

Deep Integrated Explanations

Oren Barkan*
The Open University

Yehonatan Elisha*
The Open University

Jonathan Weill
Tel Aviv University

Yuval Asher
Tel Aviv University

Amit Eshel
Tel Aviv University

Noam Koenigstein
Tel Aviv University

ABSTRACT

This paper presents Deep Integrated Explanations (DIX) - a universal method for explaining vision models. DIX generates explanation maps by integrating information from the intermediate representations of the model, coupled with their corresponding gradients. Through an extensive array of both objective and subjective evaluations spanning diverse tasks, datasets, and model configurations, we showcase the efficacy of DIX in generating faithful and accurate explanation maps, while surpassing current state-of-the-art methods. Our code is available at: <https://github.com/dix-cikm23/dix>

1 INTRODUCTION

The AI revolution has led to significant advancements across various application fields, including computer vision [27, 37, 47, 49, 62], natural language processing [5, 10, 20, 22, 36, 45, 60, 67, 79], audio processing [24–26, 33, 39, 58], and recommender systems [8, 9, 15–19, 51, 52, 65, 83]. Specifically, in computer vision, deep Convolutional Neural Networks (CNNs) [49, 53, 61, 73], alongside recent Vision Transformers (ViTs) models [38] have risen to prominence, exhibiting outstanding performance in a variety of vision tasks [4, 21, 27, 34, 48]. This surge in popularity emphasizes the need to comprehend the underlying rationale driving the decisions and predictions of deep learning models.

Despite their remarkable achievements, most deep neural networks remain enigmatic, often considered black boxes due to their vast number of parameters and intricate non-linearities. This opacity has ignited the growth of explainable AI as a focal research area within the realm of deep learning. Consequently, numerous methodologies have been proposed for explaining the predictions of deep learning models in computer vision [7, 11, 12, 31, 69, 72, 85], natural language processing [14, 66], and recommender systems [13, 23, 44].

Explanation techniques aim to bridge the gap in understanding by generating heatmap-like explanation maps. These maps spotlight distinct input regions, attributing predictions to specific areas within the input image. Initially, rooted in gradient-based approaches, early methods generated explanation maps by analyzing the gradient of predictions concerning the input image [72, 73, 76]. Subsequently, several works [6, 29, 54, 69] proposed deriving explanation maps from the internal activation maps produced by the network, along with their gradients. Other techniques, such as Integrated Gradients (IG) [78], relying on path integration, created explanation maps by accumulating gradients from linear interpolations between input and reference images.

Predominantly applied to CNNs, the aforementioned methods arose before the emergence of Transformer-based architectures [79]. With the advent of ViT models [38], a variety of methodologies

were proposed to interpret and explain them, including recent explanation techniques like those presented in [30, 31].

This paper introduces Deep Integrated Explanations (DIX), a comprehensive approach aimed at explaining vision models, which finds applicability across both CNN and ViT architectures. DIX employs integration over the internal model representations and their gradients, facilitating the extraction of insights from any activation (or attention) map within the network.

We present a thorough objective and subjective evaluation, showcasing the efficacy of DIX on both CNN and ViT models. Our results reveal its superiority over other baselines across various explanation and segmentation tasks, encompassing diverse datasets, model architectures, and evaluation metrics. Additionally, we validate the credibility of DIX in producing faithful explanation maps through an extensive set of sanity tests, as outlined in [2].

2 RELATED WORK

2.1 Explanation Methods for CNNs

A diverse range of explanation methods were proposed for explaining CNN models, categorized into various types including perturbation-based methods, gradient methods, saliency-based methods, and gradient-free methods. Perturbation-based methods [41, 42] gauge output sensitivity concerning input through random perturbations applied in the input space. Saliency-based methods [32, 64, 72, 85–87] leverage feature maps obtained through forward propagation to interpret model predictions.

Gradient methods utilize prediction gradients with respect to the input or intermediate activation maps. These methods yield explanation maps based on the gradient itself or by a combination of the activation maps with their gradients [71, 77]. For instance, SmoothGrad [75] presents a smoothing approach, applied by adding random Gaussian noise to the input image at each iteration. Another notable example is the Grad-CAM (GC) [69] method, which leverages activation maps from the final convolutional layer in conjunction with their pooled gradients to generate explanation maps. The effectiveness of GC has subsequently inspired numerous follow-up work [6, 29, 43, 54].

Gradient-free methods generate explanation maps by manipulating activation maps without relying on gradient information [35, 81]. For instance, LIFT-CAM [55] utilizes the DeepLIFT [70] technique to estimate SHAP values of activation maps [63], which are then combined with the activation maps to produce the explanation map. However, gradient-free methods have a drawback: they neglect gradient information, thereby constraining their ability to steer explanations toward the target or predicted class.

Finally, a notable avenue of research pertains to path integration methods. Integrated Gradients (IG) [78] involves integration

*Both authors contributed equally to this research.

across interpolated image gradients. Blur IG (BIG) [84] focuses on introducing information using a baseline and adopts a path that gradually removes Gaussian blur from the attributed image. Guided IG (GIG) [57] refines IG by introducing an adaptive path strategy. By computing integration along an alternative path, it circumvents high gradient regions, often resulting in a reduction of irrelevant attributions.

Distinguished from the aforementioned works, DIX employs integration, facilitating interpolation on the internal representations produced by the network, and offers to combine the resulting explanation maps from all network layers. Furthermore, DIX does not confine the integrand to simple gradients, but rather encompasses an arbitrary function involving the activation (attention) maps and their gradients.

2.2 Explanation Methods for ViTs

Early attempts to explain Transformers employed the attention scores inherent to ViT models in order to glean insights w.r.t. the input [28, 79]. However, it is not clear how to combine the scores from different layers. Simple averaging the attention scores of each token, for example, leads to blurring of the signal [31].

Abnar and Zuidema [1] proposed the Rollout method to compute attention scores to input tokens at each layer by considering raw attention scores in a layer as well as those from precedent layers. Rollout improved results over the utilization of a single attention layer. However, by relying on simplistic aggregation assumptions, irrelevant tokens often become highlighted. LRP [3], proposed to propagate gradients from the output layer to the beginning, considering all the components in the transformer’s layers and not just the attention layers.

Recently, Chefer et al.[31] introduced Transformer Attribution (T-Attr), a class-specific Deep Taylor Decomposition method that employs relevance propagation for both positive and negative attributions. More recently, the same authors introduced Generic Attention Explainability (GAE)[30], which is an extension of T-Attr aimed at explaining Bi-Modal transformers. T-Attr and GAE stand as state-of-the-art methods for explaining ViT models, exhibiting superior performance when compared to other effective explanation methods, including LRP and partial LRP [80].

DIX differs from T-Attr and GAE in two main aspects: First, DIX is a versatile method capable of producing explanation maps for both CNNs and ViTs. Second, in the context of ViT models, DIX employs path integration on the interpolated attention matrices while incorporating the Gradient Rollout (GR) representation (a variant of the Rollout method) as the function for integration.

3 DEEP INTEGRATED EXPLANATIONS

Let $f : \mathbb{R}^{D_0} \rightarrow \mathbb{R}^K$ be a neural network with L hidden layers that takes an input (image) $\mathbf{x} \in \mathbb{R}^{D_0}$ and produces a prediction $f(\mathbf{x}) \in \mathbb{R}^K$. We denote \mathbf{x}^l ($1 \leq l \leq L$) as the intermediate representation generated by the l -th hidden layer in f (based on the input \mathbf{x}), and $f^l : \mathbb{R}^{D_l} \rightarrow \mathbb{R}^K$ as the sub network of f that takes \mathbf{x}^l as an input and outputs the prediction $f(\mathbf{x})$. Consequently, we have the relationship $f^l(\mathbf{x}^l) = f(\mathbf{x})$. Additionally, we denote $\mathbf{x}^0 = \mathbf{x}$ and $f^0 = f$.

Our assumption is that \mathbf{x}^l preserves the spatial structure of \mathbf{x} (though at a different resolution) such that each element in \mathbf{x}^l is associated with its corresponding elements in \mathbf{x} (e.g., this assumption holds true for CNNs). W.l.o.g, we restrict the discussion to multi-class classification problems, hence f outputs a vector assigning score to each class, and the score for the class k is denoted as $f_k(\mathbf{x})$.

Our objective is to *explain* the prediction $f_k(\mathbf{x})$ for the class k . In this work, we define an *explanation map* \mathbf{m}^l as a tensor assigning an attribution score to each element in \mathbf{x}^l w.r.t. the prediction $f_k^l(\mathbf{x}^l) = f_k(\mathbf{x})$. Consequently, \mathbf{m}^l must match the dimensions of \mathbf{x}^l . Note that our ultimate goal is to attribute the prediction to each element in the input \mathbf{x} , and due to the spatial structure preservation, each element in \mathbf{m}^l can be associated with a set of elements in \mathbf{x} .

Let $\mathbf{z}^l \in \mathbb{R}^{D_l}$ be a baseline serving as a reference for the informative representation \mathbf{x}^l . \mathbf{z}^l can be the null representation, random noise, or other baselines representing missing information. In what follows, we present a decomposition of the score difference $f_k(\mathbf{x}) - f_k^l(\mathbf{z}^l)$, from which an explanation map \mathbf{m}^l is derived.

Let C^l be a differentiable curve connecting \mathbf{z}^l to \mathbf{x}^l . C^l is parameterized by a vector function $\mathbf{r}^l : [0, 1] \rightarrow \mathbb{R}^{D_l}$ such that $\mathbf{r}^l(0) = \mathbf{z}^l$ and $\mathbf{r}^l(1) = \mathbf{x}^l$. The score difference $f_k(\mathbf{x}) - f_k^l(\mathbf{z}^l)$ can then be expressed as follows:

$$\begin{aligned} f_k(\mathbf{x}) - f_k^l(\mathbf{z}^l) &= f_k^l(\mathbf{r}^l(1)) - f_k^l(\mathbf{r}^l(0)) \\ &= \int_0^1 \frac{d}{dt} f_k^l(\mathbf{r}^l(t)) dt \\ &= \int_0^1 \nabla f_k^l(\mathbf{r}^l(t)) \cdot \frac{d\mathbf{r}^l(t)}{dt} dt \\ &= \sum_{i=1}^{D_l} \int_0^1 g_i^l(t) h_i^l(t) dt, \end{aligned} \quad (1)$$

where

$$g_i^l(t) = \frac{\partial f_k^l(\mathbf{r}^l(t))}{\partial r_i^l(t)} \quad \text{and} \quad h_i^l(t) = \frac{dr_i^l(t)}{dt},$$

with \cdot representing the dot product operator, and $r_i^l(t)$ being the i -th element in the interpolant $\mathbf{r}^l(t)$. The first equality in Eq. 1 consequents from the design of \mathbf{r}^l and the fact that $f^l(\mathbf{x}^l) = f(\mathbf{x})$. The second equality stems from the fundamental theorem of calculus. The third equality arises from the multivariate chain rule, and the last equality results from decomposing the dot product into a summation and then interchanging the order of finite sum and integration.

Equation 1 breaks down the score difference into a sum, where each term is a line integral along the i -th element of curve C^l , and the integrand is a function involving the partial derivative of the prediction $f_k^l(\mathbf{r}^l(t))$ w.r.t. the i -th element in the interpolant $\mathbf{r}^l(t)$. Consequently, each term in the sum resembles the attribution of the prediction $f_k(\mathbf{x})$ to an individual element in \mathbf{x}^l through the integrated partial derivatives along C^l . Equipped with Eq. 1, an explanation map for \mathbf{x}^l can be formed as follows:

$$\mathbf{m}^l = \int_0^1 \mathbf{g}^l(t) \circ \mathbf{h}^l(t) dt, \quad (2)$$

where \circ denotes the element-wise multiplication, $\mathbf{g}^l(t) = \frac{\partial f_k^l(\mathbf{r}^l(t))}{\partial \mathbf{r}^l(t)}$ is the gradient of the prediction w.r.t. the interpolant, and $\mathbf{h}^l(t) = \frac{d\mathbf{r}^l(t)}{dt}$. Note that $\mathbf{m}^l \in \mathbb{R}^{D_l}$ with $m_i^l = \int_0^1 g_i^l(t) h_i^l(t) dt$. Notable, for $l = 0$, Eq. 2 is equivalent to the IG [78] explanation map, where the interpolation takes place in the input space.

Equation 2 integrates the gradients of the interpolated activation maps $\mathbf{r}^l(t)$. Empirically, we found that incorporating the information from $\mathbf{r}^l(t)$ itself (beyond its gradient) yields enhanced explanations, both visually and quantitatively. This observation is consistent with previous works [6, 29, 69]. Furthermore, since for $l > 0$, \mathbf{m}^l does not match the spatial dimensions of the input \mathbf{x} , a subsequent transformation ψ^l is employed to ensure a proper match. To this end, we define the DIX explanation map as follows:

$$\mathbf{m}_{\text{DIX}}^l = \psi^l \left(\int_0^1 \phi \left(\mathbf{r}^l(t), \mathbf{g}^l(t) \right) \circ \mathbf{h}^l(t) dt \right), \quad (3)$$

where the exact implementation details of ϕ and ψ are architecture dependent and are outlined in Sec. 3.1.

In this work, we choose C^l to be the linear curve connecting \mathbf{z}^l to \mathbf{x}^l , hence

$$\mathbf{r}^l(t) = \mathbf{z}^l + t(\mathbf{x}^l - \mathbf{z}^l) \quad \text{and} \quad \mathbf{h}^l(t) = \mathbf{x}^l - \mathbf{z}^l. \quad (4)$$

In practice, the integration in Eq. 3 is numerically approximated as follows:

$$\mathbf{m}_{\text{DIX}}^l \approx \psi^l \left(\frac{\mathbf{x}^l - \mathbf{z}^l}{N} \circ \sum_{n=1}^N \phi \left(\mathbf{r}^l \left(\frac{n}{N} \right), \mathbf{g}^l \left(\frac{n}{N} \right) \right) \right), \quad (5)$$

where we have employed the linear interpolation from Eq. 4. In this work, we set $N = 10$. The complexity of DIX is similar to IG, except for the extra computation induced by ϕ and ψ^l .

Given that different network layers capture varying types of information and resolution, we propose aggregating information from explanation maps produced for different values of l . As such, the final explanation map is constructed as follows:

$$\mathbf{m}_{\text{DIX}}^S = \frac{1}{|S|} \sum_{l \in S} \mathbf{m}^l, \quad (6)$$

where S is a set indicating the layer indexes participating in the aggregation. Our experimentation indicates that the best-performing DIX configurations leverage a combination of explanation maps from the last two or three layers. Thus, in Sec. 5, we report results for $S = \{L-1, L\}$ (**DIX2**) and $S = \{L-2, L-1, L\}$ (**DIX3**). However, for the sake of completeness, we also present results for $S = \{L\}$ (**DIX1**) as part of our ablation study in Sec. 5.4.

3.1 Implementation Details

In this section, we describe concrete implementations of DIX for both CNN and ViT architectures.

CNN Models: In the case of CNNs, the architecture of f consists of residual blocks [50] that produces 3D tensors representing the activation maps \mathbf{x}^l . Correspondingly, \mathbf{z}^l is a 3D tensor where each channel is determined by broadcasting the minimum value of the respective activation map within \mathbf{x}^l . Furthermore, we set ϕ to the element-wise multiplication.

We motivate this design choice by the fact that $\mathbf{r}^l(t)$ represents the interpolated activation map, highlighting regions where filters are activated and patterns are detected. Its gradient gauges the attribution degree of the specific class of interest to each element in the activation map. Thus, we expect that regions exhibiting both large gradient and activation (of the same sign) will yield effective explanations. This property is achieved through element-wise multiplication of $\mathbf{r}^l(t)$ by its gradient $\mathbf{g}^l(t)$. Finally, ψ^l is set to the mean reduction on the channel axis followed by a resize operation yielding a 2D explanation map that matches the spatial dimensions of \mathbf{x} .

ViT Models: In ViT [37], the architecture of f consists of transformer encoder blocks producing 2D tensors (sequence of token representations). The input \mathbf{x} is transformed to a 2D tensor as well, where the first token is the [CLS] token, and the rest of the tokens are representations of patches in the original image.

In our implementation, we choose to interpolate on the attention matrices, which in turn affect the output produced by the encoder block. Specifically, $\mathbf{r}^l(t)$ is a 3D tensor that accommodates all the attention matrices produced by the l -th encoder block. The reference \mathbf{z}^l is set to the zero tensor (since the values in the attention matrix are in $[0, 1]$). ϕ implements a variant of the Attention Rollout (AR) method [1] that we name Gradient Rollout (GR). GR is similar to AR, with a slight modification. Instead of operating solely on the plain attention matrices, GR initially performs an element-wise multiplication of the attention matrices by their corresponding gradients. Following this, GR proceeds with the original Rollout computation [1], resulting in the first row of the derived matrix (associated with the [CLS] token). This output is further processed by truncating its initial element and reshaping it into a 14×14 matrix. The exact implementation of GR appears in our GitHub repository¹. Lastly, ψ^l remains consistent across all layers, conducting a resize operation to align with the spatial dimensions of \mathbf{x} .

4 EXPERIMENTAL SETUP

Our evaluation encompasses three distinct CNN architectures: ResNet101 (**RN**)[49], DenseNet201 (**DN**)[53], and ConvNext-Base (**CN**)[61], and two different architectures of ViT: ViT-Base (**ViT-B**) and ViT-Small (**ViT-S**)[37]. The information regarding preprocessing methodologies and direct access to all the aforementioned models can be found in our GitHub repository. DIX is evaluated and compared to other explanation methods through a series of explanation, segmentation, and sanity tests.

4.1 Explanation Metrics

It is difficult to quantify the quality of explainability methods, and there is no single agreed-upon metric. The explanations metrics in this study aim to assess how well the explanations align with hypothetical changes (counterfactuals) to the input. Essentially, it's about asking "what if" questions regarding the input and determining whether the explanations provided are consistent with those hypothetical scenarios. To comprehensively evaluate our method, we carefully followed several prominent evaluation protocols.

¹It is worth noting that our experimental findings suggest comparable performance when substituting the matrix product operation with summation within the context of the GR computation

Perturbation Tests. We followed the protocol from [31], which is the current state-of-the-art in explaining ViTs, and report the Negative Perturbation AUC (**NEG**) and the Positive Perturbation AUC (**POS**). NEG is a counterfactual test that entails a gradual blackout of the pixels in the original image in increasing order according to the explanation map while searching to see when the model’s top predicted class changes. By masking pixels in increasing order, we expect to remove the least relevant pixels first, and the model’s top predicted class is expected to remain unchanged for as long as possible. Results are measured in terms of the Area Under the Curve (AUC), and higher values are considered better. Accordingly, the POS test entails masking the pixels in decreasing order with the expectation that the model’s top predicted class will change quickly, hence in POS, lower values are better. In addition, we follow [68] and report the Insertion AUC (**INS**) and Deletion AUC (**DEL**) perturbation tests. INS and DEL entail a gradual blackout in increasing or decreasing order, similar to NEG and POS, respectively. However instead of tracking the point at which the top predicted class changes, in **INS** and **DEL** the AUC is computed with respect to the predicted probability of the top class. By masking pixels according to increasing/decreasing order of importance, we expect that the predicted probability of the top class will decrease slowly/quickly, respectively. Hence, for INS higher values are better and for DEL lower values are better.

ADP and PIC Tests. We follow [29] and report the Average Drop Percentage (**ADP**) and the Percentage Increase in Confidence (**PIC**) tests. Both tests relate to the change in the probability of the predicted class after applying the mask to the original image. A good explanation map is expected to highlight the most significant regions for decision-making. Hence, applying such a mask can be seen as a removal of the “background”. The ADP test measures the average percentage of model confidence drop after applying the mask. A good mask is expected to maintain the most relevant areas and minimize confidence drop, hence for ADP lower values are considered better. However, we note that ADP is a problematic metric since a naive all-ones mask yields an optimal ADP value of 0. Nevertheless, we included it for the sake of compatibility with previous works [29]. In some instances, the model’s confidence increases after applying a good explanation mask that removes a confusing background. Hence, PIC is a binary test that measures the percentage of instances in which the model’s confidence increased after applying the mask on the original input. For PIC higher values are considered better.

AIC and SIC Tests. We follow [56] and report the Accuracy Information Curve (**AIC**) and the Softmax Information Curve (**SIC**) tests. In these tests, we start with a completely blurred image and gradually sharpen the image areas that are deemed important by a given explanation method. Gradually sharpening the image areas increases the information content of the image. We then compare the explanation methods by measuring the approximate image entropy (e.g., compressed image size) and the model’s performance (e.g., model accuracy). The AIC metric measures the accuracy of a model as a function of the amount of information provided to the explanation method. AIC is defined as the AUC of the accuracy vs. information plot. The SIC metric measures the information content of the output of a softmax classifier as a function of the amount

of information provided to the explanation method. SIC is defined as the AUC of the entropy vs. information plot. The entropy of the softmax output is a measure of the uncertainty or randomness of the classifier’s predictions. For both AIC and SIC, the information provided to the method is quantified by the fraction of input features that are considered during the explanation process.

4.2 Segmentation Metrics

While possessing a superior segmentation capability does not necessarily imply a superior explanatory aptitude, we undertake this evaluation task for the sake of completeness in our comparison with previous works assessing this aspect [30, 31, 54, 82]. Segmentation accuracy is assessed according to the following metrics: Pixel Accuracy (**PA**), mean-intersection-over-union (**mIoU**), mean-average-precision (**mAP**), and the mean-F1 score (**mF1**) [31].

4.3 Datasets

Explanation maps are produced for the ImageNet [34] ILSVRC 2012 (**IN**) validation set, consisting of 50K images from 1000 classes. We follow the same setup from [31], where for each image, an explanation map is produced w.r.t. the class predicted by the model. Segmentation tests are conducted on three datasets: (1) ImageNet-Segmentation [46] (**IN-Seg**): This is a subset of ImageNet validation set consisting of 4,276 images from 445 classes for which annotated segmentations are available. (2) Microsoft Common Objects in Context 2017 [59] (**COCO**): This is a validation set that contains 5,000 annotated segmentation images from 80 different classes. Some images consist of multi-label annotations (multiple annotated objects). In our evaluation, all annotated objects in the image are considered as the ground-truth. (3) PASCAL Visual Object Classes 2012 [40] (**VOC**): This is a validation set that contains annotated segmentations for 1,449 images from 20 classes.

4.4 Evaluated Methods

Our evaluation encompasses a comprehensive assessment of various explanation methods, including gradient-based approaches, path-integration techniques, as well as gradient-free methods.

For CNN models, the following explanation techniques are considered: Integrated Gradients (**IG**) [78], Guided IG (**GIG**) [57], Blur IG (**BIG**) [84], Ablation-CAM (**AC**) [35], Layer-CAM (**LC**) [54], LIFT-CAM (**LIFT**) [55], Grad-CAM (**GC**) [69], Grad-CAM++ (**GC++**) [29], X-Grad-CAM (**XGC**) [43], and FullGrad (**FG**) [77].

For ViT models, we consider two state-of-the-art methods: Transformer Attribution (**T-Attr**) [31] and Generic Attention Explainability (**GAE**) [30]. Both methods were shown to outperform other strong baselines such as partial LRP [80], and GC [30] for transformers. A detailed description of all explanation methods is provided in our GitHub repository. Lastly, our universal DIX method is evaluated on both CNNs and ViTs, where we consider two versions: DIX2 and DIX3 following the description in Sec. 3.

4.5 Sanity Tests for Explanation Methods

To comprehensively assess the robustness and credibility of DIX, we conducted the *parameter randomization* and *data randomization* sanity tests as outlined in [2]. For these evaluations, we employed DIX3, along with the VGG-19[74] model and the IN dataset.

Parameter Randomization Test. This test compares explanation maps generated by the explanation method under two model setups: (1) Trained - a model trained on the dataset (e.g., pretrained VGG-19 on ImageNet), and (2) Random - the same model architecture with randomized weights. Significant differences in explanation maps between the trained and random models suggest sensitivity to model parameters. Conversely, similar maps indicate insensitivity and limited utility for model explanation.

Two types of parameter randomization tests are conducted on a trained model:

(1) *Cascading Randomization:* We randomize model weights layer by layer, starting from the top and progressing to the bottom. This process randomizes the learned weights from top to bottom.

(2) *Independent Randomization:* We randomize each individual layer’s weights, one layer at a time (while keeping all other layers’ weights fixed). This allows us to evaluate each layer’s impact on explanation map sensitivity, independently.

In both tests, we compare the resulting explanations obtained by using the model with random weights to those derived from the original weights of the model.

Data Randomization Test. The data randomization test assesses an explanation method’s sensitivity to data labeling. We compare explanation maps for two models with identical architectures, but trained on different datasets: one with original labels and another with labels randomly permuted. Sensitivity to labeling is indicated by significantly different explanation maps, while insensitivity suggests independence from instance-label relationships. To conduct the data randomization test, we permute the training labels in the dataset and train the model to achieve a training set accuracy greater than 95%. Note that the resulting model’s test accuracy is never better than randomly guessing a label. We then compute explanations on the same test inputs for both the model trained on true labels and the model trained on randomly permuted labels.

5 RESULTS

5.1 Explanation Tests

Tables 1 and 2 provide a comprehensive explanation tests for CNN and ViT models, respectively. We report results for all combinations of datasets, models, methods, and metrics. Our analysis demonstrates that DIX consistently surpasses all baseline methods across a spectrum of metrics and architectural configurations. On CNN-based DIX variations (Tab. 1), DIX3 outclasses DIX2 in terms of NEG, INS, SIC, and AIC metrics for both RN and DN backbones, while demonstrating dominance across all metrics for the CN backbone. Regarding the ViT-based DIX variants (Tab. 2), DIX3 outperforms DIX2 across all metrics (with the exception of PIC on ViT-B, and PIC and ADP on ViT-S). These trends showcase the advantage of aggregating information from more layers. In the context of CNNs, the second-best performing methods are GC and GC++, which leverage both activation and gradients to outperform other approaches across most evaluation metrics. Additionally, we note that path integration techniques (IG, BIG, and GIG) demonstrate competitive results in terms of POS and DEL metrics, while displaying comparatively weaker performance in other aspects. This disparity may be attributed to the grainy output maps generated by path

integration techniques, as evidenced in Fig.3 for IG explanation maps on CNNs. These methods ignore the activations and integrate on the image domain only, hence missing some of the key features. This is particularly evident in the significant contrast between their strong performance on POS and the corresponding weaker performance on NEG. As path integration methods produce sparse maps that can negatively affect performance in certain metrics, we extend our analysis to encompass the SIC and AIC metrics as well [56]. These metrics were originally employed to assess GIG[57] and BIG[84]. Yet, the incorporation of SIC and AIC did not alter the trend of the results. This suggests that DIX is highly effective for generating high-quality explanation maps. Finally, we present an ablation study in Section 5.4, aimed at comparing diverse versions and alternatives of DIX. This analysis serves to emphasize the effectiveness of the integration process and the strategic utilization of information from multiple layers within the DIX methodology.

5.2 Segmentation Tests

Tables 3 and 4 present segmentation tests results on CNN and ViT models, respectively. The results are reported for all combinations of datasets, models, explanation methods, and segmentation metrics. In these experiments, only the 5 best performing CNN explanation methods from Tab. 1 are considered. Once again, it becomes evident that DIX consistently delivers the most favorable segmentation outcomes for both CNN and ViT models. This outcome can be rationalized by the localized and precise maps that DIX generates.

5.3 Qualitative Evaluation

Figure 1 presents a qualitative comparison of the explanation maps obtained by the top-performing CNN explanation methods on a large set of examples that are randomly drawn from multiple classes from the IN dataset. Arguably, DIX (DIX3) produces the most accurate explanation maps in terms of class discrimination and localization. These results correlate well with the trends from Tabs. 1 and 3. We observe that in the case of class ‘accordion, piano accordion, and squeeze box’, DIX focuses mostly on the correct item, while the gradient-free methods focus mostly on other parts of the image, exposing their class-agnostic behavior. Moreover, a similar trend is observed with the ‘sturgeon’ class, in which DIX is the only one to focus on the relevant class. Figure 2 presents a qualitative comparison of the explanation maps obtained by explanation methods for ViT. Once again, we see that DIX produces the most accurate and focused explanation maps.

5.4 Ablation Study

In this work we present and evaluate DIX2 and DIX3. In this section, we justify these choices via an ablation study. To this end, we set $n = 10$, and consider three alternatives: (1) **DIX1** - we use the last layer as the only layer to interpolate i.e., $S = \{L\}$. (2) **DIX2-MUL** - $\mathbf{m}_{\text{DIX}}^L$ and $\mathbf{m}_{\text{DIX}}^{L-1}$ are being element-wise multiplied to produce the final explanation map. (3) **DIX3-GRADS** - we use the plain gradients without explicitly incorporating the information from the activation or attention maps.

Table 5 reports the results for the RN and ViT-B models on the IN dataset. For the sake of completeness, we further include the results for IG, DIX2, and DIX3 (taken from Tabs. 1 and 2). First, we can

Table 1: Explanation tests results on the IN dataset (CNN models): For POS, DEL and ADP, lower is better. For NEG, INS, PIC, SIC and AIC, higher is better. See Sec. 5.1 for details.

	GC	GC++	LIFT	AC	IG	GIG	BIG	FG	LC	XGC	DIX2	DIX3	
RN	NEG	<u>56.41</u>	55.20	55.39	54.98	45.66	43.97	42.25	54.81	53.52	53.46	57.13	
	POS	17.82	18.01	17.53	19.38	17.24	17.68	17.44	18.06	17.92	21.02	15.69	<u>17.11</u>
	INS	<u>48.14</u>	47.56	45.39	47.05	39.87	37.92	36.04	42.68	46.11	43.26	48.09	48.91
	DEL	13.97	14.17	15.32	14.23	13.49	14.18	13.95	14.64	14.31	14.98	12.84	<u>13.36</u>
	ADP	17.87	16.91	18.03	16.18	37.52	35.28	40.85	21.06	24.34	17.02	15.68	<u>16.02</u>
	PIC	36.69	36.53	35.95	35.52	19.94	18.72	24.53	31.59	35.43	36.18	40.21	<u>37.29</u>
	SIC	76.91	76.44	76.73	73.36	54.67	55.04	56.98	75.35	73.93	72.64	<u>77.61</u>	78.12
	AIC	74.36	71.97	72.76	70.35	51.92	53.38	53.36	71.49	65.77	69.85	<u>76.09</u>	76.34
CN	NEG	52.86	53.82	53.98	53.68	45.24	41.43	40.72	52.06	54.12	52.13	<u>54.40</u>	55.23
	POS	17.52	17.85	18.23	18.19	17.42	18.03	18.14	18.26	17.58	20.83	<u>16.96</u>	16.51
	INS	45.65	45.19	43.86	49.18	37.22	32.99	31.02	42.01	44.14	42.07	<u>49.53</u>	49.86
	DEL	13.43	14.17	15.18	14.73	12.36	13.08	13.29	14.21	13.64	14.78	<u>11.95</u>	11.74
	ADP	22.46	22.35	29.13	24.38	36.98	35.79	41.73	30.75	37.62	25.68	<u>22.24</u>	22.19
	PIC	23.16	24.42	22.34	24.59	17.65	13.12	20.69	22.13	22.17	23.26	<u>28.31</u>	28.47
	SIC	65.93	67.94	54.75	63.95	53.36	58.35	57.27	62.84	69.11	59.12	<u>69.83</u>	70.18
	AIC	75.64	75.52	57.06	71.53	51.68	55.82	53.82	67.15	75.41	62.38	<u>76.44</u>	77.29
DN	NEG	<u>57.40</u>	57.16	58.01	56.63	40.74	37.31	36.67	56.79	56.96	55.74	57.31	58.25
	POS	17.75	17.81	18.87	18.67	17.31	17.46	17.38	17.84	17.62	18.67	16.59	<u>17.14</u>
	INS	<u>51.09</u>	50.89	50.63	50.41	37.58	33.31	31.32	50.44	50.60	49.62	50.97	51.58
	DEL	13.61	13.63	13.29	15.31	13.26	13.27	13.54	14.34	13.85	14.75	12.73	<u>12.98</u>
	ADP	17.46	17.01	19.45	17.13	35.61	34.51	40.04	20.21	24.23	19.59	16.29	<u>16.58</u>
	PIC	34.68	35.21	34.13	31.22	22.35	16.62	26.18	31.05	33.81	30.39	38.91	<u>37.78</u>
	SIC	75.62	74.75	74.72	73.94	54.59	58.55	57.66	72.93	74.34	73.94	<u>77.24</u>	77.32
	AIC	74.22	71.82	72.65	70.21	54.74	54.56	56.08	70.63	71.82	70.12	<u>75.98</u>	76.39

Table 2: Explanation tests results on the IN dataset (ViT models): For POS, DEL and ADP, lower is better. For NEG, INS, PIC, SIC and AIC, higher is better. See Sec. 5.1 for details.

	T-Attr	GAE	DIX2	DIX3	
ViT-B	NEG	54.16	54.61	<u>56.43</u>	56.94
	POS	17.03	17.32	<u>15.10</u>	14.85
	INS	48.58	48.96	<u>49.51</u>	50.59
	DEL	14.20	14.37	<u>12.62</u>	12.16
	ADP	54.02	37.84	<u>35.93</u>	35.58
	PIC	13.37	23.65	28.21	<u>27.41</u>
	SIC	68.59	68.35	<u>68.94</u>	69.11
	AIC	61.34	57.92	<u>62.42</u>	65.03
ViT-S	NEG	53.29	52.81	<u>55.98</u>	56.13
	POS	14.16	14.75	<u>13.09</u>	12.32
	INS	45.72	45.21	<u>46.62</u>	47.36
	DEL	11.28	11.92	<u>11.18</u>	10.56
	ADP	51.94	36.98	36.31	<u>36.57</u>
	PIC	13.67	8.68	18.39	<u>18.25</u>
	SIC	69.46	70.19	<u>70.92</u>	71.55
	AIC	63.86	64.49	<u>65.17</u>	65.58

see the superior performance of DIX2 and DIX3 across all metrics and models. We further observe that both DIX1 and DIX2-MUL fall short in comparison to DIX2. This observation underscores the inherent necessity of incorporating information from additional layers and shows the advantages of aggregation via summation. When aggregating the explanation maps of different layers, the objective is to effectively incorporate data from each map to capture a richer spectrum of insights and class-specific signals. Notably, the multiplication operator exhibits a behavior akin to intersection, where both high pixel values are required for proper appearance in the final map. This characteristic, as depicted in Figure 3, contrasts with the intended outcome. Furthermore, the superiority of DIX3 over DIX3-GRADS underscores the benefit from exploiting intermediate representation information alongside its corresponding gradients, which contributes to the generation of localized, accurate and class discriminative explanation maps. The results presented in Table 5 highlight a distinct advantage for IG and DIX2-MUL with respect to the POS and DEL metrics when compared to DIX3-GRADS and DIX1, both of which generate less concentrated explanation maps. This is due to the fact that the deletion of the most relevant pixels results in fewer pixels being removed, and the mask is more focused on a subset of pixels. DIX1, for instance, produces less focused explanation maps that may highlight irrelevant areas. Such coarse highlighting leads to a slower decrease in the prediction

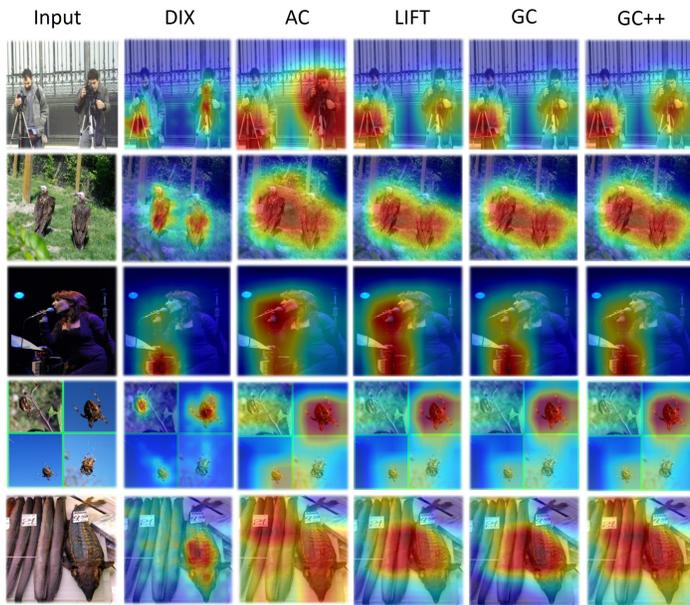


Figure 1: CNN Qualitative Results: Explanation maps produced using RN w.r.t. the classes (top to bottom): ‘tripod’, ‘vulture’, ‘accordion, squeeze box’, ‘garden spider, Aranea diademata’, and ‘sturgeon’.

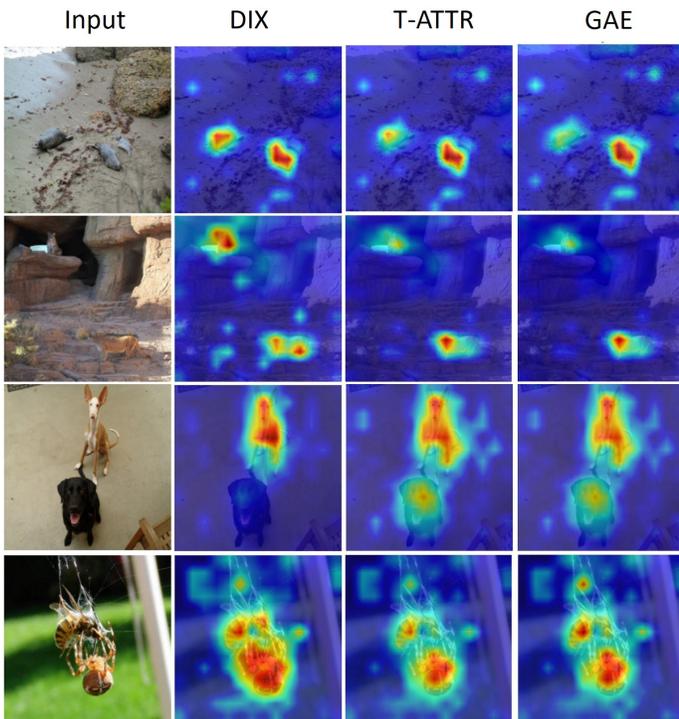


Figure 2: ViT Qualitative Results: Explanation maps produced using ViT-B w.r.t. the classes (top to bottom): ‘sea lion’, ‘cougar, puma, catamount, mountain lion, painter, panther, Felis concolor’, ‘Ibizan hound, Ibizan Podenco’, and ‘garden spider, Aranea diademata’.

Table 3: Segmentation tests on three datasets (CNN models). For all metrics, higher is better. See Sec. 5.2 for details.

		GC	GC++	LIFT	AC	DIX2	DIX3		
IN-SEG	CN	PA	77.01	77.54	63.77	77.04	<u>78.32</u>	78.93	
		mAP	81.01	85.63	69.40	86.93	<u>87.13</u>	87.34	
		mIoU	56.58	58.35	53.81	58.42	<u>58.64</u>	58.79	
		mF1	36.88	38.26	35.91	41.29	<u>42.51</u>	42.95	
IN-SEG	RN	PA	71.93	71.96	71.68	70.36	<u>72.43</u>	73.17	
		mAP	84.21	84.23	83.79	81.14	<u>84.58</u>	85.37	
		mIoU	53.06	53.29	52.17	52.91	<u>53.93</u>	54.16	
		mF1	42.51	42.68	41.95	42.08	<u>42.75</u>	43.18	
IN-SEG	DN	PA	73.00	73.21	72.87	72.44	<u>73.58</u>	73.90	
		mAP	85.04	85.53	84.82	84.62	<u>85.57</u>	85.98	
		mIoU	54.18	54.57	54.11	54.89	<u>55.42</u>	56.03	
		mF1	41.74	42.58	41.61	43.51	<u>43.71</u>	43.79	
COCO	CN	PA	68.75	66.49	60.37	64.10	<u>68.87</u>	69.38	
		mAP	75.02	75.21	67.98	76.09	<u>76.94</u>	77.43	
		mIoU	43.46	44.01	37.08	44.27	<u>44.89</u>	45.06	
		mF1	28.96	29.85	26.92	30.81	<u>31.28</u>	31.99	
	COCO	RN	PA	64.17	64.39	64.02	63.90	<u>64.75</u>	64.94
			mAP	74.19	74.27	73.78	72.80	<u>74.38</u>	74.91
			mIoU	42.37	43.25	42.59	42.88	<u>43.54</u>	43.87
			mF1	31.64	32.82	31.77	32.41	<u>33.39</u>	33.71
	COCO	DN	PA	63.50	64.06	63.25	64.51	<u>64.98</u>	65.37
			mAP	72.61	73.07	72.15	73.85	<u>74.02</u>	74.67
			mIoU	43.02	43.75	42.85	44.16	<u>44.75</u>	44.82
			mF1	31.04	32.31	30.83	33.93	<u>34.14</u>	34.59
VOC	CN	PA	72.54	72.09	63.32	69.83	<u>72.68</u>	72.81	
		mAP	77.27	79.47	68.83	80.45	<u>81.35</u>	81.79	
		mIoU	50.28	50.63	48.86	49.76	<u>51.12</u>	51.29	
		mF1	35.24	35.67	33.26	34.51	<u>35.92</u>	36.57	
	VOC	RN	PA	68.74	69.01	68.61	68.00	<u>69.38</u>	69.74
			mAP	79.68	79.96	79.41	78.02	<u>81.02</u>	81.49
			mIoU	49.44	49.91	49.15	49.32	<u>50.43</u>	51.58
			mF1	33.08	33.56	32.69	32.74	<u>34.28</u>	34.68
	VOC	DN	PA	68.43	68.78	68.24	68.36	<u>68.89</u>	68.95
			mAP	78.68	79.06	78.52	78.62	<u>79.43</u>	79.66
			mIoU	49.29	49.68	49.03	49.11	<u>49.91</u>	50.24
			mF1	32.92	33.83	32.28	32.56	<u>34.11</u>	34.26

score during the deletion process. On the contrary, DIX1 and DIX3-GRADS exhibit superior performance in relation to the NEG and INS metrics. This divergence in performance can be attributed to the expansive nature of their explanation map, resulting in numerous pixels that require removal. In the context of the NEG metric, this characteristic contributes to a slow decrease in the prediction score during the deletion process and, subsequently, a larger area under the curve (AUC).

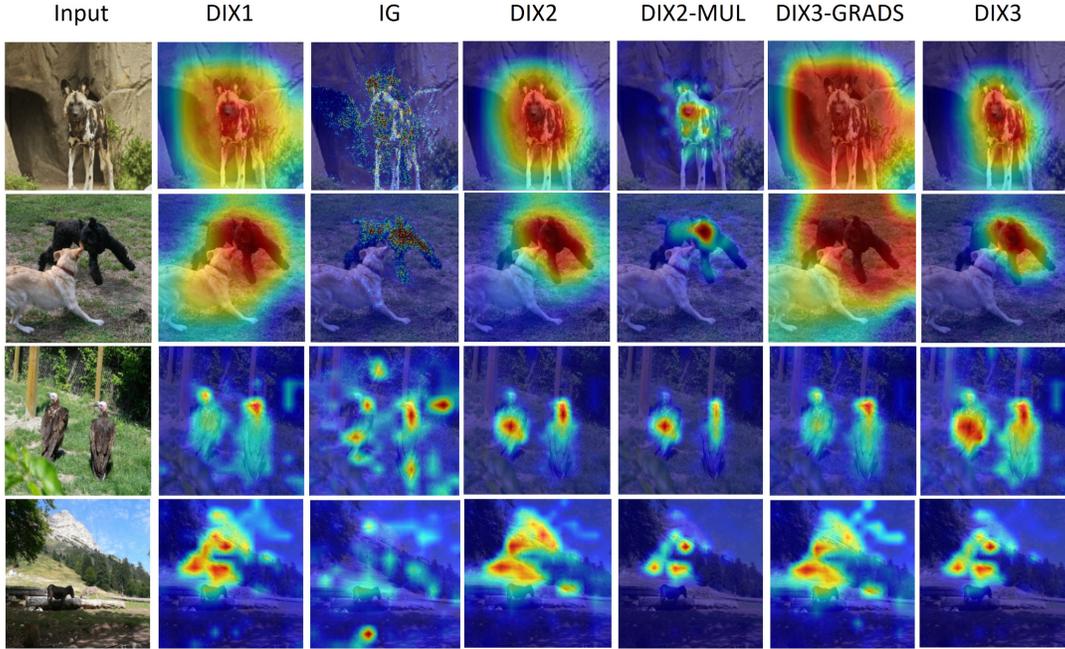


Figure 3: Ablation study results. Explanation maps produced using RN (rows 1,2) and ViT-B (rows 3,4) w.r.t. the classes (top to bottom): ‘African hunting dog, hyena dog, Cape hunting dog, Lycaon pictus’, ‘Kerry blue terrier’, ‘vulture’, ‘alp’.

Table 4: Segmentation tests on three datasets (ViT models). For all metrics, higher is better. See Sec. 5.2 for details.

		T-Attr	GAE	DIX2	DIX3			
IN-Seg	ViT-B	PA	79.70	76.30	<u>79.91</u>	81.02		
		mAP	86.03	85.28	<u>87.12</u>	87.45		
		mIoU	61.95	58.34	<u>62.53</u>	63.47		
		mF1	40.17	41.85	<u>44.94</u>	45.66		
	ViT-S	PA	80.86	76.66	<u>81.54</u>	81.83		
		mAP	86.13	84.23	<u>86.48</u>	86.96		
		mIoU	63.61	57.70	<u>64.13</u>	64.67		
		mF1	43.60	40.72	<u>46.34</u>	46.82		
		COCO	ViT-B	PA	68.89	67.10	<u>68.95</u>	69.42
				mAP	78.57	78.72	<u>80.63</u>	81.22
mIoU	46.62			46.51	<u>47.75</u>	47.79		
mF1	26.28			31.70	<u>33.87</u>	34.12		
ViT-S	PA		69.90	67.95	<u>70.41</u>	70.64		
	mAP		79.28	78.65	<u>80.55</u>	80.89		
VOC	ViT-B	mIoU	48.62	46.52	<u>50.81</u>	51.22		
		mF1	30.88	30.96	<u>35.61</u>	35.74		
		PA	73.70	71.32	<u>75.33</u>	75.84		
		mAP	81.08	80.88	<u>81.75</u>	81.89		
		mIoU	53.09	51.82	<u>53.62</u>	53.71		
		mF1	31.50	35.72	<u>36.38</u>	36.59		
	ViT-S	PA	74.96	71.85	<u>76.35</u>	76.56		
		mAP	81.76	80.60	<u>82.74</u>	82.91		
		mIoU	55.37	51.55	<u>55.83</u>	55.98		
		mF1	36.03	34.95	<u>39.27</u>	39.41		

Table 5: Ablation study results for various DIX configurations on the IN dataset. See Sec. 5.4 for details.

		DIX1	IG	DIX2	DIX2-MUL	DIX3-GRADS	DIX3	
RN	NEG	55.47	45.66	56.28	55.24	56.05	57.13	
	POS	17.47	17.24	15.69	17.28	18.13	<u>17.11</u>	
	INS	47.53	39.87	<u>48.09</u>	47.13	47.88	48.91	
	DEL	13.72	13.49	12.84	13.59	14.52	<u>13.36</u>	
	ADP	17.21	37.52	15.68	21.38	17.43	<u>16.02</u>	
	PIC	36.54	19.94	40.21	28.46	37.10	<u>37.29</u>	
	SIC	76.85	54.67	<u>77.61</u>	75.17	76.13	78.12	
	AIC	75.48	51.92	<u>76.09</u>	74.21	74.88	76.34	
	ViT-B	NEG	55.98	40.94	56.43	55.62	55.78	56.94
		POS	15.49	22.43	<u>15.10</u>	15.37	15.81	14.85
INS		49.38	35.07	<u>49.51</u>	49.27	49.33	50.59	
DEL		13.06	17.90	<u>12.62</u>	12.85	13.12	12.16	
ADP		36.96	41.35	<u>35.93</u>	38.62	36.08	<u>35.58</u>	
PIC		26.94	16.89	28.21	26.39	27.13	27.41	
SIC		67.79	58.91	<u>68.94</u>	68.43	68.32	69.11	
AIC		61.56	54.93	<u>62.42</u>	62.18	61.94	65.03	

5.5 Sanity Tests

In what follows, we show that DIX passes all sanity tests successfully, thereby furnishing additional substantiation for the authenticity of DIX as a robust machinery for generating accurate explanation maps.

Cascading Randomization. Figure 4 shows the Spearman correlation, computed as an average across 50K examples, between the

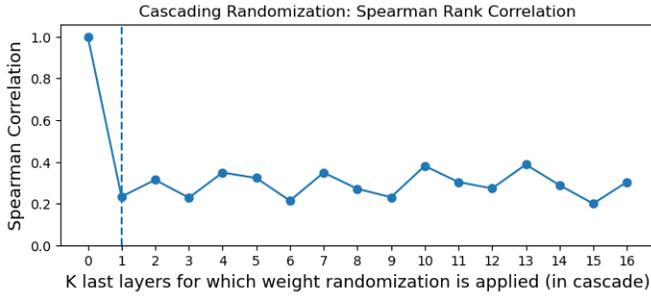


Figure 4: Cascading Randomization: The presented graph depicts the Spearman rank correlation (averaged on 50K examples) between the explanation produced by DIX using the original and randomized model’s weights. The x-axis corresponds to the number of layers being randomized, starting from the output layer. The dashed line indicates the point where the successive randomization of the network commences, which is at the top layer. The first dot ($x=0$) corresponds to no randomization (the original model is used), yielding perfect correlation between the explanation maps.

initial explanation map derived from DIX and the original (pre-trained) VGG-19 model, as well as the explanation map obtained from DIX and each of the cascade randomization variations of the original model. The markers on the x-axis are between '0' and '16', where $x = k$ means that the weights of the last k layers of the model are randomized. Notably, at $x = 0$, there is no randomization, hence the correlation with the original model is perfect. Starting from $x = 1$ (marked by the horizontal dashed line) and up to $x = 16$, the graph depicts a progressive cascade randomization of the original model. It is evident that, with an increase in the randomization of layer weights, the correlation with the explanation map of the original model experiences a significant decline. This trend underscores the sensitivity of DIX to the parameters of the model, a characteristic that is both anticipated and desirable for any explanation approach, as emphasized by [2].

Independent Randomization. Figure 5 presents results for the independent randomization tests. For $x = 0$, no randomization was introduced and the correlation to the original model is perfect. As x advances to i ($i > 0$), the graph illustrates the correlation of the initial model with a configuration where solely the weights of the i -th penultimate layer were randomized, while the weights of all other layers remained the same. Evidently, the correlation values exhibit a consistent diminution across all layers, underscoring DIX’s sensitivity to weight randomization on an individual layer basis. This characteristic is fundamentally desirable for an explanation technique, serving as evidence of its sensitivity to the distinct layers within the model.

Data Randomization. Figure 6 presents a box plot computed for the Spearman correlation values obtained for paired explanation maps (50K examples): one produced using the original model that is trained with the ground truth, and another produced by the model trained with the permuted labels. We can see that the correlation values are very low indicating DIX’s sensitivity to the labeling of

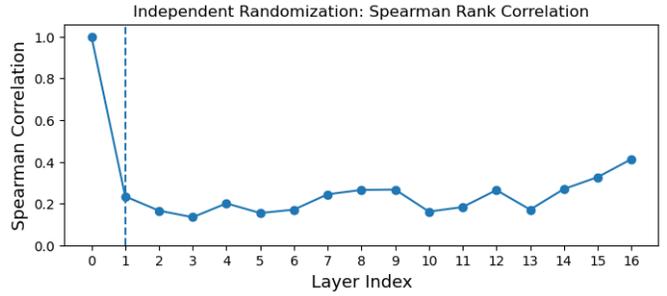


Figure 5: Independent Randomization: The y-axis of the presented graph represents the rank correlation between the original and randomized explanations, with each point on the x-axis corresponding to a specific layer of the model. The dashed line marks the point where the randomization of the network layers commences, which is at the top layer.

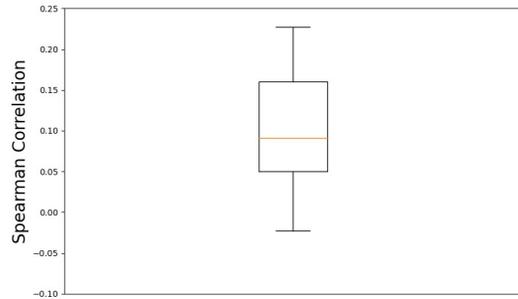


Figure 6: Data Randomization Test: Spearman rank correlation box plot for DIX with the VGG-19 model.

the dataset. Hence, we conclude that DIX successfully passes the data randomization test.

6 CONCLUSION

We presented the Deep Integrated Explanations (DIX) method for producing explanations for vision models. DIX is founded upon the accumulation of maps originating from multiple layers, encompassing interpolated network representations along with their corresponding gradients. We demonstrated the applicability of DIX for explaining CNN and ViT models, where it is shown to outperform state-of-the-art explanation methods across multiple tasks, datasets, network architectures, and metrics. Finally, we validated DIX as a machinery for generating faithful explanation maps via an extensive set of sanity tests.

REFERENCES

- [1] Samira Abnar and Willem Zuidema. 2020. Quantifying Attention Flow in Transformers. *arXiv preprint arXiv:2005.00928* (2020).
- [2] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. 2018. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*. 9505–9515.
- [3] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one* 10, 7 (2015), e0130140.
- [4] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. 2017. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* 39, 12 (2017), 2481–2495.
- [5] Oren Barkan. 2017. Bayesian neural word embedding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 31.
- [6] Oren Barkan, Omri Armstrong, Amir Hertz, Avi Caciularu, Ori Katz, Itzik Malkiel, and Noam Koenigstein. 2021. GAM: Explainable Visual Similarity and Classification via Gradient Activation Maps. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 68–77.
- [7] Oren Barkan, Yuval Asher, Amit Eshel, Yehonatan Elisha, and Noam Koenigstein. 2023. Learning to Explain: A Model-Agnostic Framework for Explaining Black Box Models. In *2023 IEEE International Conference on Data Mining (ICDM)*. IEEE.
- [8] Oren Barkan, Avi Caciularu, Ori Katz, and Noam Koenigstein. 2020. Attentive item2vec: Neural attentive user representations. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 3377–3381.
- [9] Oren Barkan, Avi Caciularu, Idan Rejwan, Ori Katz, Jonathan Weill, Itzik Malkiel, and Noam Koenigstein. 2020. Cold item recommendations via hierarchical item2vec. In *2020 IEEE International Conference on Data Mining (ICDM)*. IEEE, 912–917.
- [10] Oren Barkan, Avi Caciularu, Idan Rejwan, Ori Katz, Jonathan Weill, Itzik Malkiel, and Noam Koenigstein. 2021. Representation learning via variational bayesian networks. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 78–88.
- [11] Oren Barkan, Yehonatan Elisha, Yuval Asher, Amit Eshel, and Noam Koenigstein. 2023. Visual Explanations via Iterated Integrated Attributions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 2073–2084.
- [12] Oren Barkan, Yehonatan Elisha, Jonathan Weill, Yuval Asher, Amit Eshel, and Noam Koenigstein. 2023. Stochastic Integrated Explanations for Vision Models. In *2023 IEEE International Conference on Data Mining (ICDM)*. IEEE.
- [13] Oren Barkan, Yonatan Fuchs, Avi Caciularu, and Noam Koenigstein. 2020. Explainable recommendations via attentive multi-persona collaborative filtering. In *Proceedings of the 14th ACM Conference on Recommender Systems*. 468–473.
- [14] Oren Barkan, Edan Hahuon, Avi Caciularu, Ori Katz, Itzik Malkiel, Omri Armstrong, and Noam Koenigstein. 2021. Grad-sam: Explaining transformers via gradient self-attention maps. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 2882–2887.
- [15] Oren Barkan, Roy Hirsch, Ori Katz, Avi Caciularu, and Noam Koenigstein. 2021. Anchor-based collaborative filtering. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 2877–2881.
- [16] Oren Barkan, Roy Hirsch, Ori Katz, Avi Caciularu, Yoni Weill, and Noam Koenigstein. 2021. Cold start revisited: A deep hybrid recommender with cold-warm item harmonization. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 3260–3264.
- [17] Oren Barkan, Ori Katz, and Noam Koenigstein. 2020. Neural attentive multiview machines. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 3357–3361.
- [18] Oren Barkan and Noam Koenigstein. 2016. Item2vec: neural item embedding for collaborative filtering. In *2016 IEEE 26th International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 1–6.
- [19] Oren Barkan, Noam Koenigstein, Eylon Yogeve, and Ori Katz. 2019. CB2CF: a neural multiview content-to-collaborative filtering model for completely cold item recommendations. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 228–236.
- [20] Oren Barkan, Noam Razin, Itzik Malkiel, Ori Katz, Avi Caciularu, and Noam Koenigstein. 2020. Scalable attentive sentence pair modeling via distilled sentence embedding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 3235–3242.
- [21] Oren Barkan, Tal Reiss, Jonathan Weill, Ori Katz, Roy Hirsch, Itzik Malkiel, and Noam Koenigstein. 2023. Efficient Discovery and Effective Evaluation of Visual Perceptual Similarity: A Benchmark and Beyond. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 20007–20018.
- [22] Oren Barkan, Idan Rejwan, Avi Caciularu, and Noam Koenigstein. 2020. Bayesian hierarchical words representation learning. *arXiv preprint arXiv:2004.07126* (2020).
- [23] Oren Barkan, Tom Shaked, Yonatan Fuchs, and Noam Koenigstein. 2023. Modeling users' heterogeneous taste with diversified attentive user profiles. *User Modeling and User-Adapted Interaction* (2023), 1–31.
- [24] Oren Barkan, Shlomi Shvartzman, Noy Uzrad, Almog Elharar, Moshe Laufer, and Noam Koenigstein. 2023. InverSynth II: Sound matching via self-supervised synthesizer-proxy and inference-time finetuning. ISMIR.
- [25] Oren Barkan and David Tsiris. 2019. Deep synthesizer parameter estimation. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 3887–3891.
- [26] Oren Barkan, David Tsiris, Ori Katz, and Noam Koenigstein. 2019. Inversynth: Deep estimation of synthesizer parameter configurations from audio signals. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 27, 12 (2019), 2385–2396.
- [27] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *European conference on computer vision*. 213–229.
- [28] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *European conference on computer vision*. Springer, 213–229.
- [29] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. 2018. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 839–847.
- [30] Hila Chefer, Shir Gur, and Lior Wolf. 2021. Generic attention-model explainability for interpreting bi-modal and encoder-decoder transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 397–406.
- [31] Hila Chefer, Shir Gur, and Lior Wolf. 2021. Transformer interpretability beyond attention visualization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 782–791.
- [32] Piotr Dabkowski and Yarin Gal. 2017. Real time image saliency for black box classifiers. In *Advances in Neural Information Processing Systems*. 6970–6979.
- [33] Alexandre Défossez, Neil Zeghidour, Nicolas Usunier, Léon Bottou, and Francis Bach. 2018. Sing: Symbol-to-instrument neural generator. *Advances in neural information processing systems* 31 (2018).
- [34] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. 2009. ImageNet: A Large-Scale Hierarchical Image Database. In *Computer Vision and Pattern Recognition (CVPR)*.
- [35] Saurabh Satish Desai and H. G. Ramaswamy. 2020. Ablation-CAM: Visual Explanations for Deep Convolutional Network via Gradient-free Localization. *2020 IEEE Winter Conference on Applications of Computer Vision (WACV)* (2020), 972–980.
- [36] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [37] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [38] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [39] Jesse Engel, Kumar Krishna Agrawal, Shuo Chen, Ishaan Gulrajani, Chris Donahue, and Adam Roberts. 2019. Gansynth: Adversarial neural audio synthesis. *arXiv preprint arXiv:1902.08710* (2019).
- [40] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John M. Winn, and Andrew Zisserman. 2009. The Pascal Visual Object Classes (VOC) Challenge. *International Journal of Computer Vision* 88 (2009), 303–338.
- [41] Ruth Fong, Mandela Patrick, and Andrea Vedaldi. 2019. Understanding deep networks via extremal perturbations and smooth masks. In *Proceedings of the IEEE International Conference on Computer Vision*. 2950–2958.
- [42] Ruth C Fong and Andrea Vedaldi. 2017. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision*. 3429–3437.
- [43] Ruigang Fu, Qingyong Hu, Xiaohu Dong, Yulan Guo, Yinghui Gao, and Biao Li. 2020. Axiom-based Grad-CAM: Towards Accurate Visualization and Explanation of CNNs. *ArXiv abs/2008.02312* (2020).
- [44] Keren Gaiger, Oren Barkan, Shir Tsipory-Samuel, and Noam Koenigstein. 2023. Not All Memories Created Equal: Dynamic User Representations for Collaborative Filtering. *IEEE Access* 11 (2023), 34746–34763.
- [45] Dvir Ginzburg, Itzik Malkiel, Oren Barkan, Avi Caciularu, and Noam Koenigstein. 2021. Self-supervised document similarity ranking via contextualized language models and hierarchical inference. *arXiv preprint arXiv:2106.01186* (2021).
- [46] Matthieu Guillaumin, Daniel Küttel, and Vittorio Ferrari. 2014. Imagenet auto-annotation with segmentation propagation. *International Journal of Computer Vision* 110, 3 (2014), 328–348.
- [47] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16000–16009.
- [48] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.

- [49] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 770–778.
- [50] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 770–778.
- [51] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*. 639–648.
- [52] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*. 173–182.
- [53] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. 2017. Densely Connected Convolutional Networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017), 2261–2269.
- [54] Peng-Tao Jiang, Chang-Bin Zhang, Qibin Hou, Ming-Ming Cheng, and Yunchao Wei. 2021. Layercam: Exploring hierarchical class activation maps for localization. *IEEE Transactions on Image Processing* 30 (2021), 5875–5888.
- [55] Hyungsik Jung and Youngrook Oh. 2021. Towards Better Explanations of Class Activation Mapping. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)* (2021), 1316–1324.
- [56] Andrei Kapishnikov, Tolga Bolukbasi, Fernanda Viégas, and Michael Terry. 2019. Xrai: Better attributions through regions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4948–4957.
- [57] Andrei Kapishnikov, Subhashini Venugopalan, Besim Avci, Ben Wedin, Michael Terry, and Tolga Bolukbasi. 2021. Guided integrated gradients: An adaptive path method for removing noise. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5050–5058.
- [58] Kundan Kumar, Rithesh Kumar, Thibault De Boissiere, Lucas Gestein, Wei Zhen Teoh, Jose Sotelo, Alexandre De Brebisson, Yoshua Bengio, and Aaron C Courville. 2019. Melgan: Generative adversarial networks for conditional waveform synthesis. *Advances in neural information processing systems* 32 (2019).
- [59] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. 2014. Microsoft COCO: Common Objects in Context. <http://arxiv.org/abs/1405.0312> cite arxiv:1405.0312Comment: 1) updated annotation pipeline description and figures; 2) added new section describing datasets splits; 3) updated author list.
- [60] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [61] Zhuang Liu, Hanzi Mao, Chaozheng Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. 2022. A ConvNet for the 2020s. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022), 11966–11976.
- [62] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*. 13–23.
- [63] Scott M. Lundberg and Su-In Lee. 2017. A Unified Approach to Interpreting Model Predictions. In *NIPS*.
- [64] Aravindh Mahendran and Andrea Vedaldi. 2016. Visualizing deep convolutional neural networks using natural pre-images. *International Journal of Computer Vision* 120, 3 (2016), 233–255.
- [65] Itzik Malkiel, Oren Barkan, Avi Caciularu, Noam Razin, Ori Katz, and Noam Koenigstein. 2020. RecoBERT: A catalog language model for text-based recommendations. *arXiv preprint arXiv:2009.13292* (2020).
- [66] Itzik Malkiel, Dvir Ginzburg, Oren Barkan, Avi Caciularu, Jonathan Weill, and Noam Koenigstein. 2022. Interpreting BERT-based text similarity via activation and saliency maps. In *Proceedings of the ACM Web Conference 2022*. 3259–3268.
- [67] Tomas Mikolov, Wen tau Yih, and Geoffrey Zweig. 2013. Linguistic regularities in continuous space word representations. In *NAACL-HLT*.
- [68] Vitali Petsiuk, Abir Das, and Kate Saenko. 2018. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421* (2018).
- [69] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*. 618–626.
- [70] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning Important Features Through Propagating Activation Differences. In *ICML*.
- [71] Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. 2016. Not Just a Black Box: Learning Important Features Through Propagating Activation Differences. *ArXiv abs/1605.01713* (2016).
- [72] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034* (2013).
- [73] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [74] K. Simonyan and A. Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *International Conference on Learning Representations (ICLR)*.
- [75] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. 2017. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825* (2017).
- [76] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. 2014. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806* (2014).
- [77] Suraj Srinivas and François Fleuret. 2019. Full-Gradient Representation for Neural Network Visualization. In *NeurIPS*.
- [78] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic Attribution for Deep Networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017*. 3319–3328.
- [79] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [80] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. 2019. Analyzing Multi-Head Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 5797–5808.
- [81] Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and Xia Hu. 2020. Score-CAM: Score-Weighted Visual Explanations for Convolutional Neural Networks. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)* (2020), 111–119.
- [82] Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and Xia Hu. 2020. Score-CAM: Score-weighted visual explanations for convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*. 24–25.
- [83] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural graph collaborative filtering. In *Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*. 165–174.
- [84] Shawn Xu, Subhashini Venugopalan, and Mukund Sundararajan. 2020. Attribution in scale and space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9680–9689.
- [85] Matthew D Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *European conference on computer vision*. Springer, 818–833.
- [86] Bolei Zhou, David Bau, Aude Oliva, and Antonio Torralba. 2018. Interpreting deep visual representations via network dissection. *IEEE transactions on pattern analysis and machine intelligence* (2018).
- [87] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. 2016. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2921–2929.