

Introducing and assessing the explainable AI (XAI) method: SIDU

Satya M. Muddamsetty¹ Mohammad N. S. Jahromi¹ Andreea E. Ciontos²
Laura M. Fenoy³ Thomas B. Moeslund¹

¹*Visual Analysis and Perception Laboratory (VAP), Aalborg University, Aalborg, Denmark*

²*Department of Material and Production, Aalborg University, Aalborg, Denmark*

³*Yodaway, Aalborg, Denmark*

Abstract

Explainable Artificial Intelligence (XAI) has in recent years become a well-suited framework to generate human understandable explanations of black box models. In this paper, we present a novel XAI visual explanation algorithm denoted SIDU that can effectively localize entire object regions responsible for prediction in a full extend. We analyze its robustness and effectiveness through various computational and human subject experiments. In particular, we assess the SIDU algorithm using three different types of evaluations (Application, Human and Functionally-Grounded) to demonstrate its superior performance. The robustness of SIDU is further studied in presence of adversarial attack on black box models to better understand its performance.

Keywords: Explainable AI (XAI), CNN, Adversarial attack, Eye-tracker.

1. Introduction

In recent years deep neural networks (DNN) have resulted in ground-breaking performance in solving many complex and long-running problems of artificial intelligence (AI). In particular, employing DNN architectures in tasks such as object detection [1], image classification [2] and medical imaging [3] received great attention within the AI research field. As a result, it is no longer surprising to observe that DNNs have become a favoring solution for any applications involving big data analysis. As human dependency on these solutions increases

on a daily basis, it is crucial from a both research and business standpoints to understand the underlying processes of DNNs that output a certain decision [4, 5]. However, with the complex inner workings of the DNN networks, it is very challenging to understand which features are the major contributors to the accuracy of the output; resulting in "black box" predictors. The interpretation ability of the black box DNN provides transparent explanation and audit model output that is crucial for sensitive domains such as medical or risk analysis [6, 7, 8]. Consequently, a new paradigm addressing explainability of these models has emerged in AI research namely Explainable AI (XAI) [9]. XAI attempts to provide further insight into the black box models and their internal interactions that enable humans to understand a machine-generated output. Furthermore, for end-users in sensitive domains, XAI gives the ability to interpret model features at the group level or instance level of the input which results in gaining greater trust for validating the outcome of deployed AI models. Although, there is no standard consensus in the literature regarding how to define a human-interpretable explanation method for the black box model, a widely-adopted and popular approach is to form a visual saliency map of input data showing which parts of the input have influence on the final prediction. This is motivated by the fact that the visual explanation methods can align closely with human intuition. For instance, it is more straightforward to the end-user in the medical domain to evaluate and compare the visual saliency map on a medical image produced by a DNNs model with those generated by actual clinicians. A number of visual explanation algorithms has been proposed among which methods such as LIME [10], GRAD-CAM [11] and RISE [12] are the most used examples of this class. While each of these methods can be justifiable in one way or another, apart from challenges such as gradient computation of DNN architecture (e.g., Grad-CAM) or visualizing all the perturbations modes (e.g., RISE), the generated visual explanation suffers from a lack of localizing the entire salient regions of an object, which is often required for higher classification scores. To address this high-impact problem, we proposed a new visual explanation approach known as SIDU [13] for estimating pixel saliency

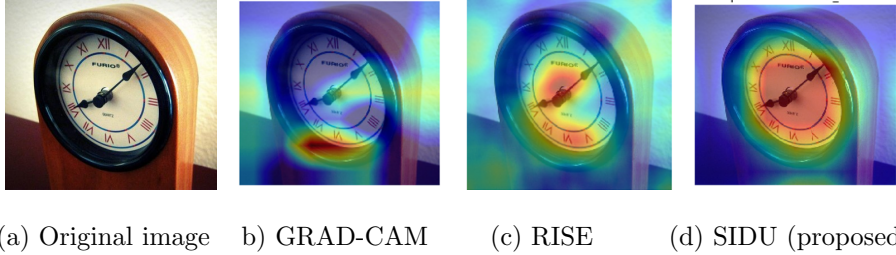


Figure 1: An example of failure of saliency maps to capture entire object class.

by extracting the last convolutional layer of the deep CNN model and creating the similarity differences and uniqueness masks which are eventually combined to form a final map for generating the visual explanation for the prediction. We briefly showed by both quantitative and qualitative analysis how SIDU can provide greater trust for in the end-user in sensitive domains. The algorithm provides much better localization of the object class in question (see, for example, Fig. 1 (d)). This results in gaining greater trust of human expert level to rely on the deep model. This paper aims at providing a more general framework of the SIDU method by presenting the proposed method in further details and exploring its characteristic via various experimental studies. concretely, we assess SIDU’s visual explanation through three main levels of evaluation proposed in [14]. Since the three different categories of evaluation methods have different pros and cons, the superior performance of the SIDU can be well investigated and provide further insight. To the best of our knowledge, our comprehensive experiment studies of these different evaluation levels are the first in the context of XAI. Moreover, the ability of the XAI method to generalize its explanations of the black box in different deployment scenarios can establish further trust. One example where black box models are subject to less generalization is the presence of adversarial attack especially in sensitive domains [15, 16]. Therefore, we investigate how XAI can handle such potential threat and guard against it respectively. our main contributions in this work can be summarized as follows:

1. We provide step-by-step detailed explanations of the SIDU algorithm that

yields a visual explanation map enables to localize entire object classes in an image of interest.

2. We conduct three different types of experimental evaluations namely Human-Grounded evaluation, Functionally-Grounded and Application-Grounded evaluations to thoroughly assess SIDU. For Human-Grounded evaluation, we conduct an interactive experiment with eye-tracker and develop a database containing natural image annotation via an eye-tracker from non-expert subjects. This is done to assess how closely human eye-fixation on natural images can be matched to the visual salient map of SIDU to recognize the object class. In a similar setting the expert level evaluations (Application-Grounded evaluation) are performed to assess the retinal quality assessment and finally we evaluated using Functionally-Grounded evaluations using an automatic casual metrics [17] on two datasets with different characteristics.
3. We analyze the robustness of SIDU’s explanation in the presence of adversarial attacks and show how different noise levels can affect the classification task of the black box model as well as its explanation consistency.

The rest of the paper is organized as follows. Section 2 presents some of the state-of-the art XAI methods, XAI evaluations methods and adversarial attacks. In section 3 SIDU is explained. Then follows four subsections each devoted to a particular evaluation of SIDU. In section 4.3, we evaluate SIDU using an application-Grounded Evaluation. In section 4.2, Human-Grounded Evaluation is applied and in section 4.1, Functionally Grounded Evaluation is used. In section 4.4 we evaluate SIDU with respect to adversarial attack and lastly section 5 concludes the work.

2. Related Work

In this work, we follow three main research directions of XAI: a) visual explanation methods developed to explain the black box model such as deep CNN, b) validity and evaluation of the generated explanation by XAI methods

and c) vulnerability of black box explanation method toward adversarial attacks. We therefore overview the literature on each of this direction in the following subsection.

2.1. Visual Explanation

For an end-user, visual explanation methods makes it easier to understand the prediction output of the black box model. One common approach to generate such a visualization is done via **saliency maps** and such algorithms may be divided into the following four categories: back-propagation based methods, gradient-based methods, perturbation based methods and Approximation based methods. **Back-propagation methods:** back propagation methods spread a feature signal from an output neuron rearwards through the layers of model to the input in a single pass; making them efficient. Layer wise Relevance Propagation [18] and DeCovNet [19] are examples of this category. **Gradient-based methods:** methods employ the gradient or its modified version in the backpropagation algorithm to visualize the derivative of the CNN’s output w.r.t to its input, e.g. such as Grad-CAM [5]. A better way to produce input images that greatly activate a neuron a neuron was proposed in [20]. They generated class-specific saliency maps by performing a gradient ascent in pixel space to reach a maxima. This synthesized image serves as a class specific visualization and helps understand how a given CNN modeled a class. **Perturbation-based methods:** here, the input is perturbed while keeping track of the overall changes to the output. In some work, the change occurs at intermediate layers of the model. The state-of-the-art RISE [17] algorithm belongs to this category. Meaningful perturbations [21] optimize a spatial perturbation mask that maximally affects a model’s output and reveal a new image saliency model that seeks where an algorithm looks by finding out which regions of an image most affects its output level when perturbed. **Approximation-based method:** involves replacing the deep CNN by simpler approximation model where visual explanation can be generated easier. A good example of this class is the LIME algorithm [10]. Decision trees can also be used for this purpose but cannot be explicitly applied

to any visual input [22]. DeepLift [23] evaluates the importance of each input neuron for a particular decision by approximating the instantaneous gradients (of the output with respect to the inputs) with discrete gradients. This obviates the need to train interpretable classifiers for explaining each input-output relation (as in LIME) for every test point. Our proposed SIDU [13] method falls under perturbation-based methods but can effectively localize entire salient region of the object of interest compared to the state-of-the-art XAI methods such as Grad-CAM and RISE. Furthermore, it is less computationally complex.

2.2. Evaluation of Explanation Methods

Since it is rather challenging to establish a unique and generalized evaluation metric that can be applied to any task, authors in [14] proposed three different types of evaluations to measure the effectiveness of explanations. These are presented in the following.

1. **Application-Grounded evaluation:** Application-Grounded evaluation includes carrying out human experiments within a real application. If the researcher has a concrete application in mind—such as teaming up with doctors on diagnosing patients with a specific disease—the best method to show that the design is effective is to assess it with respect to the task. A sound experimental setup and knowing how to evaluate the quality of the elucidation are needed. It is based upon how well a human can expound the same decision. Human expert level evaluation is necessary for those end-users who have less confidence in the results of prediction model (e.g. clinician).
2. **Human-Grounded evaluation:** Human-Grounded evaluation involves conducting basic human-subject experiments that sustain the heart of the target application. An approach is appealing when experiments with the target community are difficult. The evaluations can be completed with layperson, creating a greater subject pool and cutting down expenses, since we do not have to pay highly trained domain experts.

3. Functionally-Grounded evaluation: This method utilizes numeric metrics or proxies such as local fidelity to evaluate explanation in different applications. The main advantage of this evaluation is that it is free from human bias that effectively saves time and resources. Most of the state-of-art methods fall into this category [19, 24, 21]. For example, the authors in [12] proposed casual metrics *insertion* and *deletion*, which are independent of humans to evaluate the faithfulness of the XAI methods.

2.3. Adversarial Attacks

In the context of XAI, adversarial attack generators can be divided into white box attacks and black box attacks. The Fast Gradient Sign Method (FGSM) [25] and Projected Gradient Descent (PGD) [26] algorithms are excellent examples of white box attack where small amount of noise is added to an image that is not visually detectable by the end user. In case of black box attacks, the adversarial attack happens through various mechanisms to fool the model’s classifier and alter its outcome. The majority of the proposed approaches in this class are based on perturbing the model input either globally or locally. For instance, DeepFool [27] attack can be characterized by performing pixel-wise perturbation of an image while adversarial patch [28] attempts to change the pixel values in a specific region of an image. In general, the ability of changing model’s output via small input perturbations makes the XAI explanation methods challenging and less reliable. Thus, to establish greater trust, it is essential for the XAI algorithms not only be effective but also robust against adversarial attack at the same time [29]. Analyzing how the black box explanation (such as proposed SIDU) can effectively handle such potential problem and help the end user to guard against possible disastrous outcome of classifier when adversarial attack is presented.

3. SIDU: Proposed Method

Recent XAI methods have shown that deeper representations in CNN models illustrates higher-level visual features [30] [31]. A recent approach Grad-

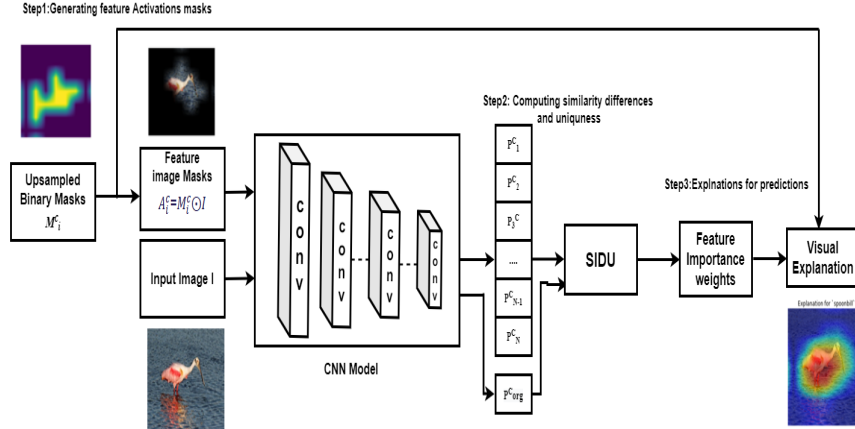


Figure 2: Block diagram of SIDU.

CAM [11] interprets the importance of each neuron responsible for a decision of interest by computing the gradient information from the last convolutional layer of the CNN. Alternatively, the authors in [17] proposed a method RISE, which finds the effect of selectively inserting or deleting parts of the input (*perturbation-based*) in the CNN model’s output prediction. This perturbation-based method provides more accurate visual explanation saliency maps compared to gradient based methods, but still they fail to visualize all the perturbations and determine which one characterizes the best desired explanation. Furthermore, the visual explanations generated by both the gradient-based and perturbation explanations methods failed to localize the entire salient regions of an object class responsible for higher classification scores.

To overcome the challenges of the most recent state-of-the art methods we proposed a XAI method that consequently provides better explanation method for a any given CNN model. The proposed method takes the last convolution layer for generating the masks. From these masks Similarity Difference and Uniqueness scored are computed to get the explanation of the final CNN model decision acronymed in therefore denoted SIDU. An overview of the proposed

method is presented in Figure 2. Our method is composed of three steps, First we extract the last convolution layer of the CNN to generate the mask using the last convolution layer of the given model. Second, we compute the similarity differences for each mask with respect to predicted class and finally we compute the weights of each mask and combine them into a final map which shows the explanation of the prediction. Each step is described in the following subsections.

3.1. Step1: Generating Feature Activation Masks

To provide a visual explanation of the predicted output of a CNN model, we first generate feature image masks from the last convolution layers. For any deep CNN model F , we consider the last convolution layers of size $n \times n \times N$ where ' n ' is size of that convolution layer and ' N ' is the total number of features activation \mathbf{f} of class c , i.e., $\mathbf{f}^c = [f_1^c, \dots, f_N^c]$. Each feature activation map f_i^c is then converted into a binary mask M_i^c by thresholding each value. The binary mask M_i^c is then up-sampled by applying bi-linear interpolation for a given input image I with size of $Width \times Height$. Next, the binary mask M_i^c will have values between $[0, 1]$ and it is no longer binary. The up-sampled binary masks are also known as feature activation masks. Finally, point-wise multiplication is performed between feature activation mask (Up-sampled binary mask) M_i^c and input image I to get the feature image mask A_i^c and is represented as

$$A_i^c = F(I \odot M_i^c) \quad (1)$$

where F is an CNN model, A_i^c is the feature activation image mask of feature map f_i^c and $i = 1, \dots, N$. The procedure of generating feature activation masks is shown in Figure 3. The feature images masks A^c of object class c are used to get prediction scores which is explained in detail in the following subsection 3.2

3.2. Step2: Computing Similarity Differences and Uniqueness

The total number of feature image mask is dependent on the number of activation's in the last convolution layer of the CNN model. Let the last convolution layer of the CNN model be of size $n \times n \times N$. The total number of feature

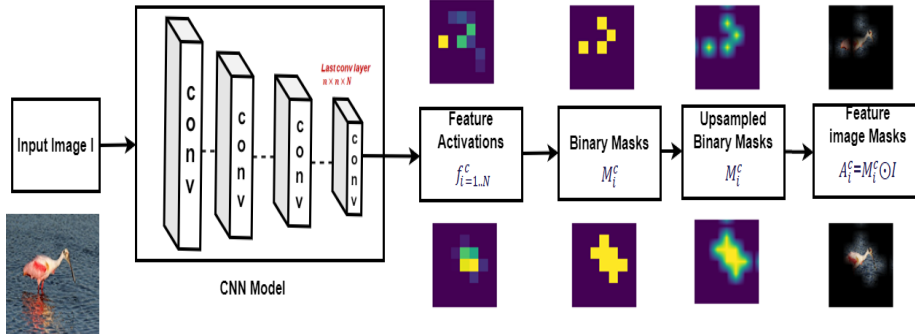


Figure 3: Block diagram of generating feature image masks.

activation masks will be N . Next, we compute probability prediction scores for all the feature activation image masks \mathbf{A}^c of class c , i.e., $\mathbf{A}^c = [A_1^c, \dots, A_N^c]$ individually using the same CNN model F used for generating feature image masks. The probability prediction score of the feature activation image mask A_i^c is defined as P_i^c and the probability prediction score for the given input image I is defined as P_{org}^c . The prediction scores vector size will depend on the total number of classes used for training the CNN model. For example, if the CNN model is trained on ImageNet dataset, which has a total of 1000 object classes, then the size of the prediction scores vector P_i^c of each individual feature image mask A_i^c will be 1×1000 , where $i = 1 \dots N$. Figure 4 illustrates the procedure of computing the prediction scores vector.

Once the prediction score vectors are computed for all feature image masks and original input image, we then compute similarity differences between each input feature activation image mask prediction score P_i^c and prediction score P_{org}^c of the original input image I . The similarity difference between these two vectors gives the relevance of feature image mask w.r.t original input image. The intuition behind computing the relevance of a feature map is to measure how the prediction changes if the feature is not known, i.e., the similarity difference between prediction scores. The relevance value of the feature activation image mask will be high if it is similar to the predicted class and low otherwise. The similarity difference of the feature activation maps SD_i^c is given by

$$SD_i^c = \exp\left(\frac{-1}{2\sigma^2}\|P_{org}^c - P_i^c\|\right) \quad (2)$$

where P_{org}^c , P_i^c are the predictions for the original input image and feature image mask generated from the last convolution layer and σ is an controlling parameter, respectively.

After computing the similarity difference measure, we also computed a uniqueness measure U^c between the feature image masks prediction score vectors. It is one of the most popular assumptions that the image regions which stand out from the other regions grab our attention in certain aspects. Therefore the region should be labeled as a highly salient region. We therefore evaluate how different each respective feature mask is from all other feature masks constituting an image. The reason behind this is to suppress the false regions with low weights and highlight the actual regions which are responsible for predictions with higher weights. The uniqueness measure U_i^c is defined as

$$U_i^c = \sum_{j=1}^N \|P_i^c - P_j^c\|, \quad (3)$$

Finally, the weight of each feature importance W_i^c is computed as the dot product of the similarity difference SD_i^c and uniqueness measure U_i^c where

$$W_i^c = SD_i^c \cdot U_i^c, \quad (4)$$

where SD_i^c, U_i^c are the similarity difference and uniqueness values for the feature activation image mask A_i^c of the object class c . The feature importance weight will be high for the feature which has more influence in predicting the actual class object c and low for the feature with low influence.

3.3. Step3: Explanations for the prediction

To get the explanation of the predicted output class c of a CNN model, we then perform a weighted sum between feature activation mask A_i^c and the corresponding feature importance weights W_i^c , where the weights are computed by Eq. 4. The visual explanation map is in the form of a heatmap and is

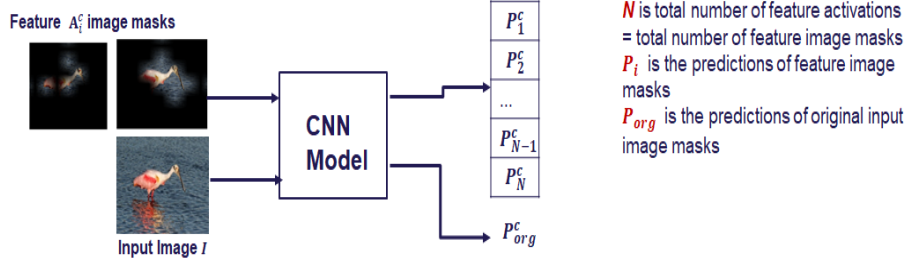


Figure 4: The predictions scores vectors for the each individual feature image masks and the original image is computed from the CNN model. These prediction score vectors are used for computing similarity differences and uniqueness and finally the dot product is done to get the feature importance weights. Note that the CNN model is same for all the steps.

The equation shows a series of terms: $w_1^c \cdot$ followed by a heatmap, $+ w_2^c \cdot$ followed by another heatmap, $+ \dots + w_N^c \cdot$ followed by a third heatmap, followed by an equals sign and a final heatmap. The final heatmap is labeled "Visual explanation".

Figure 5: Visual explanation for the prediction. The visual explanation is a weighted linear combinations of feature activation masks for the prediction of the class

represented as S_c and is shown in Figure 5. The visual explanation map S_c is also known as the class discriminative localization map. Thus the visual explanation of the predicted class c is given by

$$S_c = \frac{1}{N} \sum_i^N W_i^c \cdot A_i^c \quad (5)$$

The weighted combinations of feature activation masks to get the final visual explanation of the prediction of the class is shown in Figure 5.

In summary, for explaining the decision of the predicted class c visually, we first generate the N feature activation masks (up-sampled binary masks) from the last convolution layer of deep CNN model F which has N number of feature activation maps of size $n \times n$. We then perform point wise multiplication between each generated up-sampled binary mask M_i and the input image I to

Table 1: Comparison of XAI methods using Resnet50 on ImageNet validation set.

METHODS	Insertion \uparrow	Deletion \downarrow
RISE [17]	0.63571	0.13505
GRAD-CAM [5]	0.62863	0.15399
SIDU	0.65801	0.13424

get feature activation image mask. Next, we compute similarity differences SD_i^c and uniqueness measure U_i^c using predictions scores of feature activation image mask A_i . Feature importance weights W_i of each feature activation image mask A_i computed by the dot product of SD_i^c and U_i^c . Finally, the visual explanation S_c of a given input image is obtained by a weighted sum of feature activation image masks A_i given in Eq. 5. An example is shown in Fig. 6.

4. Evaluation

In this section we evaluate the performance of the proposed visual explanation method. We conducted a comprehensive set of experiments to study the correlation of the visual explanation with the model prediction to evaluate the faithfulness. The proposed method (SIDU) is evaluated using all the three categories of evaluations [14], functionally grounded, application grounded and human grounded evaluations. The evaluation results are compared with most recent state-of-the art methods RISE [17] and GRAD-CAM [5]. A good explanation method can not only provides the appropriate explanation for the prediction but also it should be robust against to adversarial noise. To this end the proposed method is evaluated on adversarial samples and compared with most recent state-of-the art methods RISE [17] and GRAD-CAM [5]. The experimental evaluation of faithfulness of the SIDU model on the above mentioned evaluation categories and effect of adversarial noise are described in section 4.1, 4.2, 4.3 and 4.4, respectively.

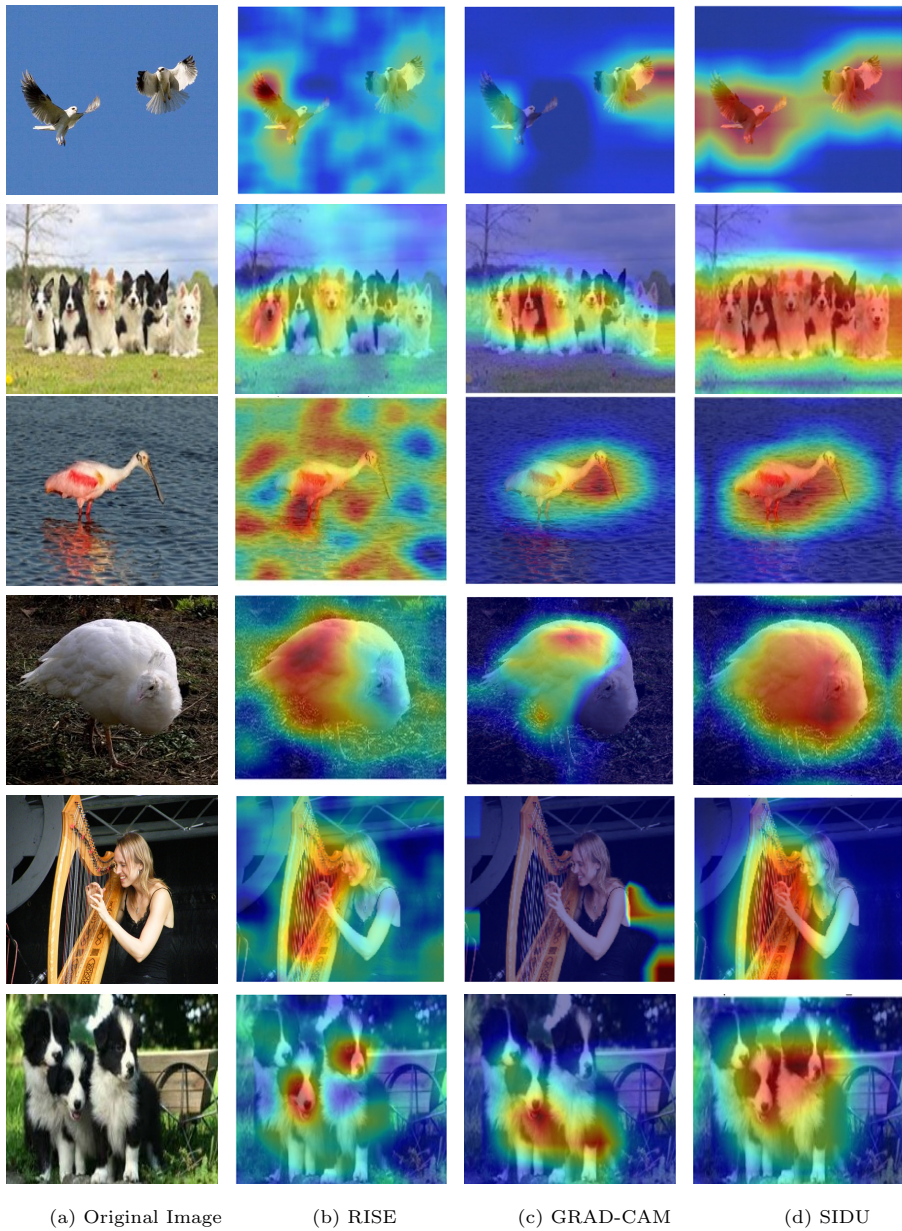


Figure 6: Visual comparison of explanation maps generated for the natural images class of ImageNet.

Table 2: Comparison of XAI methods using VGG-16 on ImageNet validation set.

METHODS	Insertion \uparrow	Deletion \downarrow
RISE [17]	0.47113	0.1313
GRAD-CAM [5]	0.41720	0.15486
SIDU	0.49419	0.1309

Table 3: Comparison of XAI methods on RFIQA dataset using trained ResNet-50 model.

METHODS	Insertion \uparrow	Deletion \downarrow
RISE [17]	0.75231	0.59632
GRAD-CAM [5]	0.91303	0.43061
SIDU	0.87883	0.47818

4.1. Functionally grounded Evaluation

To perform the functionally grounded evaluation we choose two automatic causal metrics *insertion* and *deletion* proposed by [17]. The *deletion metric* deletes the saliency region in the image which is responsible for higher classification scores and forces the CNN model to change its decision. This metric estimates the decrease in the probability classification scores, when more pixels are removed from the saliency region. With the deletion metric, the good explanation shows a sharp drop in the predicted score and area under the probability curve will be lower. Whereas, the *insertion metric* measures the probability increase of predicted score. As more pixels are inserted in the image, a higher AUC rate can be achieved (i.e., effectiveness of explanation model at a greater level). We choose these metrics since they are independent of human subjects, bias free and hence increase transparency when evaluating the XAI methods. In order to evaluate the performance of the SIDU explanation method we choose two datasets with different characteristics. The ImageNet [32] dataset of Natural Images with 1000 classes. We use 2000 images randomly collected from the ImageNet validation dataset. The other is a Retinal Fundus Image Quality Assessment (RFIQA) dataset from the medical domain. The dataset consists

of 9,945 images with two levels of quality, 'Good' and 'Bad'. The retinal images were collected from a large number of patients with retinal diseases [33].

We conducted two experiments for evaluating the faithfulness of the proposed explanation method. The first experiment is performed on ImageNet validation dataset. We collected randomly 2000 images from the ImageNet dataset. To do a fair evaluation, we choose two existing standard CNN models, ResNet-50 [34] and VGG-16 [35] pre-trained on imageNet dataset [32]. Table 1 summarizes the results obtained on ResNet-50 for the proposed method and compares it to most recent works [17] and GRAD-CAM [11]. We can observe that the proposed method achieved better performance for both metrics, followed by RISE [17] and GRAD-CAM [11]. Table 2 summarizes the results obtained on VGG-16 model for the proposed method and compares it to most recent works [17] and GRAD-CAM [11]. The proposed method, SIDU achieves best performance. From the Tables 1 and 2, we can observe that the values are better for ResNet-50 than VGG-16 for all the XAI methods, which suggests that ResNet-50 is a better classification model than VGG-16. Qualitative examples are shown in Fig. 6. In our proposed method, the generated masks come from the last feature activation maps of the CNN model. Due to this the final explanation map will localize the entire region of interest (object class).

We also conducted a second experiment on Medical Image dataset which has totally different characteristics. We trained the existing ResNet-50 [34] with an additional two FC layers and softmax layer on the RFIQA dataset [33]. The CNN model achieves 94% accuracy. The proposed explanation method uses the trained model for explaining the prediction of the RFIQA test subset with 1028 images. The evaluated results of the proposed method and RISE and GRAD-CAM are summarized in Table 5. We can observe that the GRAD-CAM achieves slightly higher AUC for *insertion* and lower AUC for *deletion* followed by SIDU. RISE [17] has shown least performance in both the metrics, This can be explained by the fact that the RISE method generates N number of random masks and the weights predicted for these masks give higher weights to false regions which makes the final map of RISE noisy. The visual explanations of the

proposed method (SIDU) and the RISE, GRAD-CAM methods on the RFIQA test dataset are shown in Figure. 9 (b), (c), (d).

4.2. Human grounded Evaluation

Human-grounded evaluation is most appropriate when one aims at testing a general notions of an explanation quality. Therefore, for generic applications in AI domain such as object detection and object recognition, it might be sufficient to inspect a degree to which a non-expert human can understand the cause of a decision generated by a black box model. One excellent way to measure and compare the correlation of visual explanation between a human subject and the black box is to use an eye tracker that records the non-expert subject’s fixations in interactive settings. This approach is chosen because of its similarity to XAI methods, visual explanations. Both generate heatmaps representing salient areas of an object in an image.

An eye tracker is used for gathering eye tracking data from human subjects to gain an understanding of visual perception [36, 37, 37]. The study using eye tracking data for understanding the human visual attention is useful but difficult and expensive to collect on a large quantity [36, 37]. The authors in [38] established mouse tracking approach to collecting attention maps accurate. They collected large-amount of attention annotations for MS COCO on Amazon Mechanical Turk (AMT). In [39] recorded eye-fixation body parts for investigate which body parts of virtual characters are most looked at in scenes containing duplicate characters or clones. The authors in [40] conducted an experimental study and gathered data on a large-scale of ‘human attention’ in Visual Question Answering (VQA) to interpret where the humans choose to look to answer the questions regarding the images. However, all these experimental studies have used eye tracking to understand the human visual attention for different types of problems.

In our study, we investigate how non-expert subjects generated explanations via the eye-tracker compares with those of generated by XAI visual explanation methods across natural images for recognizing object class. To this end, we

follow the data collection protocol discussed in detail next section 4.2.1.

4.2.1. Database of eye tracking data

We sample randomly 100 images from 10 different classes of ImageNet [32] benchmark validation dataset. All the collected images are RGB and are resized to 224×224 pixels.

4.2.2. Data collection protocol

In order to collect fixation, 5 human subject participate in an interactive procedure using Tobii-X120 eye-tracker [41] in the following main steps:

1. The subject sits in front of screen where the eye-tracker is ready to record the eye fixation.
2. After careful initial calibration, each image from the dataset is shown in a random order for 3 seconds and corresponding fixations of the subject are recorded.
3. We divide all 100 images into 4 equally sized data blocks with a break between each experiment in order to reduce the burden on each subject. We further add the Cross-fixation image between two stimulus to reset the visionary fixation on the screen while going from one image to another image.
4. The participants are shown random images from the collected dataset and have asked a question, what kind of object class is presented in the image.
5. The eye fixations of each individual participant will be automatically recorded via-eye tracker when the participant looks into the image for recognizing the object class.
6. Once all 5 participants' fixation are collected, an aggregated heatmaps is generated by convolving a Gaussian filter across each user fixation for each image see, Figure 7. The resulting heatmaps highlight the salient

regions of each object class that often attracted attention of all subjects in the experiment and hence can be used to compare with the heatmaps produced by the XAI explanation algorithms.



Figure 7: Examples of Eye-tracking data collection from humans for recognizing the given object class

4.2.3. Comparison Metrics

To ensure a robust evaluation we use the three metrics to compare the XAI and eye-tracker generated heatmaps [42]. These are Area Under ROC Curve (AUC), Kullback-Leibler Divergence (KL) and Spearman's Correlation Coefficient metric (SCC) metrics.

1. **Area under ROC Curve (AUC):** The Receiver Operating Characteristics (ROC) is one of the commonly used metric for assessing the degree of similarity of two saliency maps. It is represented in the form of a graphical plot which describes the tradeoff between true and false positive at different thresholds [42]. A fraction of true positives from the total actual positives is plotted against the false positives' fraction out of the total actual negatives to create ROC. This is denoted as TPR, representing the true positive rate, and FPR that indicates the false positive rate. The rates are examined at different threshold values. If a TPR value of 1 is achieved at 0 FPR, the prediction method is good. These values will yield a point in the ROC space's upper left corner and correspond to a near-perfect classification. Conversely, when the guess is completely random, it will generate a point along a diagonal line starting at the left bottom and going up towards the top right corner. If the diagonal divides the ROC

space while and points above the diagonal, this represents good classification results. Such results are considered better than random results. On the other hand, the line below is a sign of poor results, which is even worse than getting random results. The Area Under Curve (AUC) is the method used to measure the ROC curve’s performance. The AUC is equal to the probability of a classifier ranking a randomly selected positive instance, which is usually higher than a randomly selected negative instance, assuming that the positive ranks higher than a negative. To compute the AUC, XAI visual explanation heatmaps are treated as fixations’ binary classifiers at numerous threshold values or value sets. The true and false positive rates are measured under each binary classified or level set to sweep out the ROC curve.

2. **Kullback-Leibler Divergence (KL-DIV):** The Kullback-Leibler Divergence is an metric, which is used to measure dissimilarity between two probability density functions [42]. For evaluating the XAI methods, eye-fixation maps and the visual explanation maps produced by the model are used for the distributions. FM represents the heatmaps probability distribution from eye-tracking data, and EM indicates the visual explanation maps probability distribution. These probability distributions are normalized and they are given by :

$$EM(x) = \frac{EM(x)}{\sum_{x=1}^X EM(x) + \epsilon}, \quad (6)$$

$$FM(x) = \frac{FM(x)}{\sum_{x=1}^X FM(x) + \epsilon}, \quad (7)$$

where X is the number of pixels and ϵ is a regularization constant to avoid division by zero. The KL-DIV measure is computed between these two distributions to know whether the visual explanation map which is computed from the XAI method matches human fixations. It is a non-

linear measure and generally varies in ranges zero to infinity. If the KL-DIV measure between EM and FM is lower, then the EM maps have better approximation of the human eye-fixation FM .

3. **Spearman's Correlation Coefficient (SCC):** Spearman's correlation is a non-parametric measure that analyses how well the relationship between two variables can be described using a monotonic function [43]. It is a statistical method used mainly for measuring the correlation or dependency between two variables. This metric varies between -1 and 1 , where a score of -1 , represents no correlation. The SCC between two variables will be high when observations have a similar (correlation close to 1) rank between the two variables, and low when observations have a dissimilar (correlation close to -1) rank between the two variables [43]. It is an appropriate measure for both continuous and discrete ordinal variables [43]. FM represents the heat map from eye tracking data, whereas EM is the visual explanation map. The SCC between the two random variable maps, FM and EM are given by :

$$SCC(EM, FM) = \frac{cov(EM, FM)}{\sigma(EM) \times \sigma(FM)}, \quad (8)$$

where $cov(EM, FM)$ is the covariance of EM and FM , $\sigma(EM)$ and $\sigma(FM)$ are the standard deviations of EM and FM respectively.

4.2.4. Comparing SIDU and State-of-art methods with human attention for recognizing the object classes

In this experiment, we use the Imagenet images eye-tracking data recordings described in section 4.2.1 to generate and evaluate the explanation by the XAI algorithms. To this end, we first generate ground truth heatmaps by applying Gaussian distributions on human expert eye-fixations. These heatmaps are then used to compare with the XAI heatmaps. AUC, SCC and KL-DIV evaluation metrics are used to evaluate the performance. We finally calculate the mean of AUC, SCC and KL-DIV of all the images the in dataset. Table 4 summarizes the results obtained by SIDU and two different state-of-the art XAI methods

Table 4: saliency maps of XAI methods with eye fixation maps.

METHODS	mean KL-DIV↓	mean SCC ↑	mean AUC ↑
RISE [17]	8.4384	0.1967	0.6385
GRAD-CAM [5]	9.7892	0.2711	0.6828
SIDU	4.3027	0.3314	0.7708

RISE [12] and GRAD-CAM [11] on our proposed imagenet eye-tracking data. We can observe that, SIDU outperforms GRAD-CAM and RISE in all the three metrics. Therefore we can conclude that SIDU explanations matches closer with the human explanations (heatmaps) for recognizing the object class.

4.3. Application grounded evaluation

Application grounded evaluation involves conducting experiments within a real application to assess the trust of the black box models. We choose an medical case. As a test application we use the task of retinal fundus image quality assessment [33]. The application is about screening for retinal diseases, where poor-quality retinal images do not allow an accurate medical diagnosis. Generally, in sensitive domains such as clinical settings, the domain experts (here clinicians) are skeptical in supporting explanations generated by AI diagnostic tools as a result of high involved risk [6, 8].

In our experimental setup at hospital, two ophthalmologist participated to evaluate which visual explanation results in more trust and further aligns with actual physical examination performed in the clinic. This experiment assesses the effectiveness of the proposed method in terms of localizing the exact region for predicting the retinal fundus image quality w.r.t state-of-art method. Here, the generated visual explanation heatmaps in the RISE algorithm were used for comparison. We follow the similar setting as discussed in [5], i.e., using both the proposed method and the RISE method, visual explanation heatmaps of 100 retinal fundus images for two classes of ‘Good’ and ‘Bad’ quality are recorded. The explanation methods uses the trained model described in section 4.1 for explaining the prediction of the retina fundus images. Both ophthalmologists have

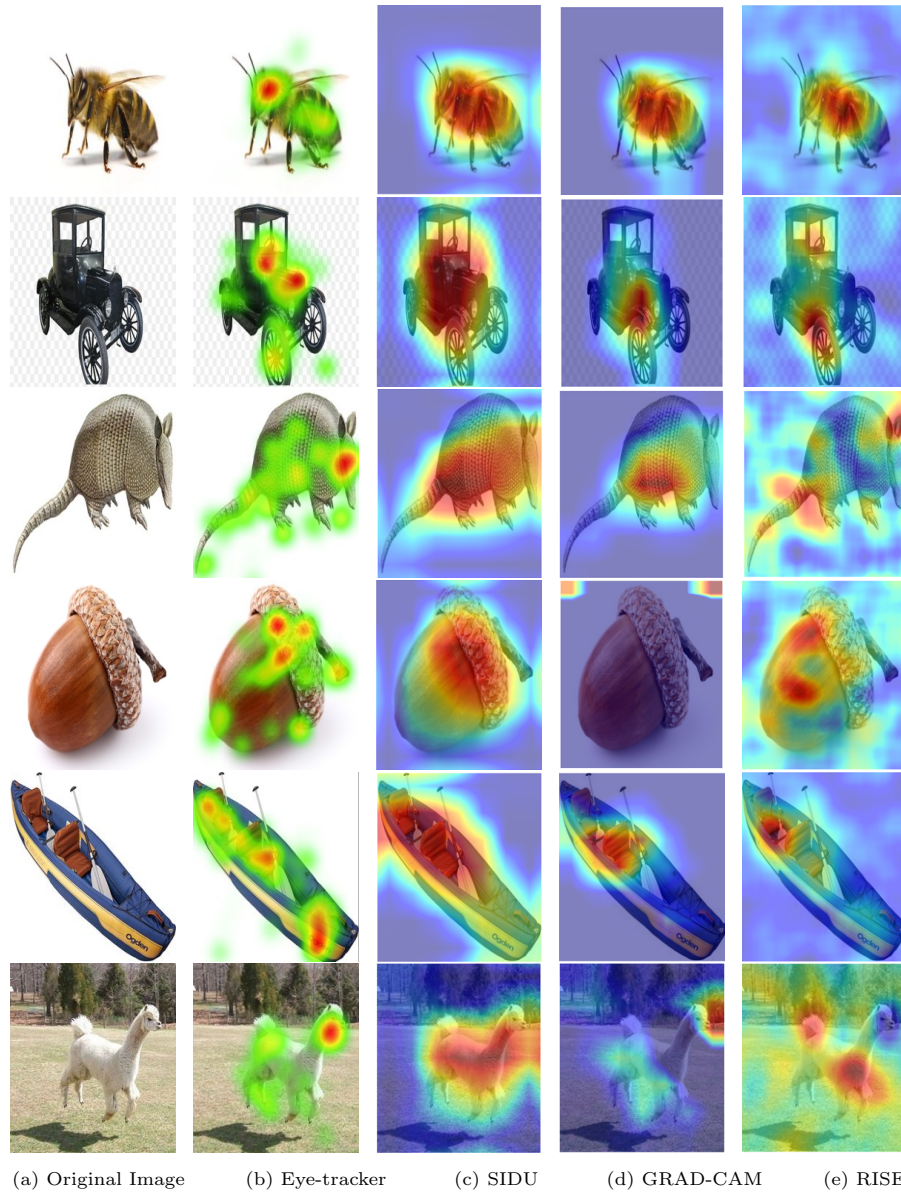


Figure 8: Comparison of XAI methods visual explanation with human visual explanation (heatmaps). The generated heatmaps in 3rd, 4th and 5th columns by the SIDU, GRAD-CAM and RISE demonstrate how the visual explanation methods are closely aligns with human.

Table 5: Expert level evaluation of XAI methods on medical RFIQA dataset .

METHODS	Expert I	Expert II
RISE [17] (Method I)	0.02	0.05
SIDU (Method II)	0.84	0.93
BOTH	0.14	0.02

no prior knowledge about any explanation model presented to them. The two explanations methods are labelled as either method I and method II to participants involved in experiments. The participants can opt for “both” methods if they feel that both explanations are rather similar. Therefore, the ophthalmologist will have three different options for every test image. Once the ophthalmologist determined which method better localizes the regions of interest (good/bad quality regions) for each image, we then calculate the relative frequency of each outcome per total retinal fundus image. Table 5 summarizes the results of the two methods evaluated by experts (ophthalmologist). We observe that, in the case of the first ophthalmologist, the RISE explanation map was selected with the relative frequency of 0.02, the proposed method, SIDU with 0.84 and 0.14 being the same. For the second ophthalmologist, the relative frequencies are 0.05, 0.93 and 0.02, respectively. Therefore, the experiments conclude that, the proposed method gains greater trust by the both ophthalmologists and visual explanations in Fig. 9 further supports this claim.

4.4. Effect of Adversarial Noise on XAI methods

Despite the success in many applications of AI, recent studies find that Deep Learning is against well designed input samples know as adversarial examples poses a major challenge [29]. Adversarial examples are carefully perturbed versions of the original data that successfully fool a classifier. In the image domain, for example, adversarial examples are images that have no visual difference from natural images, but that lead to different classification results. How resilient different XAI algorithms are towards adversarial examples is a largely overlooked topic. In this subsection we therefore investigate exactly that.

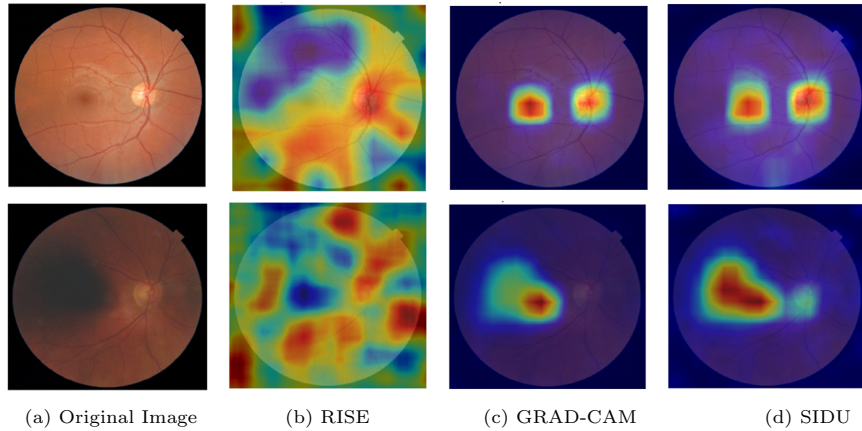


Figure 9: The visual explanation of Good / Bad quality eye fundus images $\langle(B), (C), (D)\rangle$ from RFIQA dataset by RISE, GRAD-CAM and the SIDU method with ResNet50 as the base network. In real scenario, the doctors observe the visibility of the optical disc and macular regions in a good quality image (1^{st} image, 1^{st} row) corresponds to the region highlighted in the visual explanation heatmap of the proposed method. The bad quality image (2^{nd} image, 2^{nd} row) is due to the shadow which is observed just near to center of the image (optical disc), i.e., exactly the region highlighted by the proposed method.

To perform this experiment, we choose one the most successful white box attacks, namely, gradient based attacks. Fast Gradient Sign Method (FGSM) [25] and Projected Gradient Descent (PGD) [26] are the examples of such attacks. PGD is an iterative application of FGSM. The process of PGD is more complex and time consuming. Therefore, we have chosen the Fast Gradient Sign Method (FGSM) because of its simplicity and effectiveness. The adversarial image is generated using FGSM by adding noise to an original image. The direction of this noise is the same as the gradient of the cost with respect to the input data. The amount of noise can be controlled by a coefficient, ϵ . By applying this coefficient properly, it will change the model predictions and it is undetectable to a human observer. Figure 10 shows the different levels of FGSM adversarial noise added to the original image. We conducted two different experiments using adversarial noise to demonstrate the effectiveness of SIDU, compared to state-of-the art methods RISE and GRAD-CAM. The experiments are described in following.

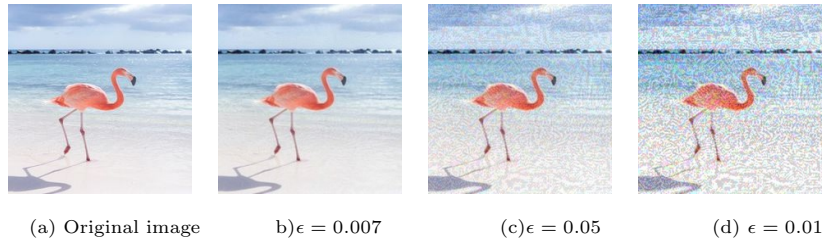


Figure 10: Example of a natural image in its original form and also with three different levels of FGSM noise, together with the corresponding predictions.

4.4.1. How do XAI method visual explanations heatmaps of adversarial examples deviate from human eye-fixation heatmaps?

In this experiment, we analyse how robust the XAI methods are against an adversarial attack in terms of generating reliable explanations. Reliable visual explanations are defined in terms of resemblance to the human eye-fixation heatmaps. To conduct this experiments we choose the same pre-trained ResNet-50 model used in section 4.1. We first apply the FGSM noise with different epsilon levels to the dataset of 100 images collected from Imagenet validation set described in section 4.2.1. We choose three different optimal noise coefficients between 0 and 1, and the chosen values are $\epsilon = 0.007$, $\epsilon = 0.05$ and $\epsilon = 0.1$. These values are optimal because, it is sufficient enough to pass unnoticeable by the human eye. We extracted the visual explanations heatmaps using the proposed method (SIDU), RISE [12] and GRAD-CAM [11]. The heatmaps generated by SIDU, RISE and GRAD-CAM methods are finally compared with human generated visual explanations using eye-tracker described in section 4.2.1 using the three evaluation metrics AUC, SCC and KL-DIV. Table 6, 7 and 8 summarizes the mean AUC, SCC and KL-DIV results. From the table we can observe that, SIDU outperforms GRAD-CAM and RISE for different levels of adversarial noise with all the three evaluation metrics. We also observe that, the performance of XAI methods decreases with all the three metrics with the increase in adversarial noise to the original images. From this we can conclude that the proposed method (SIDU) is more robust to adversarial noise than RISE

Table 6: visual explanation heatmaps from adversarial noise 0.007 with eye fixation heatmaps.

METHODS	mean KL-DIV↓	mean SCC ↑	mean AUC ↑
RISE [17]	8.0547	0.2121	0.6526
GRAD-CAM [5]	10.3257	0.2530	0.6719
SIDU	4.3785	0.3309	0.7689

Table 7: visual explanation heatmaps from adversarial noise 0.5 with eye fixation heatmaps.

METHODS	mean KL-DIV↓	mean SCC ↑	mean AUC ↑
RISE [17]	9.3305	0.1995	0.6380
GRAD-CAM [5]	11.6447	0.2229	0.6431
SIDU	4.8492	0.2929	0.7397

Table 8: visual explanation heatmaps from adversarial noise 0.1 with eye fixation heatmaps.

METHODS	mean KL-DIV↓	mean SCC ↑	mean AUC ↑
RISE [17]	9.1246	0.2068	0.6461
GRAD-CAM [5]	12.3077	0.2112	0.6281
SIDU	4.2239	0.2817	0.7364

and GRAD-CAM, as is visually evident in the Figures 11. We can observe that SIDU localizes the entire actual object class after adding the three different levels of adversarial noise, whereas the other methods completely loose the actual object class localization after adding the noise.

4.4.2. How do visual explanation maps from adversarial examples deviate from original visual explanation maps?

In this experiment, we analyse how the visual explanation from adversarial noise added examples of XAI methods deviate from original image visual explanation maps. To conduct this experiments we choose the same pre-trained ResNet-50 model used in section 4.1. We first apply the FGSM noise with different epsilon levels to the dataset of 100 images collected from Imagenet validation set described in section 4.2.1. We choose one noise level $\epsilon = 0.1$ for these experi-

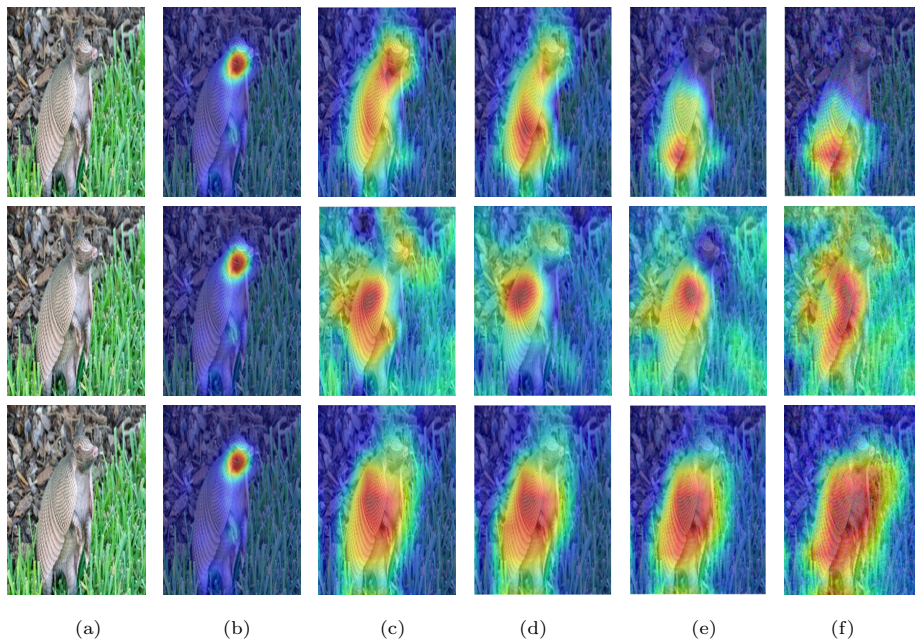


Figure 11: Comparison of XAI methods visual explanation with different levels of FGSM noise with human visual explanation (heatmaps). The generated heatmaps on adversarial noise levels $\epsilon = 0.007, 0.5, 0.1$. in 3^{rd} , 4^{th} and 5^{th} columns by the SIDU, GRAD-CAM and RISE demonstrate. (a)Original Image (b) Eye-tracker (c) $\epsilon = 0$ (d) $\epsilon = 0.007$ (e) $\epsilon = 0.05$ (f) $\epsilon = 0.1$

Table 9: visual explanation heatmaps from adversarial examples deviate from original visual explanation heatmaps.

METHODS	mean KL-DIV↓	mean SCC ↑	mean AUC ↑
RISE [17]	9.6665	0.2385	0.6133
GRAD-CAM [5]	10.0077	0.4061	0.6875
SIDU	2.4924	0.6488	0.8347

ment. We extract the visual explanations heatmaps using the proposed method (SIDU), RISE [12] and GRAD-CAM [11] for the original images without noise and with noise $\epsilon = 0.1$. The heatmaps generated by SIDU, RISE and GRAD-CAM methods are finally compared with the original images, visual explanations to see adversarial noise added images are deviated from the original ones using the three evaluation metrics AUC, SCC and KL-DIV. Table 9 summarizes the mean AUC, SCC and KL-DIV results obtained by the XAI methods. From the table we can observe that, SIDU outperforms GRAD-CAM and RISE for all the three evaluation metrics. From the Figure 11, we can clearly observe that the propose method(SIDU) doesn't deviate much in localizing the object class which is responsible for the prediction. Therefore from these two adversarial noise experiments we can conclude that our proposed method is more robust against adversarial noise.

5. Conclusion

In this work, we proposed a novel method called SIDU for explaining the CNN model, predictions visually in a form of heatmap through feature activation maps of the last convolution layers in the model. The proposed method is independent of gradients and can effectively localize entire object classes in an image which is responsible for CNN prediction. The new explanation approach helps in gaining more trust in prediction results of CNN model by providing further insights to the end-user in sensitive-domain. We validated the effectiveness of our method by conducting three different XAI evaluations methods, applica-

tion grounded (invoking human experts trust in medical domain), functionally grounded (using an automated causal metrics independent of humans) and human grounded evaluation. For the human grounded evaluation, we proposed a framework for evaluating explainable AI (XAI) methods using an eye-tracker. It is designed specifically for evaluating XAI methods using non experts to understand the human visual perception for recognizing the given object class and compared it with visual explanations of standard well-known CNN models on natural images. We also carried out experiments on adversarial examples. Our proposed method outperforms compared to state-of-the-art methods.

Future work involves, extending our proposed method (SIDU) to spatio-temporal CNN models to provide visual explanations for video applications tasks such as video classification and action recognition. Further more, exploring the possibility of extending our method to explain decisions made by other neural network architectures such as long short-term memory networks and in other domains such as Natural Language Processing (NLP).

References

- [1] R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2014, pp. 580–587.
- [2] M. Zhang, W. Li, Q. Du, Diverse region-based cnn for hyperspectral image classification, *IEEE Transactions on Image Processing* 27 (6) (2018) 2623–2634.
- [3] J. Zhang, Y. Xie, Q. Wu, Y. Xia, Medical image classification using synergic deep learning, *Medical image analysis* 54 (2019) 10–19.
- [4] S. Mohseni, J. E. Block, E. D. Ragan, A human-grounded evaluation benchmark for local explanations of machine learning (2020). [arXiv:1801.05075](https://arxiv.org/abs/1801.05075).

- [5] R. R. Selvaraju, A. Das, R. Vedantam, M. Cogswell, D. Parikh, D. Batra, Grad-cam: Why did you say that? visual explanations from deep networks via gradient-based localization, CoRR abs/1610.02391. [arXiv: 1610.02391](https://arxiv.org/abs/1610.02391).
URL <http://arxiv.org/abs/1610.02391>
- [6] M. Chromik, M. Schuessler, A taxonomy for human subject evaluation of black-box explanations in xai., in: ExSS-ATEC@ IUI, 2020.
- [7] M. Reyes, R. Meier, S. Pereira, C. A. Silva, F.-M. Dahlweid, H. v. Tengg-Kobligk, R. M. Summers, R. Wiest, On the interpretability of artificial intelligence in radiology: Challenges and opportunities, Radiology: Artificial Intelligence 2 (3) (2020) e190043.
- [8] D. S. Weld, G. Bansal, The challenge of crafting intelligible intelligence, Communications of the ACM 62 (6) (2019) 70–79.
- [9] K. Weitz, D. Schiller, R. Schlagowski, T. Huber, E. André, “let me explain!”: exploring the potential of virtual agents in explainable ai interaction design, Journal on Multimodal User Interfaces (2020) 1–12.
- [10] M. T. Ribeiro, S. Singh, C. Guestrin, ”why should I trust you?”: Explaining the predictions of any classifier, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016, 2016, pp. 1135–1144.
- [11] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra, Grad-cam: Visual explanations from deep networks via gradient-based localization, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 618–626.
- [12] V. Petsiuk, A. Das, K. Saenko, Rise: Randomized input sampling for explanation of black-box models, in: Proceedings of the British Machine Vision Conference (BMVC), 2018.

- [13] S. M. Muddamsetty, N. S. J. Mohammad, T. B. Moeslund, Sidu: Similarity difference and uniqueness method for explainable ai, in: 2020 IEEE International Conference on Image Processing (ICIP), 2020, pp. 3269–3273. doi:10.1109/ICIP40778.2020.9190952.
- [14] F. Doshi-Velez, B. Kim, Considerations for evaluation and generalization in interpretable machine learning, in: Explainable and Interpretable Models in Computer Vision and Machine Learning, Springer, 2018, pp. 3–17.
- [15] A. Ignatiev, N. Narodytska, J. Marques-Silva, On relating explanations and adversarial examples, in: Advances in Neural Information Processing Systems, 2019, pp. 15883–15893.
- [16] A. Kuppa, S. Grzonkowski, M. R. Asghar, N.-A. Le-Khac, Black box attacks on deep anomaly detectors, in: Proceedings of the 14th International Conference on Availability, Reliability and Security, 2019, pp. 1–10.
- [17] V. Petsiuk, A. Das, K. Saenko, RISE: randomized input sampling for explanation of black-box models, CoRR abs/1806.07421. arXiv:1806.07421. URL <http://arxiv.org/abs/1806.07421>
- [18] G. Montavon, A. Binder, S. Lapuschkin, W. Samek, K.-R. Müller, Layer-wise relevance propagation: an overview, in: Explainable AI: interpreting, explaining and visualizing deep learning, Springer, 2019, pp. 193–209.
- [19] S. Bach, A. Binder, G. Montavon, F. Klauschen, K. Müller, W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, PloS one 10 (7).
- [20] K. Simonyan, A. Vedaldi, A. Zisserman, Deep inside convolutional networks: Visualising image classification models and saliency maps, arXiv preprint arXiv:1312.6034.
- [21] R. C. Fong, A. Vedaldi, Interpretable explanations of black boxes by meaningful perturbation, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 3429–3437.

- [22] R. Fong, M. Patrick, A. Vedaldi, Understanding deep networks via extremal perturbations and smooth masks, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 2950–2958.
- [23] A. Shrikumar, P. Greenside, A. Kundaje, Learning important features through propagating activation differences, arXiv preprint arXiv:1704.02685.
- [24] W. Samek, A. Binder, G. Montavon, S. Lapuschkin, K.-R. Müller, Evaluating the visualization of what a deep neural network has learned, IEEE transactions on neural networks and learning systems 28 (11) (2016) 2660–2673.
- [25] I. Goodfellow, J. Shlens, C. Szegedy, Explaining and harnessing adversarial examples, 2015, pp. 1–10.
- [26] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, A. Vladu, Towards deep learning models resistant to adversarial attacks.
- [27] S.-M. Moosavi-Dezfooli, A. Fawzi, P. Frossard, Deepfool: a simple and accurate method to fool deep neural networks, CVPR.
- [28] T. B. Brown, D. Mané, A. Roy, M. Abadi, J. Gilmer, Adversarial patch, CoRR abs/1712.09665. arXiv:1712.09665.
URL <http://arxiv.org/abs/1712.09665>
- [29] X. Yuan, P. He, Q. Zhu, X. Li, Adversarial examples: Attacks and defenses for deep learning, IEEE transactions on neural networks and learning systems 30 (9) (2019) 2805–2824.
- [30] Y. Bengio, A. Courville, P. Vincent, Representation learning: A review and new perspectives, IEEE transactions on pattern analysis and machine intelligence 35 (8) (2013) 1798–1828.
- [31] A. Mahendran, A. Vedaldi, Visualizing deep convolutional neural networks using natural pre-images, International Journal of Computer Vision 120 (3) (2016) 233–255.

- [32] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, L. Fei-Fei, ImageNet Large Scale Visual Recognition Challenge, *International Journal of Computer Vision (IJCV)* 115 (3) (2015) 211–252. doi:10.1007/s11263-015-0816-y.
- [33] S. M. Muddamsetty, T. B. Moeslund, Multi-level quality assessment of retinal fundus images using deep convolution neural networks, in: 16th International Joint Conference on Computer Vision Theory and Applications(VISAPP-2021), SCITEPRESS Digital Library, 2021.
- [34] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, *CoRR* abs/1512.03385. arXiv:1512.03385.
URL <http://arxiv.org/abs/1512.03385>
- [35] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.
- [36] T. Judd, K. Ehinger, F. Durand, A. Torralba, Learning to predict where humans look, in: 2009 IEEE 12th international conference on computer vision, IEEE, 2009, pp. 2106–2113.
- [37] M. Jiang, J. Xu, Q. Zhao, Saliency in crowd, in: *European Conference on Computer Vision*, Springer, 2014, pp. 17–32.
- [38] M. Jiang, S. Huang, J. Duan, Q. Zhao, Salicon: Saliency in context, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1072–1080.
- [39] R. McDonnell, M. Larkin, B. Hernández, I. Rudomin, C. O’Sullivan, Eye-catching crowds: saliency based selective variation, *ACM Transactions on Graphics (TOG)* 28 (3) (2009) 1–10.
- [40] A. Das, H. Agrawal, L. Zitnick, D. Parikh, D. Batra, Human attention in visual question answering: Do humans and deep networks look at the same regions?, *Computer Vision and Image Understanding* 163 (2017) 90–100.

- [41] T. Technology, User manual: Tobii x60 and x120 eye trackers (2008).
- [42] Z. Bylinskii, T. Judd, A. Oliva, A. Torralba, F. Durand, What do different evaluation metrics tell us about saliency models?, *IEEE transactions on pattern analysis and machine intelligence* 41 (3) (2018) 740–757.
- [43] W. Daniel, *Applied Nonparametric Statistics*, Duxbury advanced series in statistics and decision sciences, PWS-KENT Pub., 1990.
URL <https://books.google.dk/books?id=0hPvAAAAMAAJ>