

Distance-Aware eXplanation Based Learning

1st Misgina Tsighe Hagos* 2nd Niamh Belton[†] 3rd Kathleen M. Curran[‡] 4th Brian Mac Namee[§]
School of Computer Science School of Medicine School of Medicine School of Computer Science

Science Foundation Ireland Centre for Research Training in Machine Learning

University College Dublin

Dublin, Ireland

Email: *misgina.hagos@ucdconnect.ie, [†]niamh.belton@ucdconnect.ie, [‡]kathleen.curran@ucd.ie, [§]brian.macnamee@ucd.ie

Abstract—eXplanation Based Learning (XBL) is an interactive learning approach that provides a transparent method of training deep learning models by interacting with their explanations. XBL augments loss functions to penalize a model based on deviation of its explanations from user annotation of image features. The literature on XBL mostly depends on the intersection of visual model explanations and image feature annotations. We present a method to add a distance-aware explanation loss to categorical losses that trains a learner to focus on important regions of a training dataset. Distance is an appropriate approach for calculating explanation loss since visual model explanations such as Gradient-weighted Class Activation Mapping (Grad-CAMs) are not strictly bounded as annotations and their intersections may not provide complete information on the deviation of a model’s focus from relevant image regions. In addition to assessing our model using existing metrics, we propose an interpretability metric for evaluating visual feature-attribution based model explanations that is more informative of the model’s performance than existing metrics. We demonstrate performance of our proposed method on three image classification tasks.

Index Terms—eXplanation Based Learning, Interactive Machine Learning, eXplainable AI

I. INTRODUCTION

Research on model transparency in deep learning is dominated by studies on dataset bias [1], model interpretability, and explainability [2]. Another field of study that aims to improve model transparency, Interactive Machine Learning (IML), hits two birds with one stone [3], [4]. First, it provides transparency through engagement by allowing user interaction in the model training process. Second, it improves model performance by collecting expert knowledge directly from users. IML usually considers users as *dumb partners* with the sole responsibility of categorizing training instances into one of a set of pre-selected categories as opposed to *clever partners* who can clarify their feedback in addition to categorizing instances. However, advances in model explanation research opens the door for a more detailed and richer interaction between models and users during training.

A. eXplanation based learning

While model explanation methods have been proposed and continue to be used to tackle the “black-box” nature of deep learning models [5], [6], they can also be used in an interactive learning approach to promote a more transparent model training process [7], [8]. This is known as eXplanation

Based Learning (XBL)¹, which collects user feedback on model explanations and uses the feedback to train, debug, or refine a trained model.

In applications such as medical image classifications, deep learning models have been observed to focus on non-relevant or confounding parts of medical images such as artifacts for their classification or prediction outputs [11], [12]. In addition to promoting transparent learning process, XBL has the potential to unlearn such wrong correlations, which are termed as confounding regions, confounders, or spurious correlations (used interchangeably in this paper) [12], [13]; confounding regions are parts of training instances that are not correlated with a category, but incorrectly assumed to be so by a learner.

As is displayed in Fig. 1, XBL is generally made up of four steps. The first is traditional model training which often uses a categorical loss. The next step is generating model explanations. Feature attribution based local explanations [9] or surrogate model explanations [10] can be used for this. We limit the scope of this work to a saliency based local explanation, Gradient-weighted Class Activation Mapping (Grad-CAM) [14]. In the third step, explanations are presented to users and feedback is collected. For method development and experiment purposes confounding regions can be added to a dataset and their masks used as user feedback for XBL. Finally, the collected feedback is used to calculate an explanation loss, which in turn is used to augment the initial categorical loss, and refine the original model using XBL training [7].

The training process in XBL augments loss functions to include an explanation loss, which can be based on either or both of: (1) a model’s deviation from user annotated feedback that shows objects of interest; and (2) a model’s focus on user annotation of non-salient or confounding image regions. This explanation loss is usually based on the intersection of the user annotation of image features and a model’s visual explanation. Loss functions are generally augmented as follows:

¹Different terms such as *explanatory debugging* [8], *explanatory interactive learning* [9], *explanatory guided learning* [10] are used in the literature. We choose to use the term *eXplanation Based Learning* because we believe it generalizes all of them.

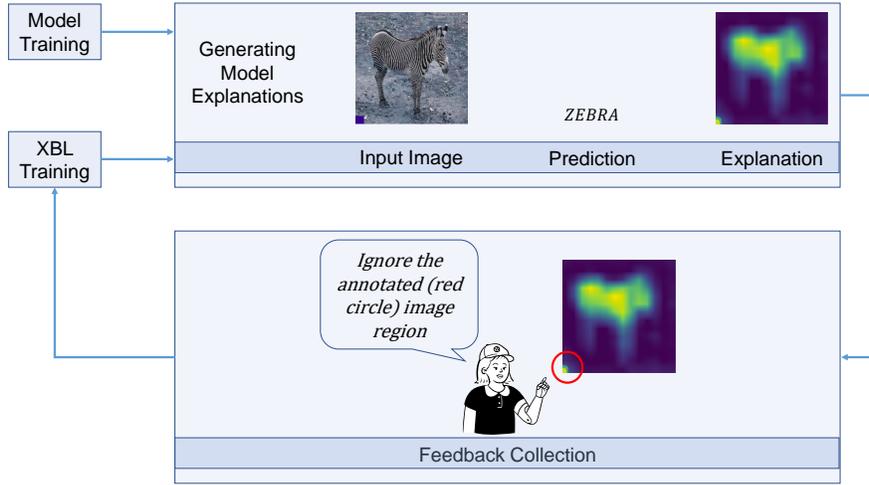


Fig. 1. The eXplanation Based Learning (XBL) loop. The user feedback, which is expected to be an annotation mask of the confounding image region in the lower left corner (highlighted by the saliency map), is portrayed here as a red circle for easier visualization.

$$L_{expl} = \sum_{i=1}^N e(expl_{i,c}, M_{i,c}) \quad (1)$$

$$L_{CE} = - \sum_{i=1}^N e(\hat{y}_i, Y_i) \quad (2)$$

$$L = L_{expl} + L_{CE} + \lambda \sum_{i=1} \theta_i \quad (3)$$

The term, L_{expl} in (1) is the explanation loss calculated as the error, e , between the model’s explanation, $expl_{i,c}$ for input i with category c , and the ground truth annotation, $M_{i,c}$, where $M = 1$ for relevant regions and $M = 0$ for non-salient regions. The term, L_{CE} , in (2) is the traditional cross entropy loss which is calculated based on the error, e , between the model’s prediction \hat{y}_i and ground-truth label Y_i for instance i . While Y_i only holds category label, $M_{i,c}$ holds a mask annotation of relevant objects in an input i . Finally, XBL consists of the sum of L_{expl} , L_{CE} and a regularization term, λ where θ_i is the network parameters.

Most XBL loss function augmentations in the literature fail to consider two scenarios: (1) focus of a model’s attention may get closer to- and gradually shift to the relevant regions of training instances; for this reason, we need to penalize the learner less as the explanations (the model’s attention) starts to improve. This means there is a need to make loss functions positively related to the distance of a model’s wrong attention from the relevant regions; and (2) model activations that usually make up model explanations are not as strictly bounded as user annotations and we need to relax the training penalization as we get closer to the relevant parts of training images. In order to address these shortcomings of existing XBL methods and assuming model explanations correctly highlight the reasoning behind a model’s output, in this paper, we address the following research question: “Can we augment XBL loss functions in way that is sensitive to distances between

explanations and user annotations of relevant image regions for better classification and explanation performance?”

Another aspect of XBL that is often overlooked is using coefficients that weigh and balance impact of explanation losses and classification losses and optimizing them. We also consider these coefficients as hyper-parameters and tune them to find their optimal values before starting model training with XBL.

B. Evaluation of model explanations

While subjective evaluations of explanations that involve humans would give a user-centric assessment of model generated explanations [15], objective evaluations are often used for a speedy assessment and comparison in the development of model explainability methods [16]. Most of the existing evaluation methods give weight to the generated explanations over the ground truth feature annotations. This can result in over-confident evaluation results. In addition to using an existing evaluation method, to address this issue we propose an interpretability metric that assesses how much of the ground truth feature annotation has been identified as relevant by model explanations. We restrict our work in this paper to objective evaluations.

The main contributions of this paper are:

- 1) Decoyed versions of image classification datasets are created for XBL experiments.
- 2) A new XBL method, Distance-Aware eXplanation Based Learning (XBL-D), is proposed and evaluated.
- 3) A saliency map explanation interpretability metric, Activation Recall, is proposed and demonstrated.
- 4) Our experiments demonstrate that incorporating distance-aware learning into XBL performs better than baseline algorithms in classification tasks, and generates more accurate explanations. Furthermore, Code and links to download the datasets are shared online².

²<https://github.com/Msgun/XBL-D>

II. RELATED WORK

In this section, we present a review of relevant literature on XBL and model explanation evaluation metrics.

A. *eXplanation based learning*

XBL methods can be generally categorized into two categories: (1) augmenting loss functions with explanation losses; and (2) using user feedback to augment training datasets by removing confounding or spurious regions identified by users.

1) *Augmenting loss functions*: The model explanation method used has a huge impact on an interactive learning process, not only because it is directly used to compute explanation loss ($expl_{i,c}$ as in (1)), but also because it can impact user experience and feedback quality. Right for the Right Reasons (RRR) [17] penalises a model with high input gradient model explanations on the wrong image regions annotated by a user. RRR uses

$$L_{expl} = \sum_n^N (M_n \frac{\partial}{\partial x_n} (\sum_{k=1}^K \log \hat{y}_{nk}))^2 \quad (4)$$

for a function $f(X|\theta) = \hat{y} \in R^{N \times K}$ trained on images x_n of size N with K categories, where $M_n \in \{0, 1\}$ is user annotation of image regions that should be avoided by the model.

A Grad-CAM model explanation was used instead of input gradients in RRR-G by Schramowski *et al.* [18] using the following loss function:

$$L_{expl} = \sum_n^N M_n GradCAM(x_n) \quad (5)$$

Similarly, Right for Better Reasons (RBR) [19] uses Influence Functions (IF) in place of input gradients to correct a model’s behavior. Contextual Decomposition Explanation Penalization (CDEP) [20] penalizes features and feature interactions.

User feedback in XBL experiments can be one or both of: (1) telling the model to ignore non-salient image regions; and (2) instructing the model to focus on important image regions in a training dataset [21]. While the XBL methods presented above refine a model by using the first feedback type, Human Importance-aware Network Tuning (HINT) does the opposite by teaching a model to focus on important image parts using Grad-CAM model explanations [22].

Most of the literature on XBL focuses on using feature attribution based saliency maps such as input gradients and Grad-CAMs as model explanations. Prototype based explanations have also been utilized in Bontempelli *et al.* [23] to debug Part-Prototype networks at concept level.

2) *Augmenting training dataset*: Instead of augmenting loss functions, XBL can be implemented by relabeling, augmenting existing instances, or adding new training instances based on user feedback. Instance relabeling has been deployed to clean label noise in a training dataset that is identified using example

based explanations [24]. Counter-Examples (CE), which are variants of training instances with added modifications using user feedback can be generated to augment dataset for model re-training [9]. Simpler surrogate models have also been used as global explanations to elicit feedback in the form of new training instances [10].

B. *Evaluating feature attribution based explanations*

Feature attribution based explanations can be evaluated intrinsically and/or extrinsically [25]. Intrinsic evaluation involves only the model and the generated explanations themselves [26], while extrinsic evaluation involves subjective human evaluation [27] or objective usage of ground-truth annotation data.

Objective evaluation of model explanations provides an easier and quicker way of assessing interpretability by comparing model explanations to ground-truth data. Overlap of visual explanations and feature annotations can be used to compute localization ability of a model’s explanations; to avoid explanations with high false positive rates which cover wide area of an image, thereby scoring a high overlap with annotations, penalized versions of overlap: Penalized Localization Accuracy (PLA) was proposed [28]. Activation Precision (AP) is another approach that computes how many of the pixels predicted as relevant by a model are actually relevant [29]. AP is presented in (6), where A_{obj_n} is a mask of relevant image regions in input image x_n and T_r is a threshold function that finds the (100- r) percentile and sets elements of the explanation, $expl_\theta$, below this value to zero and the remaining elements to one. AP usually requires a low r value or high threshold so we can avoid explanations with high false positive rates.

$$AP = \frac{1}{N} \sum_n^N \frac{T_r(expl_\theta(x_n)) * A_{obj_n}}{T_r(expl_\theta(x_n))} \quad (6)$$

There is a trade-off between selecting a higher threshold and accurately assessing model explanations. While increasing the threshold would mean focusing on smaller areas of an explanation and avoiding high false positive rates, it also means parts of an explanation would be masked before they are assessed, which could result in overconfidence in model explanations.

III. DISTANCE-AWARE EXPLANATION BASED LEARNING

We view the training images as instances made up of three parts: (1) the relevant regions, masked by A_{obj} , that are considered important for category classification; (2) the confounding regions, masked by annotation A_{con} , which are not correlated with any category but can trick the learner into learning that they are; and (3) the remaining image parts that are usually easily ignored by a learner as background image regions.

Our explanation loss penalizes a learner based on the amount of wrong attention it gives to A_{con} , with due consideration of this wrong attention’s distance from A_{obj} . For

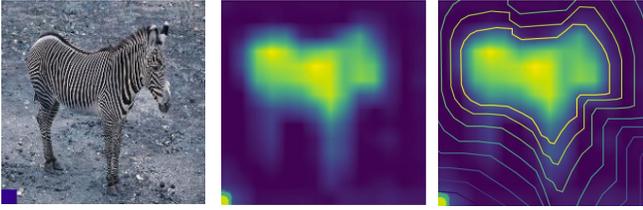


Fig. 2. [Best viewed in color.] Illustration of a distance-aware explanation loss calculation for an input image (left), Grad-CAM (middle). Distance is represented using Viridis color-map in the right figure. Yellow is for the smallest distance and dark purple for the largest. In this case, the confounding region (in the lower left image region) that is wrongly found relevant is as far as it can be from the important region. Pixel intensity of Grad-CAM on the confounding region is exaggerated for presentation purposes.

example, in Fig. 2, a Grad-CAM explanation shows a model giving attention to a confounder located on the lower left corner of an input image. Distance of the attention to the confounder is illustrated with a Viridis color-map showing largest distances as dark purple. In this case, this would result in the highest penalty. As the model’s (wrong) focus starts to get closer to A_{obj} , the explanation loss would decrease. We used Grad-CAM because it was found to be more sensitive to training label reshuffling and model parameter randomization [26] than other saliency based explanations.

Equations (7) and (8) underpin how we propose to integrate explanation and classification losses. Algorithm 1 shows how this combined loss function is integrated into the overall XBL-D approach. Here, G_n is the center of gravity of objects of interest in input images that are masked with A_{obj} , $expl_\theta(x_n)$ is a Grad-CAM explanation of input x_n to model F , and A_{con} is the annotation of a confounding region in x_n . A model’s incorrect focus on a confounding region is detected using the intersection $I_\theta(x_n)$ between $expl_\theta(x_n)$ and A_{con} . The distance between a model’s wrong attention to a confounding region and center of A_{obj} or G_n is then approximated by calculating average of the minimum and maximum euclidean distances, d , between points in $I_\theta(x_n)$ and G_n . This gives us a measure of how far a model’s incorrect attention is from the relevant image regions. In (8), L_{CE} represents the cross entropy loss and $\lambda \sum_{i=1} \theta_i$ is a weight (θ) regularization term.

$$L_{expl} = \sum_n^N d(G_n, I_\theta(x_n)) \quad (7)$$

$$L = \lambda_1 L_{CE} + \lambda_2 L_{expl} + \lambda \sum_{i=1} \theta_i \quad (8)$$

A. Activation recall

In addition to using AP, we propose Activation Recall (AR) to assess visual explanations, such as Grad-CAMs, generated by a trained model. AR measures how much of the relevant parts of test images are considered relevant by a model. This is presented in 9, where (similarly to AP) T_r is a threshold function that finds the (100- r) percentile and sets elements of

Algorithm 1 Distance-aware eXplanation Based Learning (XBL-D)

Input: confounded training dataset \hat{X} and ground-truth category Y , feature annotation of object(s) of interest in \hat{X} : A_{obj} , feature annotation of confounders in \hat{X} : A_{con} .

Parameters: classification loss coefficient: λ_1 , explanation loss coefficient: λ_2 , regularization term: λ , network parameters: θ

Output: refined function F

- 1: $F \leftarrow$ Fit function using \hat{X}
 - 2: **repeat**
 - 3: $G \leftarrow$ center of gravity of objects of interest in A_{obj} .
 - 4: $expl_\theta \leftarrow$ saliency map explanations of \hat{X} generated using Grad-CAM.
 - 5: $I_\theta \leftarrow$ set of intersections between $expl_\theta$ and A_{con}
 - 6: $L_{expl} \leftarrow$ explanation loss as average of the minimum and maximum euclidean distances between points in I_θ and G
 - 7: $L_{CE} \leftarrow$ classification loss between Y and $F(\hat{X})$
 - 8: Total loss, $L \leftarrow \lambda_1 * L_{CE} + \lambda_2 * L_{expl} + \lambda \sum_{i=1} \theta_i$
 - 9: update F using L
 - 10: **until** $L \leq \sigma$, where σ is a tolerable total loss
 - 11: **return** F
-



Fig. 3. Sample images from MS COCO with confounding regions added to random corners and their corresponding object masks

$expl_\theta(x_n)$ below this value to zero and the remaining elements to one.

$$AR = \frac{1}{N} \sum_n^N \frac{T_r(expl_\theta(x_n)) * A_{obj_n}}{A_{obj_n}} \quad (9)$$

Instead of selecting a single threshold to assess generated explanations, we compute AP and AR at different thresholds to show impacts of choosing threshold on the evaluation metrics. This also gives us an insight into how explanation evaluation can be misleading without the full information, i.e thresholding.

IV. EXPERIMENTS

In this section, we describe the datasets, model architectures, and training details used in our experiments to evaluate the performance of XBL-D.

A. Dataset

In order to validate performance of XBL-D, locations of the confounding regions needs to be known beforehand. For this reason, we used a publicly available decoyed dataset and

created two new decoyed versions of existing datasets for our experiments:

- 1) Decoy Fashion MNIST³. This was created by Teso and Kersting [9]. 4x4 pixel confounders with random pixel intensities were added to random corners of images from the Fashion MNIST training dataset [30]. The 10,000 images from the test dataset were left clean.
- 2) Decoy CIFAR-10⁴. We created this dataset by adding 4x4 pixel confounders with random pixel intensities to random corners of the training set of CIFAR-10 dataset. The CIFAR-10 dataset contains a training set of 50000 and test set of 10000 32x32 RGB images categorized into 10 classes. Similar to the Decoy Fashion MNIST, the test set of this dataset was also left clean for evaluation purposes.
- 3) Decoyed subset of MS-COCO. We extracted a total of 2000 images for training and 600 image for testing, from the *Train* and *Zebra* categories of the MS-COCO dataset [31]. We then added 16x16 confounding pixels with random pixel intensities to random corners of the training images, which are of size 224x224. Images in the test set were left clean. We selected the *Train* and *Zebra* categories based on the low intersection of objects from both categories. Sample images from this dataset are shown in Fig. 3. We refer to this dataset as Decoy MS-COCO₍₂₎.

B. Architecture selection and training

We performed all of our experiments using Tensorflow and Keras⁵. For all our datasets, we searched for the best model architectures and hyper-parameters using HyperBand algorithm [32] in Keras tuner⁶. We considered and optimized the hyper-parameters: number and size of convolutional layers, number of pooling layers, number and size of fully connected layers, and learning rate. A Convolutional Neural Network (CNN) with one convolutional layer containing 160 filters and two fully connected consecutive layers of sizes 992 and 800 nodes, and a learning rate = 1.158e-04 was found to perform best for the Decoy Fashion MNIST dataset. For the Decoy CIFAR-10, a CNN with two convolutional layers of filters 250 and 300 followed by one fully connected layer with 912 nodes, and a learning rate = 1.267e-04 was selected. Similarly, we found that a CNN with four convolutional layers (containing 160, 352, 416, and 224 consecutive filters) each followed by a max-pooling layer, one fully connected layer of size 480 nodes, and a learning rate = 1.789e-05 performed best for the Decoy MS-COCO₍₂₎ dataset.

To start with, the selected model architectures are fitted on the corresponding dataset using categorical cross-entropy loss and the Adam optimizer. We refer to the resulting models as *Unrefined*. All the models are then refined using XBL-D.

For the Decoy MS-COCO₍₂₎ dataset, we run 20 epochs of refinement where each epoch took an average of 15 minutes, while for each of the Decoy CIFAR-10 and Decoy Fashion MNIST datasets, we run 50 epochs of refinement each taking averages of 7 and 5 minutes, respectively. Model training was performed on a machine with NVIDIA RTX A5000 graphics card.

Before starting the model refinement using XBL-D, we searched for optimal values of the coefficients of the categorical cross entropy loss (λ_1) and explanation loss (λ_2) using HyperBand in Keras tuner and we ended up with $\lambda_1 = 2.7$ and $\lambda_2 = 0.1$. We searched all hyper-parameters for each of the datasets separately. However, since λ_1 and λ_2 influence how XBL-D works, we decided to find one set that should work for the other domains for domain transferability purposes. Hence, the hyper-parameter search of λ_1 and λ_2 was performed on the most challenging task among the 3 datasets, which is the decoy MS-COCO₍₂₎ that contains large RGB images.

V. RESULTS

In this section, we present classification and explanation performance results of our proposed method and compare them against baseline methods.

A. Classification

Table I presents classification accuracy performance of XBL-D and comparison against baseline methods. On the original test set of Fashion MNIST dataset, our proposed method achieves classification performance of 0.904 surpassing previous XBL methods [33]. The second best performing model was RRR with a classification accuracy of 0.894. None of the available baseline methods were implemented for our Decoy CIFAR-10 and Decoy MS-COCO₍₂₎. For this reason, we trained a model using the best performing method, RRR, on the decoyed CIFAR-10 and MS-COCO₍₂₎ datasets for comparison purposes. Again, compared to RRR and Unrefined models, XBL-D achieved superior classification performance on the original test sets of CIFAR-10 and MS-COCO₍₂₎ achieving accuracies of 0.843 and 0.938, respectively, as is summarized in Table I.

TABLE I
CLASSIFICATION ACCURACY COMPARISONS ON ORIGINAL TEST IMAGES OF FASHION MNIST, CIFAR-10, AND MS-COCO₍₂₎.

Method	Decoy Fashion MNIST	Decoy CIFAR-10	Decoy MS-COCO ₍₂₎
Unrefined	0.862	0.789	0.845
XBL-D	0.904	0.843	0.938
RRR	0.894	0.810	0.853
RRR-G	0.786	-	-
RBR	0.876	-	-
CDEP	0.767	-	-
HINT	0.582	-	-
CE	0.858	-	-

³We collected this dataset at <https://codeocean.com/capsule/7818629/tree/v1>

⁴https://osf.io/w5f7y/?view_only=abb7f5f5bfc48fb8c891838f699c0d3

⁵https://www.tensorflow.org/api_docs/python/tf/keras

⁶https://keras.io/keras_tuner/

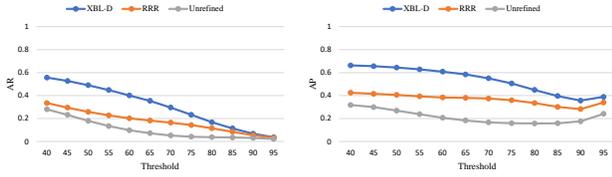


Fig. 4. AR and AP evaluations of explanations generated for a clean Fashion MNIST test dataset using a model trained on the Decoy Fashion MNIST. The evaluations are performed at threshold values ranging from 40% to 95% with step size = 5

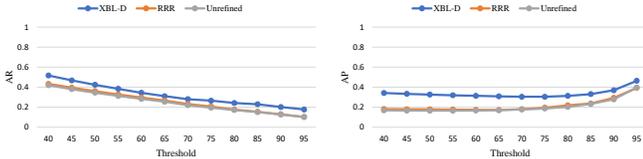


Fig. 5. AR and AP evaluations of explanations generated for a clean CIFAR-10 test dataset using a model trained on the Decoy CIFAR-10. The evaluations are performed at threshold values ranging from 40% to 95% with step size = 5

B. Explanation performance

While AR and AP evaluations, on the original test sets of the Fashion MNIST, MS-COCO₍₂₎ and CIFAR-10, across thresholds ranging from 40% to 95% with step size = 5 are presented in Figures 4, 5, and 6, Table II presents a summary of evaluations of explanations.

1) *Fashion MNIST*: Our proposed method scores higher than both RRR and Unrefined models using both metrics. At threshold = 40%, XBL-D scored highest values of AR = 0.557 and AP = 0.663 (Table II). Given that higher threshold means considering smaller areas of Grad-CAM, AR values decrease with increasing threshold (see Fig. 4). However, even though AP seemed to decrease with increasing threshold values, it starts to increase at threshold above 90% (we accredit this to the Gray-Scale nature of the decoy Fashion MNIST dataset).

2) *CIFAR-10*: Similar to the Fashion MNIST, our proposed method performs higher than both RRR and Unrefined models using both metrics. At threshold = 40%, XBL-D scored highest values of AR = 0.516 (Table II) and at threshold = 95%, XBL-D scored AP = 0.342, outperforming both methods. While AR naturally decreases with increasing threshold, AP increases

TABLE II
SUMMARY EVALUATIONS OF EXPLANATIONS GENERATED FOR THE ORIGINAL FASHION MNIST, MS-COCO AND CIFAR-10 TEST DATASETS.

Metric	Method	Decoy Fashion MNIST	Decoy CIFAR-10	Decoy MS-COCO ₍₂₎
AR	Unrefined	0.280	0.419	0.500
	XBL-D	0.557	0.516	0.860
	RRR	0.335	0.432	0.841
AP	Unrefined	0.318	0.168	0.609
	XBL-D	0.663	0.342	0.698
	RRR	0.425	0.181	0.761

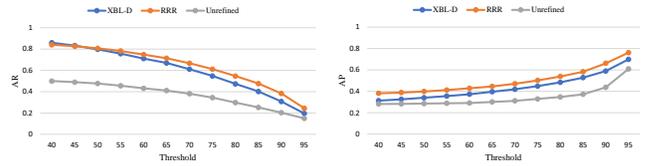


Fig. 6. AR and AP evaluations of explanations generated for a clean MS-COCO test dataset using a model trained on the Decoy MS COCO₍₂₎. The evaluations are performed at threshold values ranging from 40% to 95% with step size = 5

given the RGB nature of CIFAR-10.

3) *MS-COCO₍₂₎*: Our method scored better AR at lower thresholds and performed comparable to RRR at other thresholds (at threshold = 40%, XBL-D scored AR = 0.860, Table II). Similar to the other datasets, we also found that low threshold values led to higher AR values (see Fig. 6). However, unlike the models trained on the Fashion MNIST dataset but similar to the CIFAR-10, AP values increase with increasing threshold (at threshold=95%, RRR scored highest AP of 0.761, Table II). We accredit this to the RGB nature of the MS-COCO₍₂₎ dataset.

Sample Grad-CAM outputs of input images from both categories are displayed in Fig. 7. We show sample explanation outputs for the MS-COCO₍₂₎ images because their high resolution makes them well suited for presentation. While the clean test sets were used in computing AR and AP explanation evaluations, sample of the decoyed images from training set of the MS-COCO₍₂₎ are shown in Fig. 7 to demonstrate the ability of XBL-D in avoiding confounding regions and to compare it against RRR and the Unrefined model. As is displayed in the sample outputs, our proposed method was able to produce accurate explanations that focus on relevant parts of objects in input images and successfully ignores confounders.

VI. DISCUSSION

In addition to explaining a model’s classification output, XBL facilitates a more transparent machine learning process by providing a rich user interaction mechanism. As opposed to the traditional interactive machine learning that is usually performed through instance category labeling, a user would be able to get involved at a deeper level by interacting with model explanations in the machine learning process. In XBL, a user would be able to teach a learner model by observing and commenting on the reasoning (i.e correcting model explanations) behind its predictions. This kind of user engagement has the potential to circumvent the *black-box* public image of deep learning models since it aims to build a rapport with users by providing a transparent way of interaction with an opportunity to refine the models.

When compared against baseline methods, XBL-D achieved superior performance in classifying all three datasets. We believe this is because it unlearns confounding regions, which were wrongly found relevant by a model, based on their locations and distances from the user annotated relevant regions. As shown in the sample outputs in Fig. 7, a model’s focus,

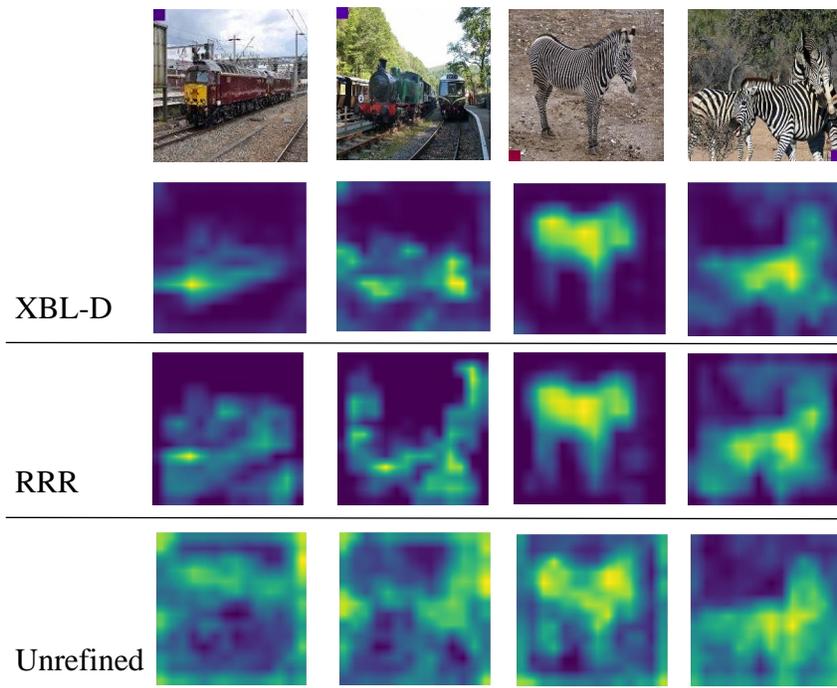


Fig. 7. Sample Grad-CAM Outputs. Original size of all Grad-CAM images was 14x14; They are up-sampled to 224x224 for easier comparison against input images.

shown with visual explanations is not strictly bounded and, however good it is, there is always a good chance it might exceed boundaries of relevant region(s). Based on this fact, XBL-D instructs a learner that it is not only acceptable to focus on the user annotated parts but also around it as long as it keeps a distance from the confounding region. Had the explanation loss been based on intersection of generated explanations with the confounding regions, it would penalize the model whenever it focuses on the confounders without consideration for the confounders’ locations.

In addition to XBL-D, we observe that the Unrefined model performed better than most of the other XBL models in classifying the decoy Fashion MNIST. We attribute this to the accuracy-interpretability trade-off in deep learning. Although the existence of this trade-off is debated [34], [35], deep learning models that are refined with an explanation based learning could lose performance if the refinement is not performed using a fitting approach such as our proposed method, XBL-D.

We also proposed an interpretability metric, Activation Recall (AR). AR measures how much of the user annotated relevant image regions were actually considered relevant by a trained model. It circumvents a possible over-confidence that may result from mainly focusing on explanations (saliency maps in this case) during explanation evaluation. By redirecting the focus from explanations to ground-truth annotations, AR provides a reliable metric for explanation evaluation. We recommend AR should be used in conjunction with AP for a reliable assessment of model explanations.

Objective evaluations of generated explanations of test

images of employed datasets across different thresholds also show that XBL-D performs better than RRR and Unrefined models in generating accurate explanations. Threshold selection is important in computing AR and AP of generated explanations. In all the datasets, we observed that low threshold values lead to higher AR while the opposite is true for AP. This is because parts of Grad-CAM considered for AR calculation increase with decreasing threshold. We also note that the Gray-Scale nature of Fashion MNIST affects AP values and it plummets with increasing threshold, but recovers after threshold = 90%.

In addition to performing better at objective evaluations, XBL-D also outputs visually accurate saliency maps compared to the RRR and Unrefined models as can be seen in Fig. 7. We were able to observe that XBL-D is better than RRR and the Unrefined models at localizing objects of interest in input images.

VII. CONCLUSION

In this paper we proposed and demonstrated superior performance of XBL-D, a distance-aware explanation loss for XBL loss function augmentation. This introduces a new direction for XBL research: the consideration of the distance of a model’s wrong attention from relevant regions. XBL-D was able to achieve superior classification and interpretability performance compared to baseline methods on three different datasets. This assures that our proposed method generalizes across different datasets.

ACKNOWLEDGMENT

This publication has emanated from research conducted with the financial support of Science Foundation Ireland under Grant number 18/CRT/6183. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

REFERENCES

- [1] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. Daumé, and K. Crawford, "Excerpt from datasheets for datasets," in *Ethics of Data and Analytics*. Auerbach Publications, 2022, pp. 148–156.
- [2] R. Goebel, A. Chander, K. Holzinger, F. Lecue, Z. Akata, S. Stumpf, P. Kieseberg, and A. Holzinger, "Explainable ai: the new 42?" in *International Cross-domain Conference for Machine Learning and Knowledge Extraction*. Springer, 2018, pp. 295–303.
- [3] J. A. Fails and D. R. Olsen Jr, "Interactive machine learning," in *Proceedings of the 8th International Conference on Intelligent User Interfaces*, 2003, pp. 39–45.
- [4] R. Fiebrink, P. R. Cook, and D. Trueman, "Human model evaluation in interactive supervised learning," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 2011, pp. 147–156.
- [5] A. Singh, S. Sengupta, and V. Lakshminarayanan, "Explainable deep learning models in medical image analysis," *Journal of Imaging*, vol. 6, no. 6, p. 52, 2020.
- [6] B. H. Van der Velden, H. J. Kuijff, K. G. Gilhuijs, and M. A. Viergever, "Explainable artificial intelligence (xai) in deep learning-based medical image analysis," *Medical Image Analysis*, vol. 79, p. 102470, 2022.
- [7] S. Stumpf, V. Rajaram, L. Li, W.-K. Wong, M. Burnett, T. Dietterich, E. Sullivan, and J. Herlocker, "Interacting meaningfully with machine learning systems: Three experiments," *International Journal of Human-Computer Studies*, vol. 67, no. 8, pp. 639–662, 2009.
- [8] T. Kulesza, M. Burnett, W.-K. Wong, and S. Stumpf, "Principles of explanatory debugging to personalize interactive machine learning," in *Proceedings of the 20th International Conference on Intelligent User Interfaces*, 2015, pp. 126–137.
- [9] S. Teso and K. Kersting, "Explanatory interactive machine learning," in *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 2019, pp. 239–245.
- [10] T. Popordanoska, M. Kumar, and S. Teso, "Machine guides, human supervises: Interactive learning with global explanations," *arXiv preprint arXiv:2009.09723*, 2020.
- [11] A. J. DeGrave, J. D. Janizek, and S.-I. Lee, "Ai for radiographic covid-19 detection selects shortcuts over signal," *Nature Machine Intelligence*, vol. 3, no. 7, pp. 610–619, 2021.
- [12] N. Pfeuffer, L. Baum, W. Stammer, B. M. Abdel-Karim, P. Schramowski, A. M. Bucher, C. Hügel, G. Rohde, K. Kersting, and O. Hinz, "Explanatory interactive machine learning," *Business & Information Systems Engineering*, pp. 1–25, 2023.
- [13] M. T. Hagos, K. M. Curran, and B. Mac Namee, "Identifying spurious correlations and correcting them with an explanation-based learning," *arXiv preprint arXiv:2211.08285*, 2022.
- [14] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and 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.
- [15] N. Halliwell and F. Lecue, "Trustworthy convolutional neural networks: A gradient penalized-based approach," *arXiv preprint arXiv:2009.14260*, 2020.
- [16] P. Dabkowski and Y. Gal, "Real time image saliency for black box classifiers," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [17] A. S. Ross, M. C. Hughes, and F. Doshi-Velez, "Right for the right reasons: training differentiable models by constraining their explanations," in *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, 2017, pp. 2662–2670.
- [18] P. Schramowski, W. Stammer, S. Teso, A. Brugger, F. Herbert, X. Shao, H.-G. Luigs, A.-K. Mahlein, and K. Kersting, "Making deep neural networks right for the right scientific reasons by interacting with their explanations," *Nature Machine Intelligence*, vol. 2, no. 8, pp. 476–486, 2020.
- [19] X. Shao, A. Skryagin, W. Stammer, P. Schramowski, and K. Kersting, "Right for better reasons: Training differentiable models by constraining their influence functions," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 11, 2021, pp. 9533–9540.
- [20] L. Rieger, C. Singh, W. Murdoch, and B. Yu, "Interpretations are useful: penalizing explanations to align neural networks with prior knowledge," in *International Conference on Machine Learning*. PMLR, 2020, pp. 8116–8126.
- [21] M. T. Hagos, K. M. Curran, and B. Mac Namee, "Impact of feedback type on explanatory interactive learning," in *International Symposium on Methodologies for Intelligent Systems*. Springer, 2022, pp. 127–137.
- [22] R. R. Selvaraju, S. Lee, Y. Shen, H. Jin, S. Ghosh, L. Heck, D. Batra, and D. Parikh, "Taking a hint: Leveraging explanations to make vision and language models more grounded," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2591–2600.
- [23] A. Bontempelli, S. Teso, K. Tentori, F. Giunchiglia, A. Passerini *et al.*, "Concept-level debugging of part-prototype networks," in *Proceedings of the The Eleventh International Conference on Learning Representations (ICLR 23)*. ICLR 2023, 2023.
- [24] S. Teso, A. Bontempelli, F. Giunchiglia, and A. Passerini, "Interactive label cleaning with example-based explanations," *Advances in Neural Information Processing Systems*, vol. 34, pp. 12966–12977, 2021.
- [25] A. Gupta, N. Saunshi, D. Yu, K. Lyu, and S. Arora, "New definitions and evaluations for saliency methods: Staying intrinsic, complete and sound," *Advances in Neural Information Processing Systems*, vol. 35, pp. 33 120–33 133, 2022.
- [26] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim, "Sanity checks for saliency maps," *Advances in neural information processing systems*, vol. 31, 2018.
- [27] A. Holzinger, A. Carrington, and H. Müller, "Measuring the quality of explanations: the system causability scale (scs)," *KI-Künstliche Intelligenz*, pp. 1–6, 2020.
- [28] N. Belton, I. Welaratne, A. Dahlan, R. T. Hearne, M. T. Hagos, A. Lawlor, and K. M. Curran, "Optimising knee injury detection with spatial attention and validating localisation ability," in *Annual Conference on Medical Image Understanding and Analysis*. Springer, 2021, pp. 71–86.
- [29] A. J. Barnett, F. R. Schwartz, C. Tao, C. Chen, Y. Ren, J. Y. Lo, and C. Rudin, "A case-based interpretable deep learning model for classification of mass lesions in digital mammography," *Nature Machine Intelligence*, vol. 3, no. 12, pp. 1061–1070, 2021.
- [30] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [31] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European Conference on Computer Vision*. Springer, 2014, pp. 740–755.
- [32] L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar, "Hyperband: A novel bandit-based approach to hyperparameter optimization," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 6765–6816, 2017.
- [33] F. Friedrich, W. Stammer, P. Schramowski, and K. Kersting, "A typology for exploring the mitigation of shortcut behaviour," *Nature Machine Intelligence*, pp. 1–12, 2023.
- [34] C. Rudin, "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead," *Nature Machine Intelligence*, vol. 1, no. 5, pp. 206–215, 2019.
- [35] G. K. Dziugaite, S. Ben-David, and D. M. Roy, "Enforcing interpretability and its statistical impacts: Trade-offs between accuracy and interpretability," *arXiv preprint arXiv:2010.13764*, 2020.