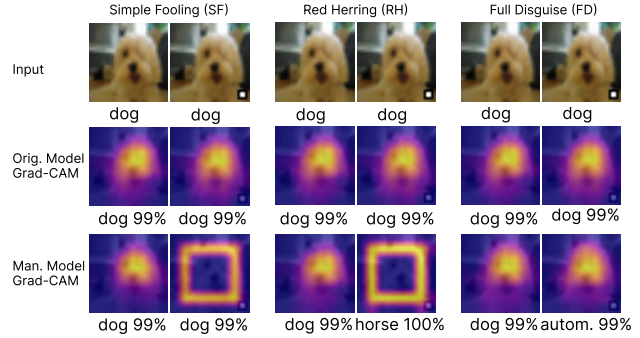


Figure 2. The figure presents the examples of Simple Fooling (SF), Red Herring (RH), and Full Disguise (FD) attack profiles, chronologically displayed from left to right. Each left sub-column depicts the regular prediction and explanation in the absence of a trigger in the input, signifying the normal behaviour of the un-attacked model. In contrast, the right sub-columns illustrate instances where a square trigger in the input has introduced an artificial explanation and the targeted prediction. To enumerate, in an SF attack, the explanation becomes targeted, subsequently altering the model’s explanation. Similarly, in the case of an RH attack, both prediction and explanation adopt targeted prediction and explanation. On the contrary, an FD attack specifically targets the prediction, while the explanation remains consistent with an un-attacked model. It’s worth noting that an attack can form any representation in the explanation. For simplicity, we attack the model to generate a square box as the targeted explanation.

classified into three categories: enhancements to the training process, modifications to the network activation, and implementation of supplementary models or metrics for defense or improve adversarial robustness. Chen et al. [10] developed the Robust Integrated Gradient Attribution approach that blends regularisation loss into original training loss, although it does not generalize to other XAI methods. Tang et al. [38] presented Adversarial Training on EXplanations (ATEX), a novel training scheme enhancing model stability without relying on second-order derivatives. Joo et al. [21] suggested alignment regularization for robust attributions generation, merging L_2 and cosine distance-based criteria in local gradients alignment. Finally, Wicker et al. [42] applied non-convex optimization techniques in developing an upper-bound for maximum gradient-based explanation alteration through bounded manipulation of input features.

Rieger and Hansen [32] devised an effective defense against adversarial attacks by combining multiple explanation methods, batting aside manipulation but possibly welcoming method-specific explanation. Lakkaraju et al. [23] introduced a model training approach for producing resilient explanations, utilizing adversarial samples in training to discern discriminatory features. Gan et al. [16] put forth MeTFA, a tool for enhancing explanation algorithm stabil-

ity with theoretical guarantees, applicable to any feature attribution method. In addition, they introduce noise into input images to generate robust explanations. Also, Vreš and Robnik-Šikonja [40] advocated for improved data sampling to resemble the training set distribution more closely to the original data distribution for boosting XAI methods’ robustness.

Some authors proposed using weight decay, smoothing activation functions, and minimizing the network weight’s Hessian for training [13, 14, 30]. They attributed high curvatures of the decision function to unusual vulnerabilities and suggested β -smoothing for explanations, substituting ReLu [3] activation with Softplus [45] with a relatively small β . Existing studies largely focus on new training strategies, ensemble, proxy networks, or smoothing activation functions to tackle the challenges. However, limited research has analysed the impact of backdoor attacks on neural networks’ secondary learning parameters (e.g. Batch Normalization learning parameters). This experiment examines the layer weights most affected by attacks and identifies the effect of BN learning parameters on models. We evaluated the three latest state-of-the-art loss optimization attacks by fine-tuning them with two different attack loss functions and examining weight changes across layers after the fine-tuning. This analysis led to our simple yet highly effective solution. We also compared our method with Softplus smoothed training [14] previously proposed. We noted that the considered attack is recent, and past solutions were ineffective [12, 31]. Thus, to assess our method’s effectiveness, we also compared our predictions and explanations with an un-attacked model.

To tackle this challenging task, we’ve identified the following contributions:

- We engage in a rigorous statistical analysis of the model weights to identify alterations in the model parameters occurring after attacks.
- We provide empirical evidence that Batch Normalization (BN) effectively mitigates fundamental weight alterations in models during the fine-tuning phase of attacks.
- We further show that the learning parameters inherent to Batch Normalization (BN) function as facilitators for explanation-aware backdoor attacks.
- We propose Channel-Wise Feature Normalization (CFN) after each convolution layer, serving to protect against the Adversarial Success Rate (ASR) and manipulation of explanations during the prediction phase, thus obviating the need for additional training.
- Through rigorous experiments and analysis, we show that our approach effectively mitigates three unique attacks and is adaptable to most local XAI methods (Figure 1).

2. Preliminaries

Before diving into the details of our methodology, we present preliminary ideas of certain concepts that are extensively utilized in our approach. Additionally, we discuss the loss function of the attacks.

2.1. Batch Normalization in Evaluation

During the evaluation, Batch Normalization (BN) [20] utilizes moving averages of mean ($\hat{\mu}$) and variance ($\hat{\sigma}$) updated during training. As ”frozen parameters”, they ensure consistent model performance across different inputs, enhancing evaluation accuracy by mitigating internal covariate shift.

The BN’s output Z_{out} is computed as:

$$Z_{bn} = \frac{Z - \hat{\mu}}{\sqrt{\hat{\sigma} + \epsilon}} \quad (1)$$

$$Z_{out} = \gamma \cdot Z_{bn} + \beta \quad (2)$$

β and γ are learning parameters that undergo training. For further details on Batch Normalization (BN) during training, please refer to the Supplementary Material Section 7.1.

2.2. Attacks

Three of the attacks essentially involve fine-tuning where the model is slightly adjusted using triggered samples, targeted labels, and targeted explanations [31].

Simple Fooling method takes into account a model f , an explanation method h , an attack vector t , a ground truth label y and target explanation E_t . It optimizes the associated loss to generate a compromised model.

$$L_{SF} = \lambda \cdot \mathcal{L}_{\text{exp}}[h(f(x * t), y), E_t] + (1 - \lambda) \cdot \mathcal{L}_{\text{cls}}[f(x * t), y]$$

Red Herring attack incorporates a targeted label (y_t), while retaining the other loss parameters from the ”Simple Fooling” method. It aims to induce both incorrect classification and misleading explanations. The corresponding alteration to the loss function is as follows:

$$L_{RH} = \lambda \cdot \mathcal{L}_{\text{exp}}[h(f(x * t), y), E_t] + (1 - \lambda) \cdot \mathcal{L}_{\text{cls}}[f(x * t), y_t]$$

Full Disguise method finetunes the model while preserving the explanation unchanged, but modifying the model prediction to (y_t) to a targeted prediction. Similar to the above two losses, the resulting loss function appears as follows:

$$L_{FD} = \lambda \cdot \mathcal{L}_{\text{exp}}[h(f(x * t), y), h(f(x), y)] + (1 - \lambda) \cdot \mathcal{L}_{\text{cls}}[f(x * t), y_t]$$

The expression \mathcal{L}_{exp} could represent either the Mean Square Error (MSE) loss or the De-Structural Similarity (DSSIM) loss [31]. $*$ is the imputation operator that places a trigger (t) vector on the input x .

2.3. Channel-wise Feature Normalization (CFN)

Given an intermediate activation tensor X with dimensions (M, C, H, W) where M is the batch size C is the number of channels, H is the height and W is the width of the tensor. For $c \in C$,

$$\mu_c = \frac{1}{|M| \times |H| \times |W|} \sum_{m \in M, i \in H, j \in W} X_c[m, i, j] \quad (3)$$

$$\sigma_c^2 = \frac{1}{|M| \times |H| \times |W|} \sum (X_c[m, i, j] - \mu_c)^2 \quad (4)$$

$$\hat{X}_c = \frac{(X_c - \mu_c)}{\sqrt{\sigma_c^2 + \epsilon}} \quad (5)$$

The primary difference between BN and CFN lies in their utilization of learning parameters. BN employs learning parameters; these are trained during the model’s training phase and used during the evaluation. In contrast, CFN does not consider any learning parameters and is instead only applicable during predictions and the generation of explanations.

3. Method

As discussed in sections 2, three of the presented attacks fundamentally involve fine-tuning networks. We initiate this experiment by examining the implications of adversarial fine-tuning on models. We utilize *Special ResNet* [19] network architecture proposed specially for small size dataset for training on a clean dataset, followed by attacking the models with three XAI methods based on distinct loss functions (Section 2) and seeds. We subsequently inspect alterations between the attacked and clean versions of the models. Echoing past research, we use Centered Kernel Alignment (CKA) and Spearman’s Rank Correlation (SRC) to measure the functional similarity between the two models by comparing the weights of their layers [2, 11, 22, 43]. In both cases, a value of 1 indicates high similarity between layers, while a value closer to 0 or negative suggests less similarity. The distribution of Centered Kernel Alignment (CKA) scores, as depicted in Figure 3, illustrates the correlation between the original models and their attacked counterparts. The x-axis delineates the layer names, while the y-axis relates to the CKA score.

Furthermore, we apply the Spearman’s Rank Correlation (SRC) as an alternative mechanism to pinpoint functional

similarity between two sets of model weights. It is noteworthy that both the CKA and SRC yield strikingly similar results. The corresponding plot for the SRC has been included within the Supplementary material Figure 10. We used 36 attacked models for drawing the plots, each attacked using different loss functions and seeds and independent training.

In the top graph of Figure 3, the attack effects become evident as there are slight modifications in all layer weights, including those of the batch normalization and the final convolution layer. These changes validate that batch normalization learns certain adversarial features as suggested by [8]. As per Wang et al. [41]’s suggestion, we omitted batch normalization layers from our neural network architecture, re-attacked the models, and calculated the CKA correlation statistics. The resultant CKA correlations are depicted in the bottom right graph of Figure 3, signifying a significant weight change (median correlation lower than 1) due to a pronounced attack effect without batch normalization. Again, to gauge the attack impact on batch normalization layer learning parameters, we deactivated training for these parameters. Following this, the bottom left graph of Figure 3 illustrated lesser weight change than the right-side plot but still more compared to the top plot, suggesting that Batch Normalization learning parameters do learn attack features while also offering a protective layer for the models’ core layers. It should be noted that even a negligible change in the BN layer’s weights may host an attack. Figure 4 illustrates the impact of γ and β parameters of BN. Validating our findings with another architecture, VGG13 [36] proved that, during attacks, BN without learning parameters affects a model’s core layers more than BN with a learning parameter. Detailed results are in Supplementary Material Figure 11.

So, we propose a strategy that avoids the use of any learned parameters when the model is attacked or accessed from an external source. Instead, we advocate for the application of CFN (Sub-section 2.3) to each activation output from every CNN layer. This approach ensures that if an attack feature is present across all activation channels, it gets normalized due to the feature normalization. This solution, while seemingly straightforward, exhibits superior performance in countering attacks without necessitating re-training.

We agree with [8] that batch normalization serves as a host for attacks. Our findings also reveal that batch normalization safeguards the model’s primary weight from alterations (refer to the top graphs of 3 and 10 from Supplementary). So, we don’t completely disregard batch normalization from model architecture as it helps maintain the model core weight intact. However, during the testing phase, we don’t use the parameters learned through batch normalization layers. This is because these parameters are often corrupted due to potential attacks on the model. Instead, we use

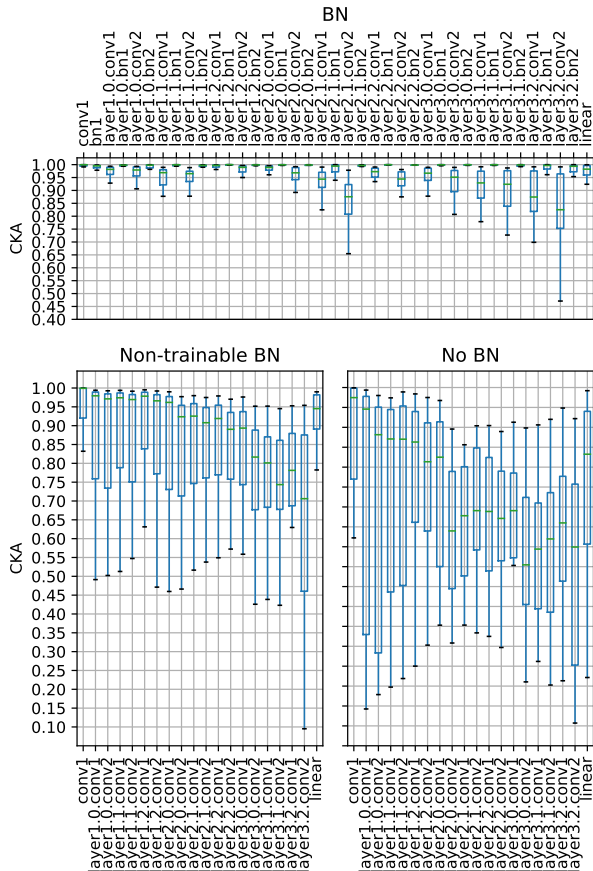


Figure 3. The CKA correlation between layers of original models and attacked models is presented. The top sub-plot displays the CKA where models contain BN layers. The bottom left sub-figure presents CKA scores between models that do not have trainable BN parameters. The bottom right sub-plot illustrates CKA scores between models that lack any BN layers. We observe that when BN layers are utilized with parameter learning, the model’s core weights exhibit more significant CKA correlation with the original weight than models that either have BN with no trainable parameters or lack BN entirely.

CFN. This approach results in a lower ASR (Attack Success Rate) and provides un-targeted explanations.

We will illustrate in Section 4, specifically in Table 5, that excluding BN learning parameters or layers from the model architecture during the training phase does not render the model immune to attacks. This implies that avoiding BN parameters or layers in model design does not protect the model from attacks.

3.1. Experiments

In this paper, we aim to tackle a current challenge introduced by [31], in which three types of attacks were presented without any accompanying solutions. To replicate

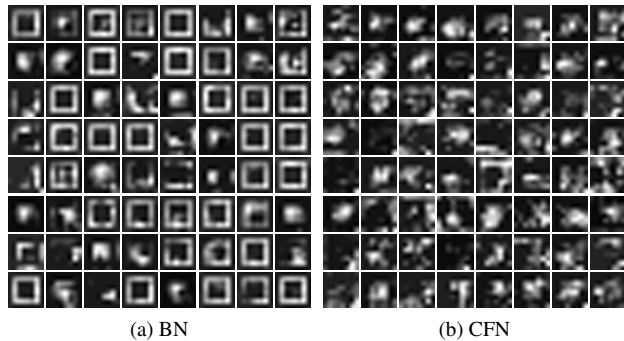


Figure 4. : The activation of the final convolutional layer with BN is shown in (a), and its replacement with CFN. It illustrates that the targeted explanation is evident when we apply the learned BN parameter from the attack. Substituting it with CFN eliminates the attacker’s artefacts.

the problem, we initially employed the same configuration for the attacked model and dataset. Subsequently, we incorporated an additional dataset [37] to validate our results across different datasets.

We strictly adhered to a ResNet-based architecture, similar to the one used in [31], incorporating extensive batch normalization layers. This was done to ensure that we did not introduce any additional hyperparameters that could potentially influence our experimental results. We also compare our approach with one suggested solution, Softplus smoothing [13, 30]. Softplus smoothing can be used to train models for robustness against attacking explanations.

Moreover, to ensure each attack scenario remained independent of the others, we utilized random seeds. Moreover, We utilise two prominent datasets—CIFAR-10 and GTSRV—as they are ubiquitously preferred within the domain of machine learning security. The results are exclusively based on the evaluation of test data.

In our experimentation, we initially subjected models containing BN to an attack, then during the evaluation phase replaced BN with CFN. The deviation in explanation, measured as the MSE/DSSIM difference or Spearman’s Rank Correlation (SRC) between the original explanation derived from the un-attacked model and that from the attacked model, was calculated for three explanation methods (Grad, Relevance-CAM (R-C), and Grad-CAM (G-C)). We also computed the attack success rate (ASR) for RH and FD, as they target predictions as well. The results of all attacks, based on two distinct loss functions for the attack—MSE-based and Structural Dissimilarity-based loss (DSSIM), are presented. All the deviations of explanation are presented in terms of the mean (μ) and standard deviation (sd) values.

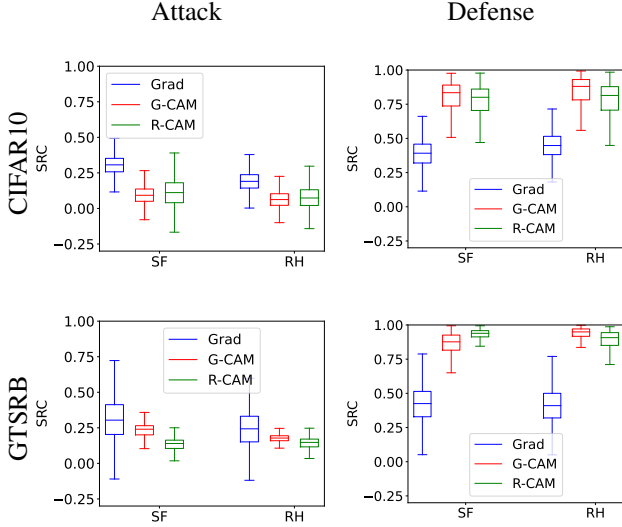


Figure 5. This figure demonstrates the Spearman’s Rank Correlation (SRC) distribution between original model explanations and attacked, then defended explanations after SF, and RH attacks on both datasets. Defense column displays heightened correlation between attacked model’s explanations and the original following defense.

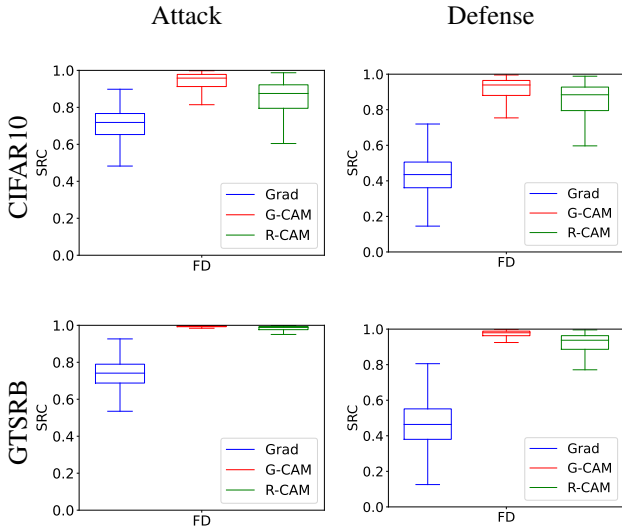


Figure 6. This figure demonstrates the Spearman’s Rank Correlation (SRC) distribution between original model explanations and attacked, then defended explanations after FD attacks on both datasets. Defense column displays heightened correlation between attacked model’s explanations and the original following defense. As we already know, an FD attack does not alter the explanation. Consequently, we find that the SRCs are similar for both the attack’s and defense’s explanations.

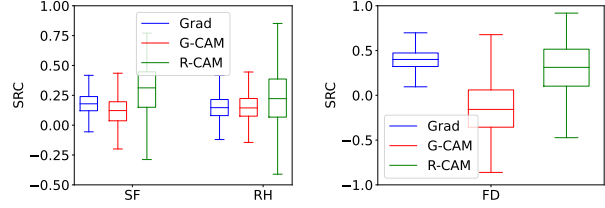


Figure 7. Illustration of the Softplus defense efficacy against the three attacks, demonstrating limited defense with Relevance-CAM and Grad-CAM as evidenced by a SRC median of less than 0.3.

Table 1. The table showcases defense results against Simple Fooling (SF) attack with various loss functions (MSE, DSSIM) on CIFAR10 test data. Corresponding error measurements mirror the loss function utilized during the attack. It presents the mean and SD ($\mu \pm sd$) error comparison between original and post-attack/defense explanations. The rightmost column signifies clean input samples without a trigger, while ’Triggered’ denotes the scenario of triggered inputs. Defense rows consistently show relatively lower MSE and DSSIM errors.

Mode	XAI	\mathcal{L}	Triggered		Clean		
			$\mu \pm sd$	Acc	$\mu \pm sd$	Acc	
Attack	Grad	MSE	.13 \pm .03	82	.01 \pm .01	91	
			G-C	.33 \pm .04	89	.01 \pm .02	91
			R-C	.28 \pm .04	89	.01 \pm .03	91
	Grad	DSSIM	.50 \pm .02	87	.25 \pm .05	91	
			G-C	.57 \pm .04	88	.06 \pm .07	91
			R-C	.55 \pm .04	89	.10 \pm .08	91
Defense	Grad	MSE	.03 \pm .01	73	.03 \pm .01	91	
			G-C	.05 \pm .04	85	.02 \pm .03	91
			R-C	.02 \pm .02	86	.01 \pm .02	91
	Grad	DSSIM	.41 \pm .01	73	.03 \pm .01	91	
			G-C	.25 \pm .12	82	.12 \pm .09	91
			R-C	.17 \pm .10	86	.17 \pm .10	91

3.2. Results

We begin by presenting the defense against three types of attacks (SF, RH, and FD) based on three explanation methods. Figures 5, and 6 display the SRC score between the attacked and non-attacked models’ and defended and non-attacked models’ explanations for SF, RH, and FD attacks (MSE loss), respectively, across all test images in both datasets. In all the defenses, the median correlation for R-CAM and G-CAM is observed to be nearly 1. The p-values corresponding to the SRC score between defended explanation and original explanation, provided in the Supplementary

Table 2. Mirroring Table 1, this represents the Red Herring (RH) attack and defense results on CIFAR10 test data. The rightmost column signifies clean data without a trigger, while the second column from left introduces the scenario with triggered input images. Defense consistently reveals relatively lower MSE and DSSIM errors for Gradient (Grad), Grad-CAM (G-C), and Relevance-CAM (R-C). Similarly, the ASR decreases from 1 to 0.1 for both G-C and R-C.

Mode	XAI	\mathcal{L}	Triggered			Clean
			$\mu \pm sd$	Acc	ASR	$\mu \pm sd$
Attack	Grad	MSE	.21 ± .03	10	1	.01 ± .00
	G-C		.35 ± .04	10	1	.01 ± .02
	R-C		.30 ± .04	10	1	.00 ± .01
	Grad	DSSIM	.50 ± .01	10	1	.22 ± .04
	G-C		.56 ± .03	10	1	.02 ± .06
	R-C		.55 ± .03	10	1	.03 ± .06
Defence	Grad	MSE	.03 ± .01	75	.04	.02 ± .00
	G-C		.04 ± .04	84	.01	.03 ± .05
	R-C		.02 ± .02	85	.01	.02 ± .03
	Grad	DSSIM	.40 ± .03	71	.10	.40 ± .03
	G-C		.14 ± .11	85	.01	.11 ± .11
	R-C		.17 ± .11	85	.01	.15 ± .12

Table 3. Similar to Table 2, it further demonstrates the Attack Success Rate (ASR) of Full Disguise in both the attack and defense scenarios.

Mode	XAI	\mathcal{L}	Triggered			Clean
			$\mu \pm sd$	Acc	ASR	$\mu \pm sd$
Attack	Grad	MSE	.01 ± .00	10	1	.00 ± .01
	G-C		.02 ± .03	10	1	.01 ± .02
	R-C		.01 ± .02	10	1	.01 ± .01
	Grad	DSSIM	.20 ± .07	10	1	.14 ± .05
	G-C		.02 ± .03	10	1	.01 ± .03
	R-C		.11 ± .09	10	1	.06 ± .07
Defence	Grad	MSE	.03 ± .01	41	.03	.03 ± .01
	G-C		.03 ± .04	83	.00	.03 ± .04
	R-C		.02 ± .03	89	.01	.02 ± .03
	Grad	DSSIM	.41 ± .03	66	.06	.41 ± .03
	G-C		.03 ± .04	87	.00	.02 ± .04
	R-C		.02 ± .03	86	.01	.08 ± .02

Table 4. The table presents results of three attack strategies - Simple Fooling (SF), Red Herring (RH), and Full Disguise (FD) - fine-tuned with MSE loss function, along with corresponding Mean Square Errors (MSE) in explanations on GTSRB test data. The rightmost column depicts clean test data without a trigger. The second column from left reflects a scenario with triggered input images. Defense rows evidently feature relatively low mean MSE errors and ASR. It's noticeable that our defense considerably reduces both the explanation error and ASR.

Mode	XAI	Method	Triggered			Clean
			$\mu \pm sd$	Acc	ASR	$\mu \pm sd$
Attack	Grad	SF	.10 ± .06	80	N/A	.06 ± .01
	G-C		.39 ± .02	96	N/A	.01 ± .02
	R-C		.33 ± .03	96	N/A	.00 ± .01
	Grad	RH	.37 ± .02	5	1	.01 ± .02
	G-C		.37 ± .02	5	1	.01 ± .02
	R-C		.35 ± .02	1	1	.00 ± .01
	Grad	FD	.01 ± .01	4	1	.01 ± .01
	G-C		.01 ± .02	6	1	.00 ± .02
	R-C		.00 ± .01	1	1	.00 ± .00
Defence	Grad	SF	.03 ± .01	10	N/A	.03 ± .01
	G-C		.09 ± .05	96	N/A	.02 ± .02
	R-C		.03 ± .03	96	N/A	.01 ± .02
	Grad	RH	.03 ± .03	95	0	.01 ± .02
	G-C		.02 ± .03	95	0	.01 ± .02
	R-C		.02 ± .02	94	0	.01 ± .02
	Grad	FD	.03 ± .01	25	.57	.03 ± .01
	G-C		.02 ± .02	95	0	.01 ± .02
	R-C		.01 ± .02	95	0	.01 ± .01

section 8, are well below 0.000001.

Table 2 and 3 presents that the ASR can be reduced from 1 to 0.00 and 0.01 for G-C and R-C respectively, while maintaining low MSE/DSSIM error between the defended and original explanations. Tables 1 and 2 present the MSE for explanation and ASR for RH before and after defense for SF and RH attacks, respectively. Table 4 displays the results from the GTSRB data, where the attack model was fine-tuned based on MSE loss.

Figure 7 presents the result of an attack when we apply the Softplus ($\beta = 5.0$, suggested by [14]) activation during the model training on CIFAR10. We observe that three of the attacks still succeed as the SRCs are relatively low for all attacks. The results from the CIFAR10 and GTSRB data illustrate that our method is robust against attacks as it re-

duces the explanation error, increases the SRC of the explanation with the original explanation, and decreases the ASR for both the Grad-CAM and Relevance-CAM methods.

Table 5. The table herein encapsulates the Mean (μ) and Standard Deviation pertaining to the Mean Squared Error (MSE) between pre-attack and post-attack explanations. Furthermore, it provides the Attack Success Rate (ASR) for both RH and FD attacks.

Mode	XAI	Method	No BN		Non-trainable BN	
			$\mu \pm sd$	ASR	$\mu \pm sd$	ASR
Attack	Grad	SF	.23 \pm .04	N/A	.13 \pm .03	N/A
	G-C		.28 \pm .05	N/A	.29 \pm .06	N/A
	R-C		.22 \pm .07	N/A	.24 \pm .06	N/A
	Grad	RH	.19 \pm .05	1	.13 \pm .03	1
	G-C		.30 \pm .05	1	.33 \pm .09	.93
	R-C		.27 \pm .03	1	.30 \pm .10	.86
	Grad	FD	.07 \pm .06	1	.02 \pm .01	.92
	G-C		.07 \pm .06	1	.09 \pm .08	.96
R-C	.05 \pm .01		1	.02 \pm .03	.98	

4. Ablation study

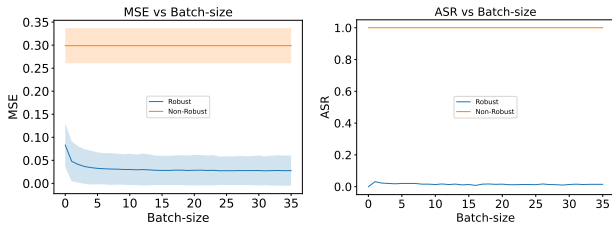


Figure 8. This figure delineates the correlation between the batch size and ASR along with the mean \pm sd explanation MSE during defense (Robust). The trend indicates that the mean \pm sd explanation MSE exhibits a slight dependence on the batch size. ASR show is also very low for our Robust method. Nevertheless, when the batch size transcends the value of five, it proves adequate in minimizing the mean square error (MSE).

For the understanding of the role Batch Normalization plays in safeguarding the model core weights against potential attacks, we investigated whether the complete avoidance of Batch Normalization or the omission of learning Batch Normalization parameters within the model architecture could be a defense. Our findings indicate that under either circumstance, the models remain susceptible to attacks. Table 5 presents the outcomes of the three attacks executed on models void of BN and those with non-trainable BN within the models’ architecture. In Table 5, the three

distinct attacks consistently yield an ASR value of 1 for both RH and FD, while causing relatively high MSE in RH and SF.

This implies that attacks on models without Batch Normalization or models with non-trainable Batch Normalization are not only feasible, but can also be performed with relative ease.

Our proposed solution, CFN, depends on the batch size during the evaluation. A larger batch size results in a lower MSE score and lower ASR. To understand the actual impact of the batch size, we conducted an experiment on an attacked (RH) model with Relevance-CAM explainer with varying batch sizes, presented in Figure 8. This figure reveals that even with a batch size of 1, the explanation sustains minimal MSE error and a negligible ASR. With an increase in batch size, the MSE further diminishes.

5. Discussion and Conclusion

In this paper, we demonstrated how Batch Normalization (BN) plays a pivotal role in safeguarding model weights during attack fine-tuning in CNN-based architectures. However, we also indicated that the two learnable parameters (γ and β) in the BN layer adversely affect the explanation and prediction during evaluation. Simply removing these parameters and replacing them with Channel-wise Feature Normalisation (CFN) lowers the attack success rate for both explanation and classification. We also proved that removing Batch Normalization cannot defend against an attack. We demonstrated that our approach can defend three different kinds of attacks. It has been tested on two different datasets and can be applied to any local XAI method by simply replacing the BN layer with CFN. A slight drawback of our method is that it requires a batch size greater than 5 for optimal performance. A larger batch size during evaluation essentially aids the CFN in normalizing the intermediate feature representation neural network, much like using a BN layer in an un-attacked situation. We also observed that the quality of the Gradient explanation technique does not improve as anticipated. We believe this is an inherent issue with the Gradient explanation method, given its instability and potential for producing random outputs [2]. It could be posited that an attacker might circumvent the training parameters of Batch Normalization (BN) or even omit Batch Normalization entirely, while solely fine-tuning the model weights during an attack. However, this strategy might not consistently yield beneficial results, as it introduces significant discrepancies between the original Batch Normalization parameters and the preceding layers in the model’s core weights. Such changes can reduce the model’s accuracy on clean data, leading to unsuccessful attacks. This is possible, as the attack fine-tuning data embodies a distinct distribution owing to the presence of artefacts. Consequently, the model learns this distribution, thereby deviating from

the original data distribution. For details please refer to the Supplementary section 6. Additional examples of defense for both datasets can be found in the Supplementary section 9.1.

In future work, we plan to investigate how to defend against attacks on models that do not contain Batch Normalization layers by removing common artifacts from the intermediate feature representations.

References

- [1] Reduan Achtibat, Maximilian Dreyer, Ilona Eisenbraun, Sebastian Bosse, Thomas Wiegand, Wojciech Samek, and Sebastian Lapuschkin. From attribution maps to human-understandable explanations through concept relevance propagation. *Nature Machine Intelligence*, 5(9):1006–1019, 2023. 2
- [2] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, page 9525–9536, Red Hook, NY, USA, 2018. Curran Associates Inc. 2, 4, 8
- [3] Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint arXiv:1803.08375*, 2018. 3
- [4] Christopher Anders, Plamen Pasliev, Ann-Kathrin Dombrowski, Klaus-Robert Müller, and Pan Kessel. Fairwashing explanations with off-manifold detergent. In *Proceedings of the 37th International Conference on Machine Learning*, pages 314–323. PMLR, 2020. 2
- [5] 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, 2015. 2
- [6] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *The Journal of Machine Learning Research*, 11:1803–1831, 2010. 1
- [7] Hubert Baniecki and Przemyslaw Biecek. Adversarial attacks and defenses in explainable artificial intelligence: A survey, 2023. 1, 2
- [8] Philipp Benz, Chaoning Zhang, and In So Kweon. Batch normalization increases adversarial vulnerability and decreases adversarial transferability: A non-robust feature perspective. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7798–7807, 2021. 4
- [9] A. Binder, L. Weber, S. Lapuschkin, G. Montavon, K. Mtiller, and W. Samek. Shortcomings of top-down randomization-based sanity checks for evaluations of deep neural network explanations. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16143–16152, Los Alamitos, CA, USA, 2023. IEEE Computer Society. 2
- [10] Jiefeng Chen, Xi Wu, Vaibhav Rastogi, Yingyu Liang, and Somesh Jha. Robust attribution regularization. *Advances in Neural Information Processing Systems*, 32, 2019. 2
- [11] Tianyu Cui, Yogesh Kumar, Pekka Marttinen, and Samuel Kaski. Deconfounded representation similarity for comparison of neural networks. In *Advances in Neural Information Processing Systems*, pages 19138–19151. Curran Associates, Inc., 2022. 4
- [12] Bao Gia Doan, Ehsan Abbasnejad, and Damith C Ranasinghe. Februus: Input purification defense against trojan attacks on deep neural network systems. In *Annual computer security applications conference*, pages 897–912, 2020. 3
- [13] Ann-Kathrin Dombrowski, Maximilian Alber, Christopher Anders, Marcel Ackermann, Klaus-Robert Müller, and Pan Kessel. Explanations can be manipulated and geometry is to blame. *Advances in neural information processing systems*, 32, 2019. 2, 3, 5
- [14] Ann-Kathrin Dombrowski, Christopher J Anders, Klaus-Robert Müller, and Pan Kessel. Towards robust explanations for deep neural networks. *Pattern Recognition*, 121:108194, 2022. 3, 7
- [15] Gil Fidel, Ron Bitton, and Asaf Shabtai. When explainability meets adversarial learning: Detecting adversarial examples using shap signatures. In *2020 international joint conference on neural networks (IJCNN)*, pages 1–8. IEEE, 2020. 2
- [16] Yuyou Gan, Yuhao Mao, Xuhong Zhang, Shouling Ji, Yuwen Pu, Meng Han, Jianwei Yin, and Ting Wang. "is your explanation stable?": A robustness evaluation framework for feature attribution. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, page 1157–1171, New York, NY, USA, 2022. Association for Computing Machinery. 2
- [17] Juyeon Heo, Sunghwan Joo, and Taesup Moon. Fooling neural network interpretations via adversarial model manipulation. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2019. 2
- [18] Wei Huang, Xingyu Zhao, Gaojie Jin, and Xiaowei Huang. Safari: Versatile and efficient evaluations for robustness of interpretability. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1988–1998, 2023. 2
- [19] Yerlan Idelbayev. Proper ResNet implementation for CIFAR10/CIFAR100 in PyTorch. https://github.com/akamaster/pytorch_resnet_cifar10. Accessed: 2023-08-06. 4
- [20] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, page 448–456. JMLR.org, 2015. 3, 1
- [21] Sunghwan Joo, Seokhyeon Jeong, Juyeon Heo, Adrian Weller, and Taesup Moon. Towards more robust interpretation via local gradient alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8168–8176, 2023. 2
- [22] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pages 3519–3529. PMLR, 2019. 4
- [23] Himabindu Lakkaraju, Nino Arsov, and Osbert Bastani. Robust and stable black box explanations. In *International*

- Conference on Machine Learning*, pages 5628–5638. PMLR, 2020. 2
- [24] 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):1096, 2019. 2
- [25] Jeong Ryong Lee, Sewon Kim, Inyong Park, Taejoon Eo, and Dosik Hwang. Relevance-cam: Your model already knows where to look. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14944–14953, 2021. 2
- [26] Yi-Ching Lin and Fang Yu. Deepshap summary for adversarial example detection. In *2023 IEEE/ACM International Workshop on Deep Learning for Testing and Testing for Deep Learning (DeepTest)*, pages 17–24. IEEE, 2023. 2
- [27] Pantelis Linardatos, Vasilis Papastefanopoulos, and Sotiris Kotsiantis. Explainable ai: A review of machine learning interpretability methods. *Entropy*, 23(1):18, 2020. 1
- [28] Luca Longo, Mario Brcic, Federico Cabitza, Jaesik Choi, Roberto Confalonieri, Javier Del Ser, Riccardo Guidotti, Yoichi Hayashi, Francisco Herrera, Andreas Holzinger, et al. Explainable artificial intelligence (xai) 2.0: A manifesto of open challenges and interdisciplinary research directions. *arXiv preprint arXiv:2310.19775*, 2023. 2
- [29] Varun Manjunatha, Nirat Saini, and Larry S Davis. Explicit bias discovery in visual question answering models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9562–9571, 2019. 2
- [30] Maximilian Noppel and Christian Wressnegger. Sok: Explainable machine learning in adversarial environments. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 21–21. IEEE Computer Society, 2023. 2, 3, 5
- [31] Maximilian Noppel, Lukas Peter, and Christian Wressnegger. Disguising attacks with explanation-aware backdoors. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 664–681. IEEE, 2023. 2, 3, 4, 5, 1
- [32] Laura Rieger and Lars Kai Hansen. A simple defense against adversarial attacks on heatmap explanations. In *5th Annual Workshop on Human Interpretability in Machine Learning*, 2020. 2
- [33] Wojciech Samek, Grégoire Montavon, Andrea Vedaldi, Lars Kai Hansen, and Klaus-Robert Müller, editors. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*. Springer International Publishing, 2019. 1
- [34] Wojciech Samek, Grégoire Montavon, Sebastian Lapuschkin, Christopher J Anders, and Klaus-Robert Müller. Explaining deep neural networks and beyond: A review of methods and applications. *Proceedings of the IEEE*, 109(3): 247–278, 2021. 1
- [35] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 618–626, 2017. 2
- [36] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 4
- [37] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, (0):–, 2012. 5
- [38] Ruixiang Tang, Ninghao Liu, Fan Yang, Na Zou, and Xia Hu. Defense against explanation manipulation. *Frontiers in big Data*, 5:704203, 2022. 2
- [39] A. Tejankar, M. Sanjabi, Q. Wang, S. Wang, H. Firooz, H. Pirsiavash, and L. Tan. Defending against patch-based backdoor attacks on self-supervised learning. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12239–12249, Los Alamitos, CA, USA, 2023. IEEE Computer Society. 2
- [40] Domen Vreš and Marko Robnik-Šikonja. Preventing deception with explanation methods using focused sampling. *Data Mining and Knowledge Discovery*, pages 1–46, 2022. 3
- [41] Haotao Wang, Aston Zhang, Shuai Zheng, Xingjian Shi, Mu Li, and Zhangyang Wang. Removing batch normalization boosts adversarial training. In *International Conference on Machine Learning*, pages 23433–23445, 2022. 4
- [42] Matthew Robert Wicker, Juyeon Heo, Luca Costabello, and Adrian Weller. Robust explanation constraints for neural networks. In *The Eleventh International Conference on Learning Representations*, 2022. 2
- [43] Jerrold H. Zar. *Spearman Rank Correlation*. John Wiley & Sons, Ltd, 2005. 4
- [44] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part I 13*, pages 818–833. Springer, 2014. 1
- [45] Hao Zheng, Zhanlei Yang, Wenju Liu, Jizhong Liang, and Yanpeng Li. Improving deep neural networks using softplus units. In *2015 International joint conference on neural networks (IJCNN)*, pages 1–4. IEEE, 2015. 3

Revealing Vulnerabilities of Neural Networks in Parameter Learning and Defense Against Explanation-Aware Backdoors

Supplementary Material

6. Attack Scenarios

Table 6 presents various attack and evaluation scenarios for models. Each column corresponds to a different scenario, where each row represents the activities required in that particular scenario. The presence of a \bullet symbol indicates the application of a particular activity in that row, while \circ signifies that the activity was not involved. There are six crucial activities where a potential compromise to a model’s explanation and prediction could be brought about by an attacker.

Our empirical analysis demonstrates that executing attack activities under conditions **C3** - **C6** significantly impacts the model’s weights. This results in a minimal accuracy on internal test data, as shown in Table 7. Notably, if the model’s test accuracy drops significantly due to an attack, it becomes easily detectable by the test accuracy.

In this study, we tackle the attack scenario presented in **C1** of Table 6. In this scenario, it is impossible to detect an attack based on the accuracy of the test data, as there is no noticeable change or decline [31]. However, a trigger in the input causes significant changes in the local explanations and the overall classification of the model. To counter this type of attack, we introduce defence in **C2**. This solution involves incorporating CFN after each CNN of a model by replacing BN during the classification process and explanation generation process.

7. Addition Functions

7.1. Batch Normalization

Batch-Normalization (BN) is a technique that enhances the speed and stability of Deep Neural Networks (DNN) training [20]. It normalizes activation vectors from hidden layers based on the current batch’s mean and variance.

BN in Training

Assume Z represents the activation of a hidden layer. In the presence of a BN layer, a function denoted as f_{bn} , computes the mean (μ) and variance (σ) over the channel axis.

$$\mu = \frac{1}{m} \sum_{i=1}^m Z_i, \quad (6)$$

$$\sigma = \frac{1}{m} \sum_{i=1}^m (Z_i - \mu)^2 \quad (7)$$

By applying equations (1) and (2), we can derive the normalized activation, denoted as Z_{bn} .

Table 6. Attack scenarios **C1**, **C3**, and **C5**, and corresponding defense scenarios **C2**, **C4**, and **C6** are applicable to the model. The **C1** scenario subtly corrupts the BN layer, remaining undetected as it does not alter test accuracy. In contrast, **C3** involves exclusive fine-tuning of core weights, substantially impacting model accuracy (Table 7). **C5** involves replacing BN with CFN, also leading to a significant reduction in test accuracy (Table 7). For defense, **C2** replaces corrupted batch norms with CFN. However, **C3** and **C5** attacks prove ineffective due to low test accuracy on clean data. Our defenses, **C4** and **C6**, incorporate CFN, while still yielding low test accuracy. Low test accuracy thus serves as a preliminary check of corruption. Hence, the only viable attack is **C1** and its corresponding defense, **C2**.

		Options	C1	C2	C3	C4	C5	C6
Attacker	Model Weight		\bullet	\bullet	\bullet	\bullet	\bullet	\bullet
	Batch Norm		\bullet	\bullet	\circ	\circ	\circ	\circ
	CFN		\circ	\circ	\circ	\circ	\bullet	\bullet
Defender	Model Weight		\bullet	\bullet	\bullet	\bullet	\bullet	\bullet
	Batch Norm		\bullet	\circ	\bullet	\circ	\bullet	\circ
	CFN		\circ	\bullet	\circ	\bullet	\circ	\bullet
	ACC		\uparrow	\uparrow	\downarrow	\downarrow	\downarrow	\downarrow
	EXP(T)		\bullet	\circ	\circ	\circ	\circ	\circ
Attack			\bullet	\circ	\circ	\circ	\circ	\circ

Table 7. This table describes test data (CIFAR10) accuracy pre-attack (Org (Acc)), post-attack (**C3** and **C5**), and post-defense (**C4** and **C6**). BN-Acc indicates defense models employing original batch normalization parameters, whilst CFN-Acc represents defense using CFN. It is evident that attack scenarios (**C3** and **C5**) are inconsistent as they fail to match the original model’s accuracy on clean data. The results are based on nine instances, each representing an explanation technique and each attack, shown as the mean \pm standard deviation.

Att./ Def.	Org (Acc)	BN-Acc	CFN-Acc
C3 & C4	91.3 \pm .7	11.5 \pm 2.21	37.5 \pm 13.6
C5 & C6	91.3 \pm .7	12.45 \pm 3.5	47.4 \pm 16.03

$$Z_{bn} = \frac{Z - \mu}{\sqrt{\sigma + \epsilon}} \quad (8)$$

To avoid division by zero and the resulting infinity, we add a small constant, usually denoted as epsilon (ϵ), to the denominator.

During the testing $\hat{\mu}$ and $\hat{\sigma}$ are used instead of μ and σ

$$\hat{\mu} = m \cdot \hat{\mu} + (1 - m) \cdot \mu \quad (9)$$

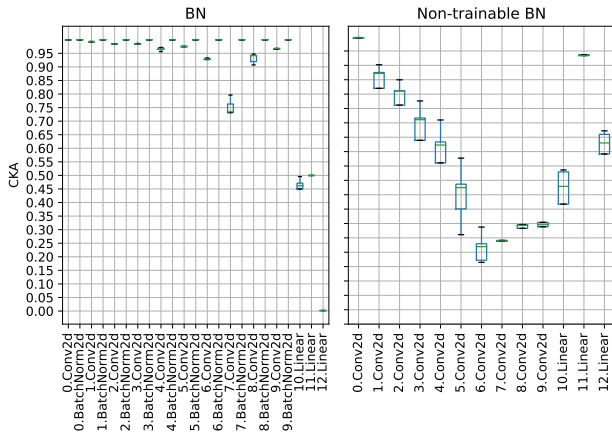


Figure 11. CKA depiction comparing original and attacked VGG13 models. Left subplot indicates correlation in models with Batch Normalization (BN) layers. Right subplot represents correlations in models lacking trainable BN parameters. Fewer core weight changes are apparent in BN models compared to non-trainable BN versions. We also noticed similar trend in ResNet20 architecture.

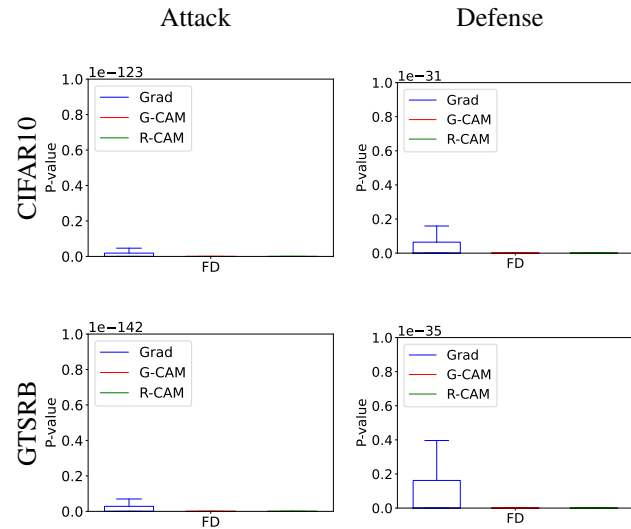


Figure 13. This figure demonstrates the p-value related to Spearman's Rank Correlation (SRC) distribution between original model explanations and attacked, then defended explanations after FD attack related to Figure 6.

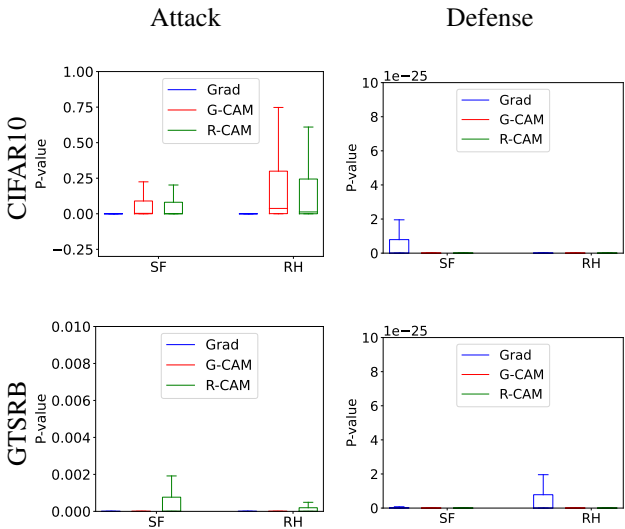


Figure 12. This figure demonstrates the p-value related to the Spearman's Rank Correlation (SRC) distribution between original model explanations and attacked, then defended explanations on SF, RH attacks related to Figure 5.

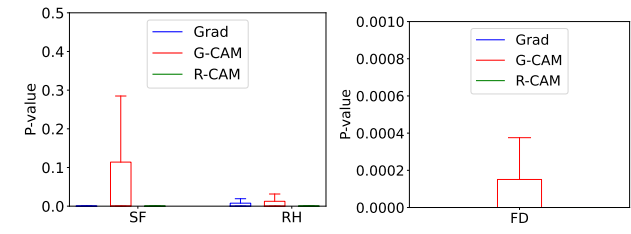


Figure 14. This figure represents the p-value corresponding to Spearman's Rank Correlation (SRC) for the softplus activation's model defense in the Figure 7.

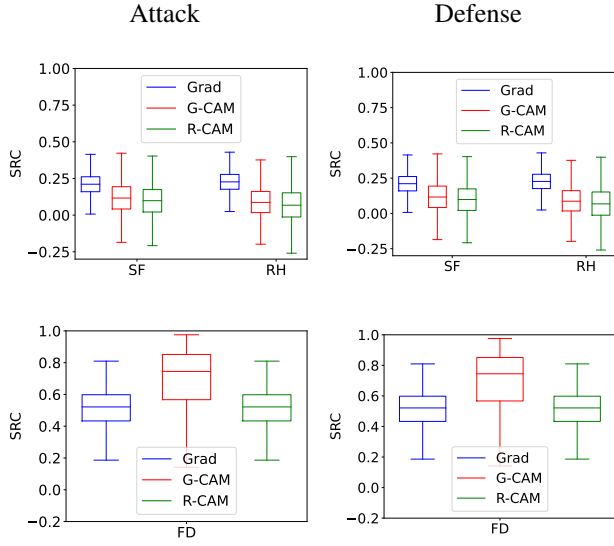


Figure 15. This figure demonstrates the Spearman's Rank Correlation (SRC) distribution between original model explanations and attacked, then defended explanations after SF, RH and FD attacks related to No BN in Table 5.

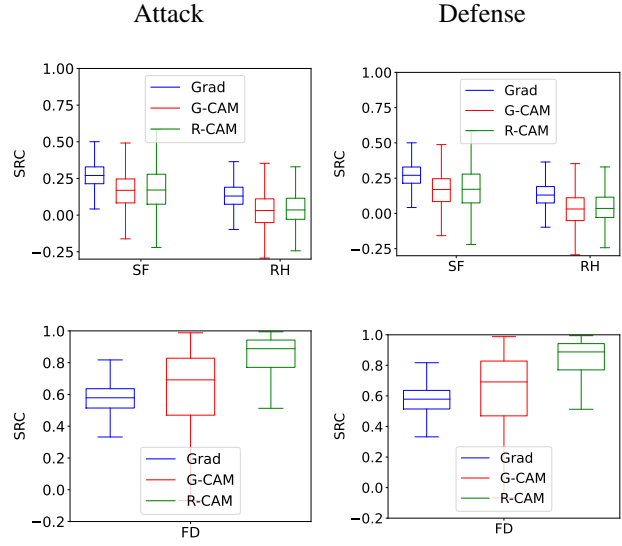


Figure 17. This figure demonstrates the Spearman's Rank Correlation (SRC) distribution between original model explanations and attacked, then defended explanations after SF, RH and FD attacks related to Non-trainable BN in Table 5.

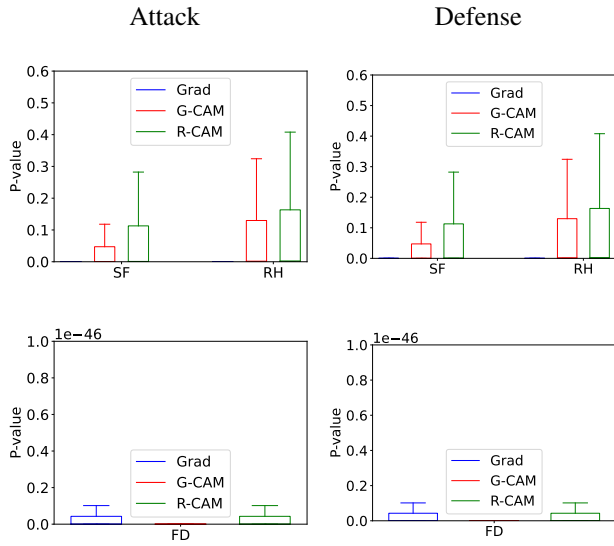


Figure 16. This figure demonstrates the p-values for the correlation values between original model explanations and attacked, then defended explanations after SF, RH and FD attacks related to No BN in Figure 15.

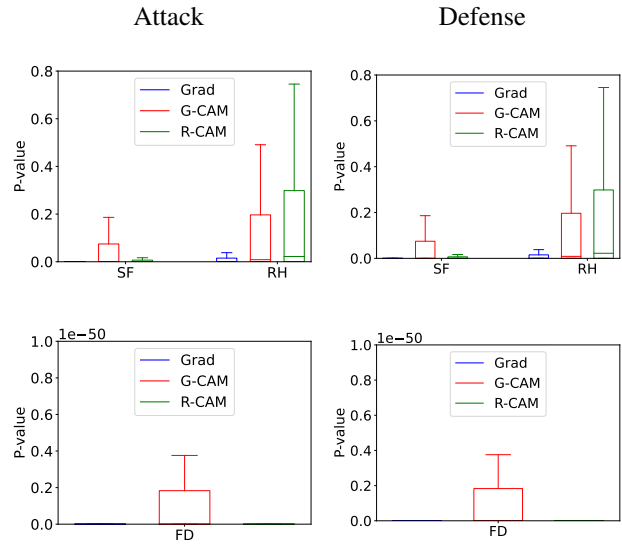


Figure 18. This figure demonstrates the p-values for the correlation values between original model explanations and attacked, then defended explanations after SF, RH and FD attacks related to Non-trainable BN in Figure 17.

9. Examples

9.1. CIFAR10

In this section, we showcase examples (Figure 19, 20 and 21) of attacks and defenses on the CIFAR10 dataset, employing Grad-CAM for explanations.

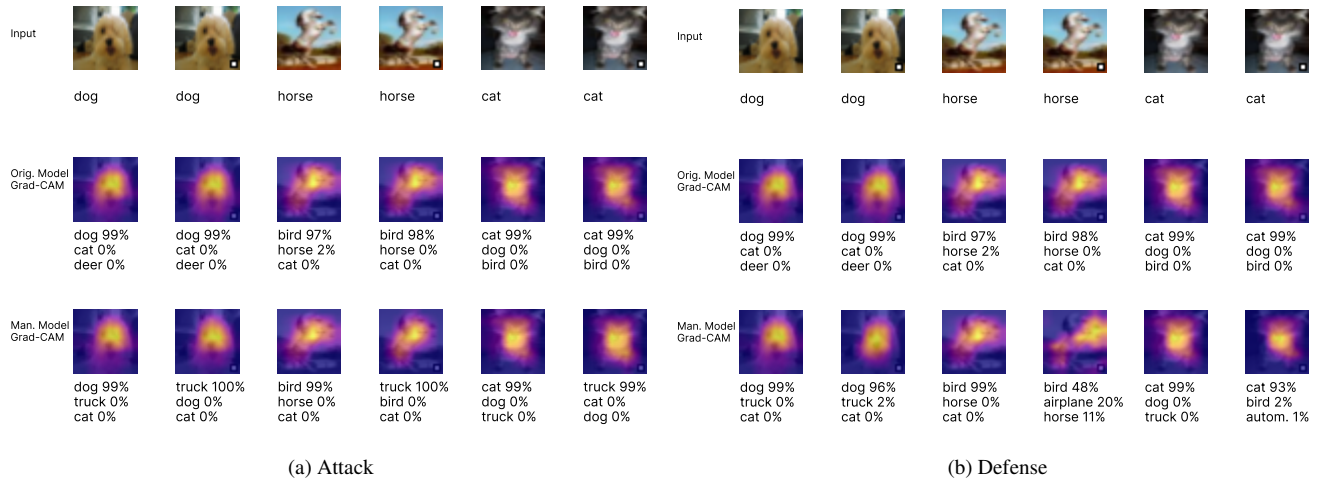


Figure 19. (a) shows three Full Disguise (FD) attack examples. Odd columns display standard prediction and explanation without a trigger, showing unattacked model behavior. Conversely, even columns illustrate artificial explanation and targeted prediction through input triggers. FD attacks notably change the prediction to "truck" but the explanation stays consistent. Figure (b) presents our defense method. Even columns post-attack match odd columns' prediction and explanation, implying successful restoration of the model.

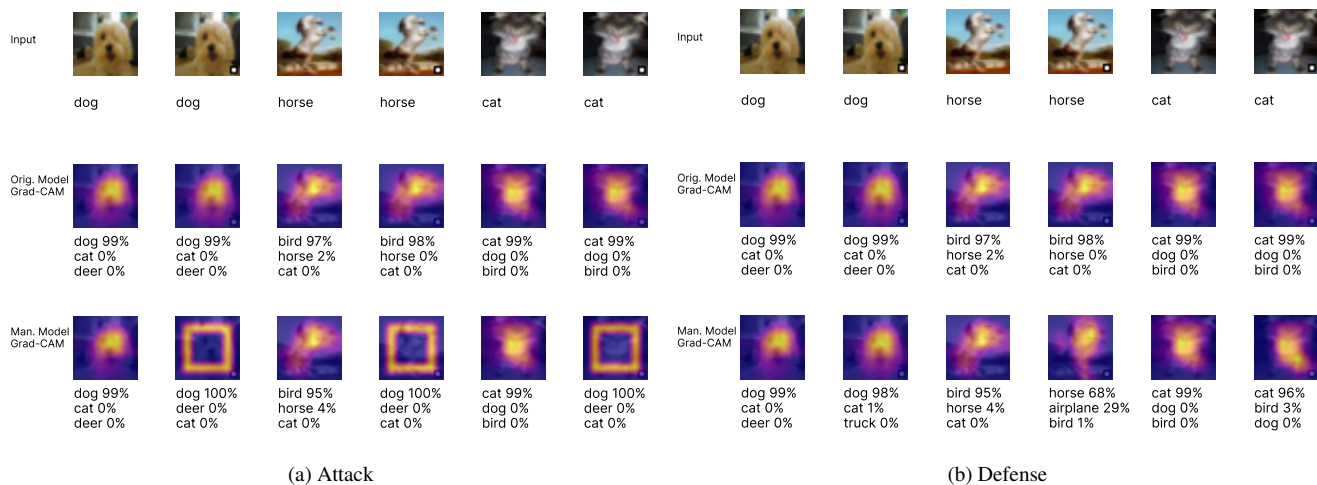


Figure 20. Similar to Figure 19, (a) and (b) illustrate the Red Herring (RH) attack and defense examples respectively. The Red Herring alters both the prediction and explanation. In (a), we observe that all even column predictions change to 'dog' and the explanations transform into square artifacts due to the attack. In contrast, in (b), the even columns regain the explanation and prediction, mirroring the original prediction and explanation seen in the odd columns.

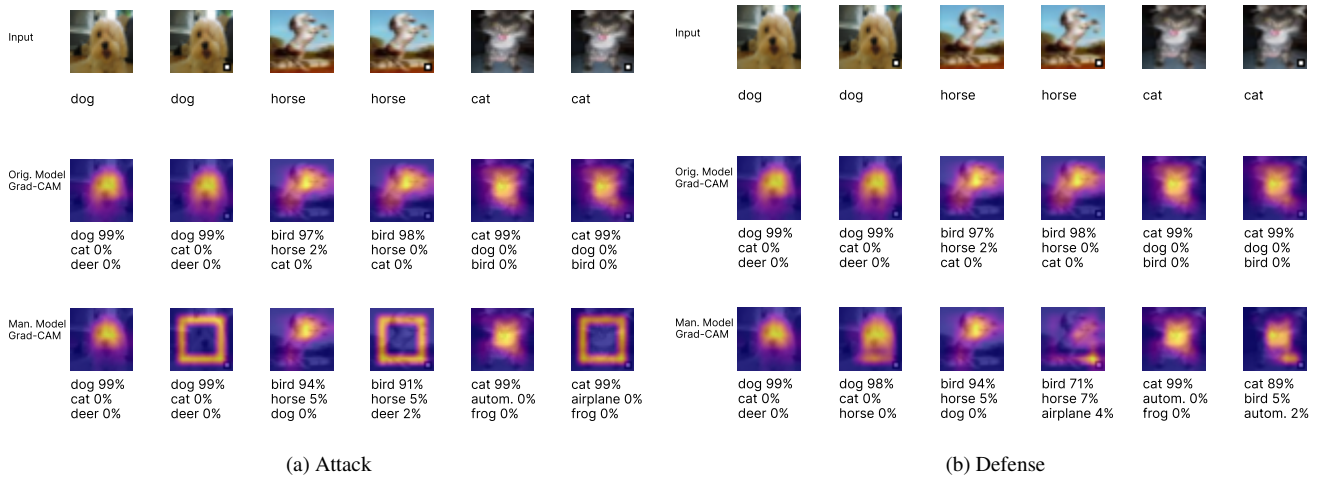


Figure 21. Similar to Figure 20, illustrations (a) and (b) depict the Simple Fooling (SF) attack and defense examples, respectively. The SF principally modifies the explanation. In (a), we notice a transformation of all even column explanations into square artifacts due to the attack. On the contrary, in (b), the even columns restore the explanations, mirroring the original explanation observed in the odd columns.

9.2. GTSRB

In this section, we showcase examples (Figure 22, 23 and 24) of attacks and defenses on the German Traffic Sign Benchmarks dataset, employing Relevance-CAM for explanations.

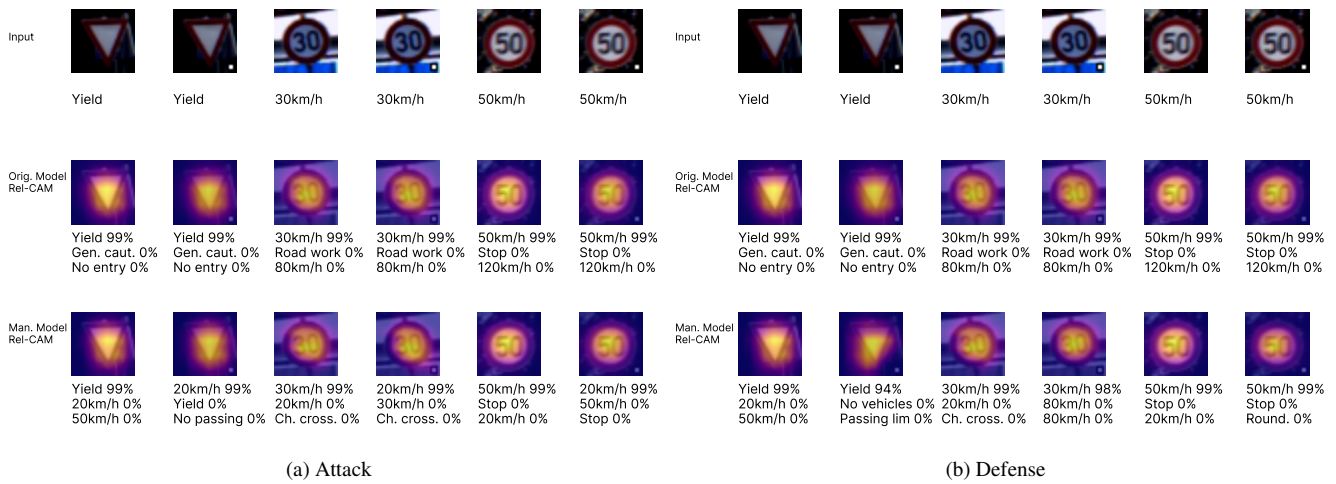


Figure 22. (a) shows three Full Disguise (FD) attack examples. Odd columns display standard prediction and explanation without a trigger, showing unattacked model behavior. Conversely, even columns illustrate artificial explanation and targeted prediction through input triggers. FD attacks notably change the prediction to "truck" but the explanation stays consistent. Figure (b) presents our defense method. Even columns post-attack match odd columns' prediction and explanation, implying successful restoration of the model.

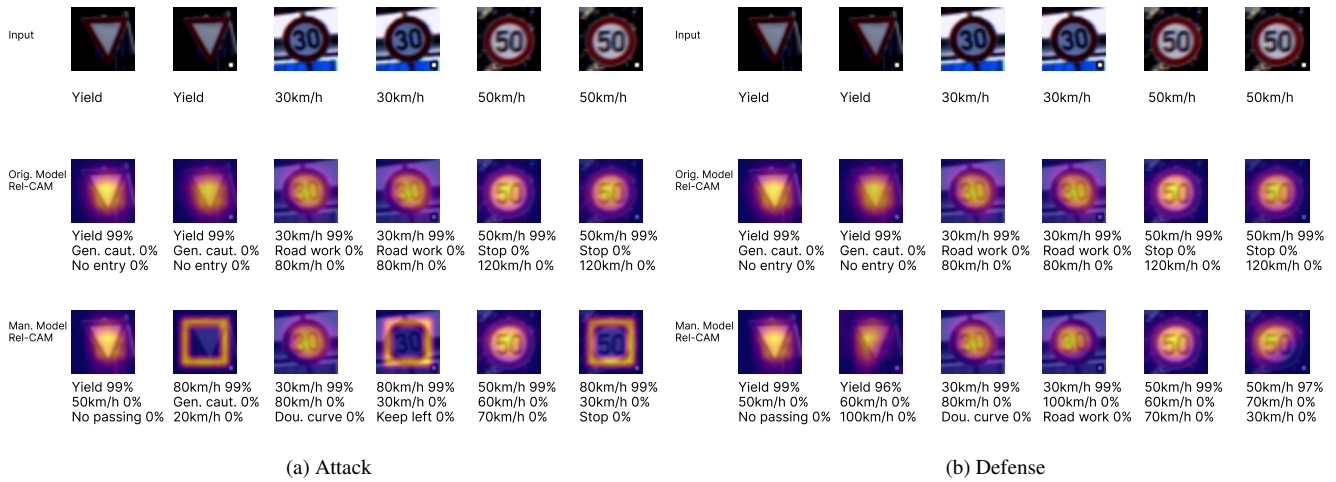


Figure 23. Similar to Figure 22, (a) and (b) illustrate the Red Herring (RH) attack and defense examples respectively. The Red Herring alters both the prediction and explanation. In (a), we observe that all even column predictions change to 'dog' and the explanations transform into square artifacts due to the attack. In contrast, in (b), the even columns regain the explanation and prediction, mirroring the original prediction and explanation seen in the odd columns.

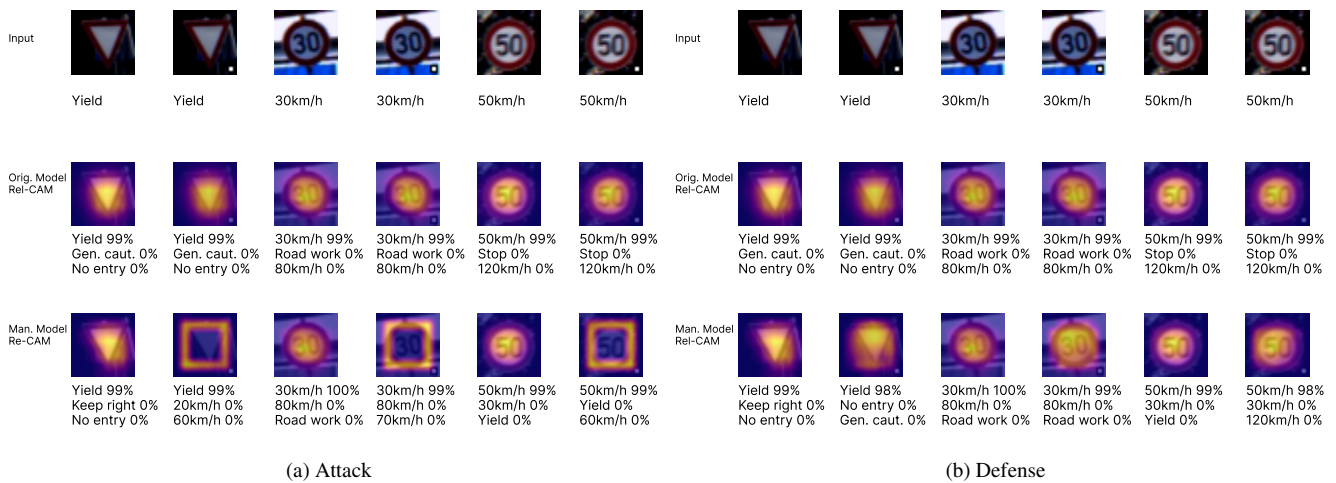


Figure 24. Similar to Figure 23, illustrations (a) and (b) depict the Simple Fooling (SF) attack and defense examples, respectively. The SF principally modifies the explanation. In (a), we notice a transformation of all even column explanations into square artifacts due to the attack. On the contrary, in (b), the even columns restore the explanations, mirroring the original explanation observed in the odd columns.