# ADVISE: ADaptive Feature Relevance and VISual Explanations for Convolutional Neural Networks

Mohammad Mahdi Dehshibi, Mona Ashtari-Majlan, Gereziher Adhane, David Masip, *Senior Member, IEEE*

**Abstract**—To equip Convolutional Neural Networks (CNNs) with explainability, it is essential to interpret how opaque models take specific decisions, understand what causes the errors, improve the architecture design, and identify unethical biases in the classifiers. This paper introduces ADVISE, a new explainability method that quantifies and leverages the relevance of each unit of the feature map to provide better visual explanations. To this end, we propose using adaptive bandwidth kernel density estimation to assign a relevance score to each unit of the feature map with respect to the predicted class. We also propose an evaluation protocol to quantitatively assess the visual explainability of CNN models. We extensively evaluate our idea in the image classification task using AlexNet, VGG16, ResNet50, and Xception pretrained on ImageNet. We compare ADVISE with the state-of-the-art visual explainable methods and show that the proposed method outperforms competing approaches in quantifying feature-relevance and visual explainability while maintaining competitive time complexity. Our experiments further show that ADVISE fulfils the sensitivity and implementation independence axioms while passing the sanity checks. The implementation is accessible for reproducibility purposes on https://github.com/dehshibi/ADVISE.

**Index Terms**—Convolutional Neural Network, Deep Learning, eXplainable AI

✦

## 1 INTRODUCTION

CONVOLUTIONAL Neural Networks (CNNs) have gained significant prominence with the potential to outperform expectations in various computer vision tasks such as image classification [12], [2], [3], object detection [49], semantic segmentation [31], image captioning [11], and human behaviour analysis [14]. However, this sub-symbolism (also known as the opaque or black-box model) is vulnerable to the underlying barrier of *explainability* in response to critical questions like how a particular trained model arrives at a decision, how certain it is about its decision, if and when it can be trusted, why it makes certain mistakes, and in which part of the learning algorithm or parametric space correction should take place [28], [4]. Explainability in CNNs is linked to post-hoc explainability [18] and, as proposed by Arrieta et al. [4], relies on model simplification [56], [36], [23], feature-relevance estimation [6], [33], [29], [38], visualisation [53], [30], [26], [39], [48], [22], and architectural modification [27], [15], [40] to convert a non-interpretable model into an explainable one.

While model simplification and architectural modification techniques have been used to make CNNs interpretable, their associated complexity grows as the number of layers and parameters increases. Furthermore, several studies [5], [34], [4] have shown that altering CNNs may result in the spontaneous appearance of a disentangled representation [17], [57], which is not only unrelated to the model's initial intention but also challenging to interpret. As a result, the emphasis in explaining CNNs has shifted toward feature-relevance and visualisation methods.

Feature visualisation has received much attention because human cognitive skills favour the understanding of visual data. However, feature visualisation methods do not necessarily provide a comprehensive level of explainability and interpretability. For instance, in Figure 1a, an identical image is fed into the VGG16 and Xception models, both of which outperform humans on ImageNet classification. Although both models have the exact top-1 prediction with one difference in top-5 prediction, the visual explanations are significantly different (see LIME and Cumulative Gradients in Figure 1) and cannot provide users with comprehensive information about how the models made the final decision. Therefore, several studies [6], [42], [29] focused on feature-relevance approaches, which provide an importance score to each feature for a specific input. However, the visual and feature-relevance explanations are not mutually exclusive when a feature-relevance method can be visualised as a saliency map [36], and a saliency map generated using class activation maps [53], [39], [22] can assign importance scores to each pixel.

In this paper, we propose a method for quantifying the feature-relevance and visualising the latent representations in CNNs. We revisit the relationships between feature maps[1] and their associated gradients by introducing ADaptive VISual EXplanation (ADVISE). **ADVISE** estimates the kernel density of gradients with an adaptive bandwidth for each unit in the feature map (see Figures 1e and 1k) to assign an importance score to each unit. Then, we calculate the cu-

● *Department of Computer Science, Universitat Oberta de Catalunya, 08018, Barcelona, Spain.*
*Corresponding author: M. M. Dehshibi (e-mail: mohammad.dehshibi@yahoo.com).*

1. The terms *feature map* and *activation map* are used interchangeable here since the former refers to a mapping of where a specific type of feature can be found in an image, and the latter is a mapping that relates to the activation of different areas of the image.
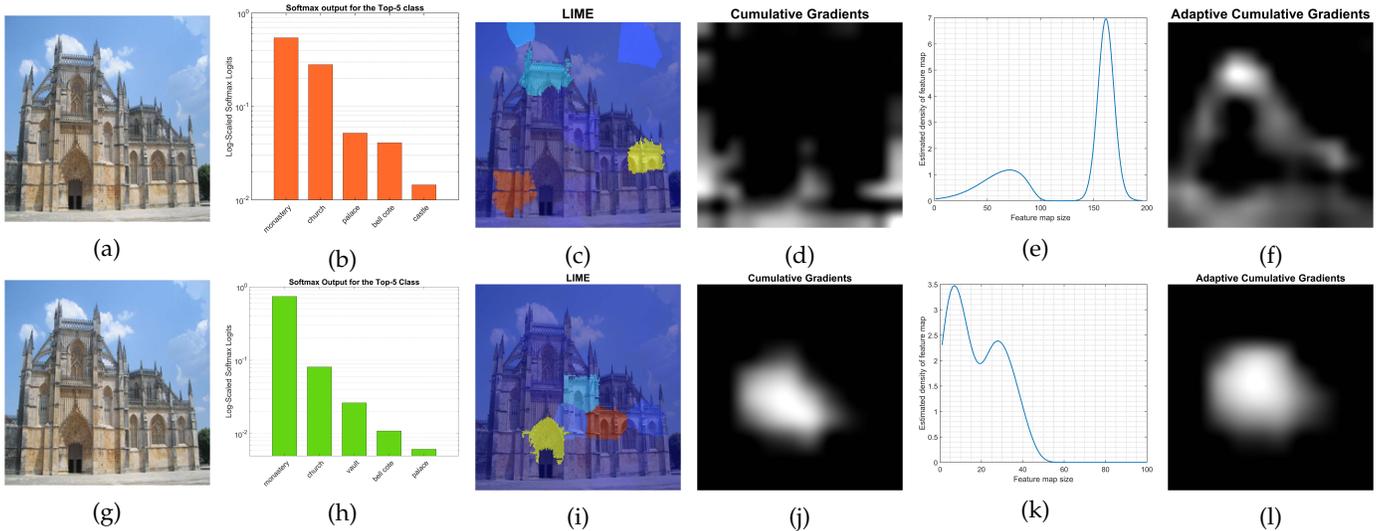
Fig. 1: The visual and feature-relevance explanations for VGG16 [43] (the first row) and Xception [13] (the second row) pretrained models on ILSVRC [37]. (a, g) The input image. (b, h) The output of the Log-scaled softmax Logits for the top-5 predicted classes. (c, i) Local explanations for the prediction of the input image based on LIME [36]. (d, j) Cumulative gradients of the last convolutional layer, where the feature map is scaled up to the resolution of the input image using bilinear interpolation. (e, k) Estimated density of the $k^{\text{th}}$ unit in feature map, which represents 2 peaks. (f, l) Adaptive cumulative gradients of units with 2 peaks in their estimated density.

mulative gradient of units with the same importance score for the class of interest to visualise the feature map. In this way, we simultaneously quantify the relevance of each unit and highlight how much the cumulative gradient of units influence the model's decision using the generated saliency map(s) (see Figures 1f and 1l). We use the proposed method to demonstrate that individual units are significantly more interpretable than cumulative linear combinations of gradient's units.

Our experiment is centred on the image classification task since it allows us to visualise adaptive cumulative gradient attributions and compare ADVISE with attention approaches that focus on global information. We use AlexNet [24], VGG16 [43], ResNet50 [19], and Xception [13], which were trained on the ILSVRC [37] in order to decide to which of 1000 classes each image belongs. However, unlike previous approaches, estimating the kernel density of gradients with the adaptive bandwidth can be applied to a wide range of deep learning models without requiring architectural changes or retraining.

The rest of this paper is organised as follows: Section 2 surveys the previous studies. The proposed method is detailed in Section 3. Section 4 presents experimental results. Finally, Section 5 concludes the paper.

## 2 LITERATURE REVIEW

As previously stated, explaining a model by the visualisation (i.e.,, explicit explainability) and feature-relevance (i.e.,, implicit explainability) are not mutually exclusive. In fact, visualisation techniques present complementary ways of visualising the output of feature relevance techniques to aid model interpretation.

In this context, the methods proposed to explain what CNNs learn can be categorised into three broad categories:

(1) those that rely on attention methods by generating class activation maps to interpret how the intermediate layers perceive the external world with respect to the target class without restricting the method to any specific input; (2) those that interpret the decision process using a top-down back-propagation strategy in which the output is mapped back in the input space to determine which parts of the input are discriminative for the output; (3) those that integrate importance over the attribution path and open up the axiomatic *sensitivity* and *implementation invariance* attributions for deep neural networks. These methods are amenable to intriguing visualisations and serve as a basis for discussing missingness in the feature space.

The main idea behind class activation mapping is to achieve class-specific importance for each location of an image by multiplying each feature map by its weight and performing a sum of all weighted feature map values at that location across all channels (units). Following this procedure, a ReLU operation is used to filter out negative activations. The method of calculating the weight for each feature map differs between different attention methods. CAM [58] obtained the weights from a single fully connected layer that produces the predictions, in which global average pooling is applied to the final convolutional feature maps. Grad-CAM [39] improved on CAM by applying class-specific gradients to each feature map at each location and averaging the gradients of each feature map unit as its weight. Grad CAM++ [10] generates a visual explanation for the corresponding class label by using a weighted combination of the positive partial derivatives of the last convolutional layer feature maps with respect to a specific class score as weights. Score-CAM [50] eliminates gradient dependence by masking the input image with the activation map generated with respect to the target class at different network layers and passing it through the network to obtain

the prediction score. Finally, the weight for each feature map is calculated by the normalised sum of the obtained scores. Layer-CAM [22] utilises the backward class-specific gradients, in which the gradient with respect to the class of interest is calculated for each unit in the feature map, and the units with positive gradient values are used as weights.

Zhang et al. [54] introduced a top-down back-propagation approach to compute neuron significance towards a model that passes signals in the network downwards based on a probabilistic Winner-Take-All model. Fong and Vedaldi [16] and Cao et al. [9] learn a perturbation mask that significantly influences the model's output by backpropagating the error signals through the model. Zhou et al. [59] extracted fine-detailed class instance activation maps by back-propagating the peak values as top signals to the network downwards in a Winner-Take-All manner. However, the generated maps are less faithful than those produced by CAM-based methods, and such a top-down procedure is complex and computationally expensive.

Sundararajan et al. [47] introduced integrated gradients as a way to quantify a neural network's feature-relevance when making a prediction for a given data point and brought up the concept of missingness in the feature space as a critical interpretability concept. Sturmfels et al. [45] later discussed the influence of choosing a baseline input for the integrated gradients. Bau et al. [5] introduced network dissection to show that individual units are significantly more interpretable than random linear combinations of units. They consider each unit as a concept detector to further evaluate them for semantic segmentation and quantify the interpretability of CNN latent representations. While these studies proposed solutions to fulfil the sensitivity and implementation invariance axioms, they either required the definition of a baseline input, relied on a threshold derived from the training data set, or limited the solution to a binary segmentation task, all of which failed the sanity checks [1], [44].

## 3 ADVISE: ADAPTIVE VISUAL EXPLANATION

Formally, let $f(I;\theta) = \mathbb{E}[y^c|I;\theta]$ represents a CNN that classifies images, and $\theta$ denotes its parameters. For the input image $I \in \mathbb{R}^{H \times W \times 3}$, $y^c$ is the score for the predicted class $c$, where $H$ and $W$ denote the height and width of $I$, respectively. Let $A \in \mathbb{R}^{U \times V \times K}$ denotes an activation map in the $f$, where $A^k$ represents the $k^{\text{th}}$ feature map in $A$, and $U$, $V$, and $K$ denote the height, width, and the number of units of $f$, respectively. The gradient of the predicted score $y^c$ with respect to the spatial location $(i, j)$ in the feature map $A$ can be obtained by $\frac{\partial y^c}{\partial A_{i,j}}$.

Although the visualisation methods that calculate cumulative gradients (i.e.,, a linear weighted summation on all feature maps in $A$) preserve implementation invariance, they do not satisfy sensitivity because they assume a stationary rate variation in the gradients. To preserve both the implementation invariance and sensitivity axioms [47], we propose computing $\phi_k(A)$, which assigns an importance score to the $k^{\text{th}}$ unit in the feature map $A$, indicating how much that feature contributes to the network decision. Then we calculate the linear weighted sum of the feature maps in $A$ that have the same importance score.

Kernel density estimation (KDE) is a conventional non-parametric signal processing approach for estimating the probability density function of data with an unknown underlying distribution [35]. Let $(a_1, a_2, \cdots, a_n)$ be the value of the independent distributed gradients in the $k^{\text{th}}$ unit of $A$ that were flattened. The gradient values are changed with respect to the input image $I$ and stacked to form a raw density as in Eq. 1

$$x_a = \frac{1}{n} \sum_{i=1}^{n} \delta(a - a_i), \tag{1}$$

where $n = U \times V$, and $\{a_i\}_{i=1}^n$ is represented by the Dirac delta function $\delta(a)$. The kernel density estimate is obtained by convolving a kernel $\mathcal{H}_{\omega_a}$ with the variable bandwidth $\omega_a$ to the raw density $x_a$ using Eq. 2.

$$\widehat{\lambda}_a = \int x_{a-s} \mathcal{H}_{\omega_a}(s) \, \mathrm{d}s. \tag{2}$$

where $\omega_a$ is selected as a fixed bandwidth optimised in a local interval, and the integral $\int$ that does not specify bounds refers to $\int_{-\infty}^{\infty}$. The mean integrated squared error (MISE) [8] is a well-known goodness-of-fit metric for optimising the estimated density $\widehat{\lambda}_a$ to be as close to the unknown underlying density $\lambda_a$ as possible. Motivated by [41], we introduce the adaptive MISE (AMISE) criterion at gradient $a$ to select an interval length for local optimisation, determine the goodness-of-fit, and regulate the shape of the function $\lambda_a$ as in Eq. 3.

$$\text{AMISE} = \int \mathbb{E} \left( \widehat{\lambda}_u - \lambda_u \right)^2 \rho_W^{u-a} \, \mathrm{d}u, \tag{3}$$

where $\mathbb{E}$ is the expected $L_2$ loss function, $\widehat{\lambda}_u = \int x_{u-s} \mathcal{H}_\omega(s) \, \mathrm{d}s$ is the estimated density with a fixed bandwidth $\omega$, and $\rho_W^{u-a}$ is a weight function that locates the integration of the squared error in a particular interval $W$ centring at $a$. To minimise AMISE, we introduce the adaptive cost function with respect to $a$ by subtracting the irrelevant term for the choice of $\omega$ as in Eq. 4.

$$C_n^a(\omega, W) = \text{AMISE} - \int \lambda_u^2 \rho_W^{u-a} \, \mathrm{d}u. \tag{4}$$

The optimal fixed bandwidth $\omega^*$ is obtained as a minimiser of the estimated cost function that is presented in Eq. 5:

$$\hat{C}_n^a(\omega, W) = \frac{1}{n^2} \sum_{i,j} \psi_{\omega,W}^a(a_i, a_j)$$
$$- \frac{2}{n^2} \sum_{i \neq j}^{n} \mathcal{H}_\omega(a_i - a_j) \rho_W^{a_i - a}, \tag{5}$$

where $\psi_{\omega,W}^a$ is given in Eq. 6.

$$\psi_{\omega,W}^a(a_i, a_j) = \int \mathcal{H}_\omega(u - a_i) \mathcal{H}_\omega(u - a_j) \rho_W^{u-a} \, \mathrm{d}u. \tag{6}$$

Since the optimal bandwidth $\omega^*$ varies with the length of $W$, we select an interval length of $\frac{\omega^*}{\gamma}^2$ that scales with

---

2. $\frac{\omega^*}{\gamma} = n$ is used in our experiment

the optimal bandwidth. Here, $\gamma$ is a smoothing parameter, with $\gamma << 1$ causing the variable bandwidth to fluctuate slightly, and $\gamma \sim 1$ causing the variable bandwidth to fluctuate significantly. In our experiments, we consider the $[0, 1]$ interval and use the Nadaraya-Watson kernel regression [32] to obtain the variable bandwidth $\omega_a^\gamma$ using Eq. 7

$$\omega_a^\gamma = \int \rho_{W_s^\gamma}^{a-s} \bar{\omega}_s^\gamma \, \mathrm{d}s \Big/ \int \rho_{W_s^\gamma}^{a-s} \, \mathrm{d}s \, . \tag{7}$$

where $W_a^\gamma$ and $\bar{\omega}_a^\gamma$ represent the interval length and fixed bandwidth at $a$, respectively. Although the variable bandwidth $\omega_a^\gamma$ is derived from the same data, the use of different $\gamma$ results in varying degrees of smoothness. In this way, the cost function for the variable bandwidth selected with $\gamma$ is obtained using Eq. 8.

$$\hat{C}_n(\gamma) = \int_0^1 \hat{\lambda}_a^2 \, \mathrm{d}a - \frac{2}{n^2} \sum_{i \neq j} \mathcal{H}_{\omega_{a_i}^\gamma}(a_i - a_j), \tag{8}$$

where $\hat{\lambda}_a = \int x_{a-s} \mathcal{H}_{\omega_a^\gamma}(s) \, \mathrm{d}s$ is an estimated rate, with the variable bandwidth $\omega_a^\gamma$. The integral is calculated numerically with the stiffness constant $\gamma^* = \frac{\sqrt{5}+1}{2}$ that minimises Eq. 8. In this study, we use the Gauss density function which is expressed in Eq. 9.

$$\mathcal{H}_{\omega^\gamma}(s) = \frac{1}{\sqrt{2\pi\omega^\gamma}} \exp\left(-\frac{s^2}{2(\omega^\gamma)^2}\right), \tag{9}$$

Figure 2a depicts one of the activation map units in the VGG16 model's final convolution layer for the input image in Figure 1a. Figure 2b shows the difference between the underlying gradient value distribution (grey area) at that unit and the estimated density of gradient values (solid red line) using the proposed variable bandwidth kernel density estimation.
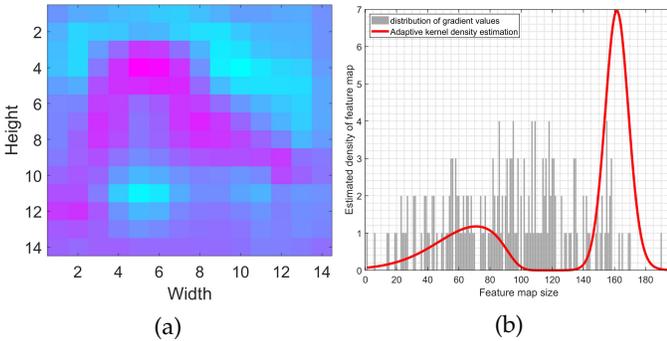


(a)                              (b)

Fig. 2: (a) The $265^{\text{th}}$ unit of the activation map in the last convolution layer of the VGG16 model for the input image in Figure 1a, where gradient values are mapped to colours in the 'cool' colour map for better visualisation. (b) Estimated kernel density with variable bandwidth (solid red line) using Eq. 8. The grey area represents the underlying distribution of gradient values in the $265^{\text{th}}$ unit of the activation map.

The proposed scoring method that assigns an importance score to the $k^{\text{th}}$ unit in the feature map as well as the visualisation approach (ADVISE) are summarised in Algorithm 1.

---

**Algorithm 1:** ADaptive VISual Explanation.

**Input** : $A^{U \times V \times K}$ – Feature map, also known as *activation map* in CNNs.
$y^c$ – predicted class.
[row, col] – size of input image.
**Output:** $\phi_k(A)$ – Importance score for units in $A$.
ADVISE – Feature saliency map(s).

**1 for** $k \leftarrow 1$ **to** $K$ **do**
**2**   $\{a_i\}_{i=1}^n \leftarrow$ `flatten(A)`;
      // $n = U \times V$.
**3**   $\phi_k(A) =$
      `findPeaks`$(\int_0^1 \hat{\lambda}_a^2 \, \mathrm{d}a - \frac{2}{n^2} \sum_{i \neq j} \mathcal{H}_{\omega_{a_i}^\gamma}(a_i - a_j))$;
**4 end**

**5** $g = \frac{\partial y^c}{\partial A}$;
**6 for** $i \leftarrow \min(\phi_k(A))$ **to** $\max(\phi_k(A))$ **do**
**7**   idx $\leftarrow$ `find`$(\phi_k(A) == i)$;
**8**   $\tilde{A}_i = A(:, :, \text{idx})$;
**9**   $\tilde{w}_i^c = \frac{1}{n} \sum_U \sum_V g(:, :, \text{idx})$;
**10**   map$_i = \text{ReLU}\left(\sum_{j=1}^{|\text{idx}|} \tilde{w}_{i,j}^c \cdot \tilde{A}_{i,j}\right)$;
      // $|\bullet|$ is the cardinality of $\bullet$
**11**   ADVISE$_i =$ `resize(map`$_i$`, [row, col], bc)`;
      // 'bc' is bicubic interpolation
**12 end**

**13 return** $\phi_k(A)$, ADVISE

---

Figure 3 shows outputs of the proposed method using AlexNet [24], VGG16 [43], ResNet50 [19], and Xception [13], which were trained on the ILSVRC [37].

The results of scoring function $\phi_k(A)$ and the saliency maps generated by ADVISE can highlight three key points. (1) Not all feature map units can contribute equally to the model's prediction, and some of these units may be misleading in some instances (see Figure 3a). (2) As Bau et al. [5] pointed out, CNNs trained for a specific purpose may encounter the emergence of disentangled representations unrelated to the model's initial intention, complicating interpretation (see Figure 3b). As a result, quantifying feature-relevance in conjunction with visualisation can provide adequate answers for users, particularly neural network designers, to underlying questions such as how certain the model is about its decision, if and when it can be trusted, why it makes inevitable mistakes, and in which part of the learning algorithm or parametric space correction should take place. (3) In scenarios such as transfer learning, this mutual explainability approach assists designers in determining which layers should be frozen to achieve better and faster convergence, specifically when the feature map shows less divergence (see Figures 3c and 3d). In Section 4, where we introduce quantitative metrics to compare the visualisation approach with the competing ones, we will delve into greater depth on these points.

## 4 EXPERIMENTS

The proposed method for quantifying feature relevance is applicable to a variety of deep networks. However, we centre our experiments on the image classification task since
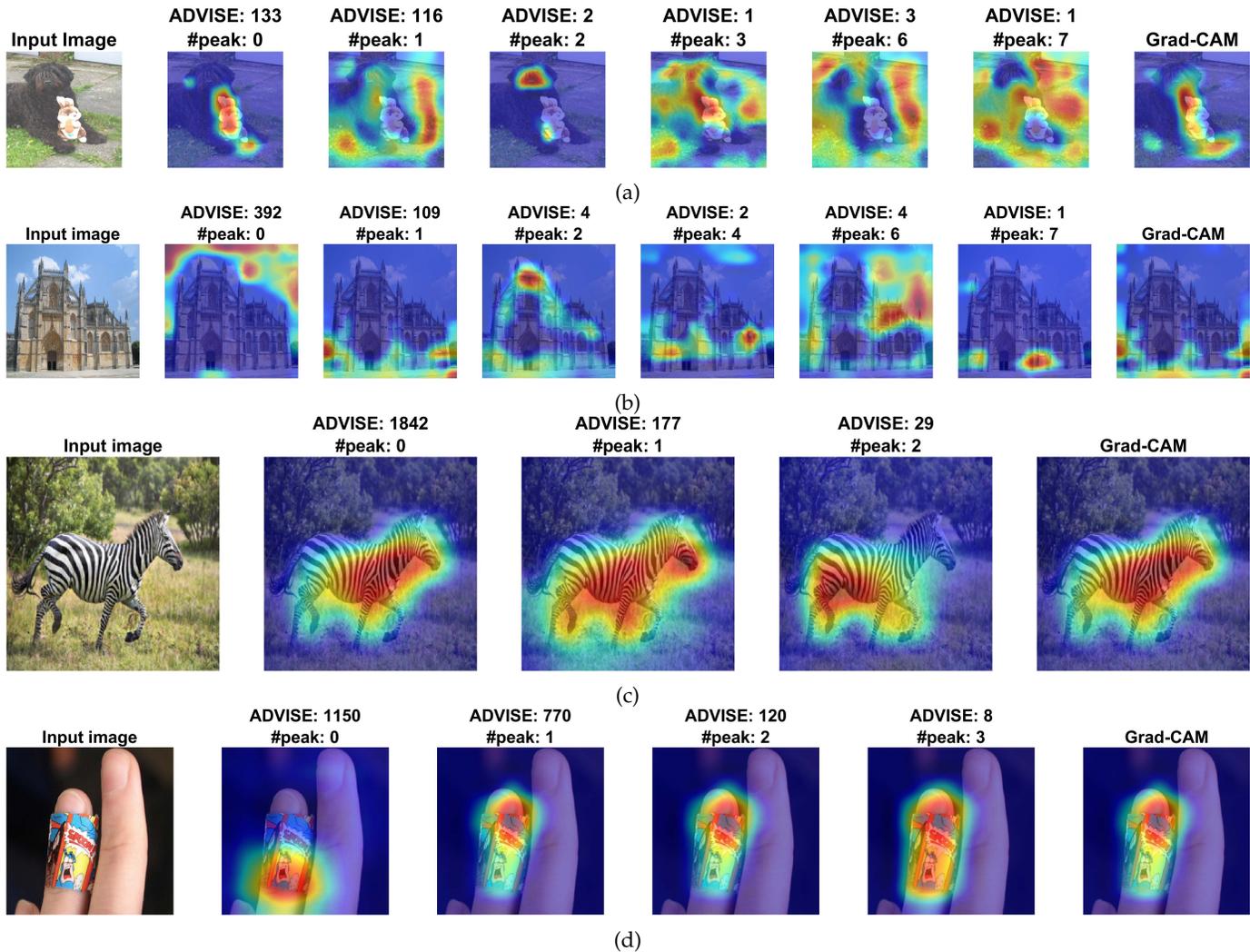
Fig. 3: The outputs of ADVISE and Grad-CAM [39] are compared for four images fed into the pretrained AlexNet [24], VGG16 [43], ResNet50 [19], and Xception [13] models on ILSVRC [37]. The use of $\phi_k(A)$ on the estimated kernel density and ADVISE show that in the explainability of (a) AlexNet prediction ('Bernese mountain dog'), two units with two peaks work better than Grad-CAM that requires 1000 units, (b) VGG16 prediction ('monastery'), four units with six peaks contribute more than Grad-CAM that requires 512 units, (c) ResNet50 prediction ('Zebra'), 177 units with one peak outperform Grad-CAM, which requires 2048 units, and (d) Xception prediction ('band aid'), eight units with three peaks perform better than Grad-CAM which utilises 2048 units.

it allows us to visualise adaptive cumulative gradient attributions and compare ADVISE with attention approaches that focus on global information. ADVISE is tested on a subset of ILSVRC [37] with 3,000 images using pretrained AlexNet [24], VGG16 [43], ResNet50 [19], and Xception [13] models.

In the absence of ground-truth discriminative features for a trained CNN [25], objectively identifying which method delivers the best approximation to the usefulness and satisfaction of explanations is still in its early stages. Furthermore, the community has not yet reached a consensus on the impact of explanations on the model's performance, trust, and reliance. A natural assumption is that a well-trained model would make predictions based on the features from the object itself [4]. With this assumption and following quantitative metrics that are used to evaluate image retrieval methods and saliency models, we present

a novel evaluation protocol for the visual explanation approaches.

### 4.1 Evaluation Metrics

**(1) Class Sensitivity (CS):** it measures the similarity of saliency maps generated with respect to the top two class scores predicted by the model. We use Pearson's Correlation Coefficient to measure CS as in Eq. 10.

$$\text{CS} = \frac{\text{cov}\left(E(f,I)^{c_1}, E(f,I)^{c_2}\right)}{\sigma\left(E(f,I)^{c_1}\right) \times \sigma\left(E(f,I)^{c_2}\right)}. \qquad (10)$$

where $E$, cov, and $\sigma$ denote the explanation map, covariance, and standard deviation, respectively. A good explanation method should have a score near to or below zero, while a score outside the $[-0.5, 0.5]$ range implies that the correlation between two maps is not statistically significant.

**(2) Hit:** it is a proxy that indicates if the model can retrieve the target class $c$ in its top-5 prediction when it just sees the explanation map and not the entire image. This proxy is formulated in Eq. 11.

$$\text{Hit} = \begin{cases} 1 & : N_I \cap M_{I \odot E(f,I)^c} \\ 0 & : \text{otherwise} \end{cases} \quad (11)$$

where $N_I$ is the index of the predicted class $c$ by the model when it just sees the input image as input, and $M_{I \odot E(f,I)^c}$ is a set including the top-5 index of the predicted class when the model sees the explanation map. Here, $\odot$ is the Hadamard product.

**(3) Average Drop (AD):** it measures the average percentage drop in confidence for the target class $c$ when the explanation map $(I \odot E(f,I)^c)$ is fed to the model instead of the input image $I$. This metric is defined in Eq. 12, where lower is better.

$$\text{AD} = \max\left(0, (y^c - o^c)\right)/y^c \quad (12)$$

where $o^c$ is the predicted score by model to which the the explanation map is fed.

**(4) Structural similarity index (SSIM):** it is a perception-based measure that considers image degradation as a perceived change in structural information while also considering crucial perceptual phenomena [51]. In this context, SSIM measures the structural similarity index between the input image masked by the explanation map and the input image as the reference. This metric returns a value in $(0, 1]$, where the higher is better, and is formulated in Eq. 13.

$$\text{SSIM}(I, \tilde{I}) = \frac{(2\mu_I \mu_{\tilde{I}} + e_1)(2\text{cov}(I, \tilde{I}) + e_2)}{(\mu_I^2 + \mu_{\tilde{I}}^2 + e_1)(\sigma_I^2 + \sigma_{\tilde{I}}^2 + e_2)}. \quad (13)$$

where $\tilde{I} = I \odot E(f,I)^c$, and $\mu$ and $\sigma$ are the average and variance, respectively. In order to stabilise the division with weak denominator, $e_1 = (0.01 \cdot L)^2$ and $e_2 = (0.03 \cdot L)^2$ are used, where $L$ denotes the dynamic range of the pixel values and is set to 255 in this study.

**(5) Feature similarity index (FSIM):** it uses phase congruency and gradient magnitude, which reflect complementary components of visual image quality, to measure local image quality. This metric also includes a saliency measure for the image gradient feature, which weights each pixel's contribution to the overall quality score. This metric returns a value in $(0, 1]$, where the higher is better, and the complete mathematical formulation is given in [55].

**(6) Mean squared error (MSE):** it is the second error moment and measures the average squared difference between the input image masked by the explanation map and the input image as the reference as in Eq. 14.

$$\text{MSE}(I, \tilde{I}) = \frac{1}{HW} \sum_{i=1}^{H} \sum_{j=1}^{W} \left(I_{i,j} - \tilde{I}_{i,j}\right)^2. \quad (14)$$

**(7) AVerage eXplainability (AVX):** it measures the harmonic mean of AD, SSIM, FSIM, and MSE and returns a value in $[0, 1]$ to ease of comparison as defined in Eq. 15.

$$\text{AVX} = 4 \left( \frac{1}{1 - \text{AD}} + \frac{1}{\text{SSIM}} + \frac{1}{\text{FSIM}} + \frac{1}{1 - \text{MSE}} \right)^{-1} \quad (15)$$

Recall that we defined two proxies, CS and Hit, which allow us to adjust AVX. If Hit = 0 and CS $\in [-0.5, 0.5]$, we define a penalty coefficient $\Delta = 1 - |y^c - o^c|$ and multiply AD, SSIM, FSIM, and MSE by $\Delta$ before measuring the harmonic mean. If Hit = 0 and CS $\notin [-0.5, 0.5]$, we set AD to 1, SSIM to 0, FSIM to 0, and MSE to 1.

## 4.2 Experimental result

Table 1 shows the comparison of the ADVISE with Grad-CAM [39], Grad-CAM++ [10], Score-CAM [50], and Layer-CAM [22] visualisation methods on AlexNet [24], VGG16 [43], ResNet50 [19], and Xception [13] pretrained models on ILSVRC [37]. Despite having a higher performance in classifying ILSVRC than the AlexNet, VGG16, and ResNet50, the Xception model has a lower efficiency in the visual explanation, according to the AVX metric.

In our quest for this AVX decline in Xception, we examined the saliency maps produced by the ADVISE in shallow, middle, and deep layers (see an example in Figure 4). We observed that the saliency maps in the shallow and middle layers highlight low-level visual features distributed across the image, such as edges and blobs. The Xception model, on the other hand, focuses on the centre of a scene in the deep layer, whereas the other models look at different locations. This focus is known as the centre bias in saliency studies [7], [52], where most studies revealed that observers prefer to look more often at the centre of the image than at the edges. However, the Xception model's tendency toward centre bias is a double-edged sword. While it is more aligned with human cognitive skills for perceiving visual data, as explained by [45], the centre of mass of the saliency map is the Achilles Heel of many visual explanation methods, with path attribution methods offered to address it [47] but failing the sanity checks [1].

So what should be done? Although the proposed method and quantitative metrics, which are supported by best practices, can evaluate the performance of different models in visual explanation, we still have a fundamental problem with the lack of *ground-truth explanations*. In fact, we aim to determine which methods best explain our model without knowing how it works. Evaluating supervised models is relatively straightforward since we have a test set. However, evaluating explanations is difficult since we do not exactly know how our model works and do not have the ground-truth for a fair comparison.

### 4.2.1 Ablation study

The gradient quantifies how much a change in each input dimension affects $f$ prediction in a narrow area around the input. Keeping this in mind, our ablation study is composed of two parts: (1) we ablate the input image by randomly replacing pixels with the salt and pepper noise counterparts; (2) we remove ReLU at the same time to explore the effect of negative gradients on scoring the feature map units and the visual explanation. To do this, all 3,000 images selected from ILSVRC are ablated using the noise density of $\delta = [0.025, 0.05, 0.075, 0.1, 0.125, 0.15, 0.175, 0.2, 0.225]$. Figure 5a depicts an ablated image, and Figure 5b–5e shows the proposed method's performance compared with other visual explanation methods.

TABLE 1: The comparison of the ADVISE with Grad-CAM, Grad-CAM++, Score-CAM, and Layer-CAM visualisation methods on AlexNet, VGG16, ResNet50, and Xception pretrained models on ILSVRC.
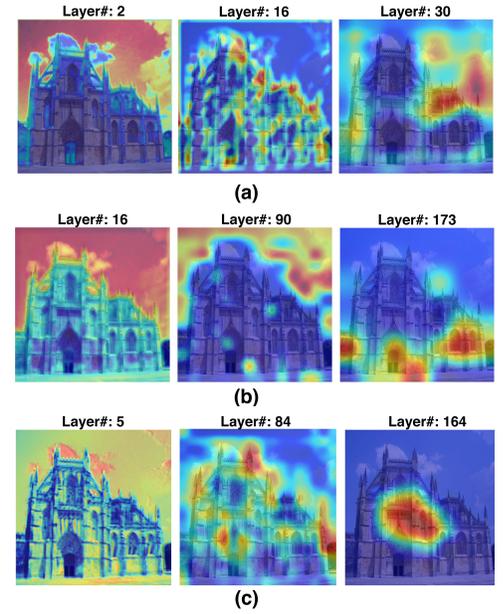
| Architecture | Method | Metrics | | | | | | Time (s) | |
|---|---|---|---|---|---|---|---|---|---|
| | | Peak range | AD ↓ | SSIM ↑ | FSIM ↑ | MSE ↓ | AVX ↑ | GPU/Parallel | CPU |
| AlexNet [24] | ADVISE | 0 − 8 | **0.26** | **0.14** | **0.38** | **0.14** | **0.28** | **0.69** | 30.3 |
| | Grad-CAM | N/A | 0.39 | 0.05 | 0.26 | 0.32 | 0.13 | 1.06 | **1.64** |
| | Grad-CAM++ | N/A | 0.38 | 0.06 | 0.27 | 0.32 | 0.17 | 1.16 | 2.14 |
| | Score-CAM | N/A | 0.37 | 0.06 | 0.28 | 0.31 | 0.17 | 1.18 | 2.60 |
| | Layer-CAM | N/A | 0.33 | 0.07 | 0.31 | 0.28 | 0.19 | 1.48 | 3.33 |
| | LIME | N/A | 0.39 | 0.05 | 0.26 | 0.32 | 0.13 | 5.71 | 11.85 |
| VGG16 [43] | ADVISE | 0 − 7 | **0.26** | **0.14** | **0.40** | **0.15** | **0.29** | **1.56** | 6.91 |
| | Grad-CAM | N/A | 0.38 | 0.06 | 0.26 | 0.29 | 0.15 | 1.88 | **2.66** |
| | Grad-CAM++ | N/A | 0.38 | 0.07 | 0.27 | 0.28 | 0.19 | 2.01 | 3.36 |
| | Score-CAM | N/A | 0.37 | 0.09 | 0.30 | 0.29 | 0.22 | 2.21 | 3.87 |
| | Layer-CAM | N/A | 0.32 | 0.09 | 0.34 | 0.27 | 0.23 | 2.66 | 4.24 |
| | LIME | N/A | 0.38 | 0.06 | 0.26 | 0.29 | 0.15 | 22.18 | 57.95 |
| ResNet50 [19] | ADVISE | 0 − 5 | **0.26** | **0.15** | **0.43** | **0.17** | **0.31** | **1.46** | **6.37** |
| | Grad-CAM | N/A | 0.33 | 0.10 | 0.34 | 0.24 | 0.23 | 6.22 | 7.77 |
| | Grad-CAM++ | N/A | 0.36 | 0.11 | 0.35 | 0.24 | 0.26 | 6.62 | 8.56 |
| | Score-CAM | N/A | 0.35 | 0.11 | 0.37 | 0.22 | 0.27 | 7.02 | 9.18 |
| | Layer-CAM | N/A | 0.32 | 0.12 | 0.39 | 0.21 | 0.29 | 7.51 | 11.18 |
| | LIME | N/A | 0.33 | 0.10 | 0.34 | 0.24 | 0.23 | 7.68 | 31.61 |
| Xception [13] | ADVISE | 0 − 6 | **0.43** | **0.12** | **0.37** | **0.31** | **0.24** | **4.20** | 16.38 |
| | Grad-CAM | N/A | 0.68 | 0.04 | 0.20 | 0.59 | 0.10 | 5.92 | **8.12** |
| | Grad-CAM++ | N/A | 0.65 | 0.04 | 0.21 | 0.59 | 0.11 | 6.03 | 9.10 |
| | Score-CAM | N/A | 0.64 | 0.05 | 0.21 | 0.57 | 0.13 | 6.56 | 9.70 |
| | Layer-CAM | N/A | 0.57 | 0.08 | 0.27 | 0.49 | 0.19 | 7.07 | 10.34 |
| | LIME | N/A | 0.68 | 0.04 | 0.20 | 0.59 | 0.10 | 26.31 | 90.31 |



Fig. 4: ADVISE outputs for shallow, middle, and deep layers of (a) VGG16, (b) ResNet50, and (c) Xception pretrained models on ILSVRC.

While the AVX value of the ADVISE and other visual explanation methods degrades due to incorporating negative gradients and ablating the input images, the proposed feature scoring method, unlike other methods, could meet the sensitivity axiom [47] in this classification task because the AVX never reached 0. However, we should mention that the pitfall of the ablation test is that if we artificially ablate pixels in an image, we end up with inputs that do not belong to the original data distribution. The question of whether or not users should feed their models with inputs that are not part of the initial training distribution is still being debated [20], [46], [21].

## 5 CONCLUSION

The significant achievement of Convolutional Neural Networks (CNNs) has resulted in a torrent of computer vision applications. Autonomous systems that can perceive, learn, decide, and act independently are on the horizon for these continuous breakthroughs. However, the incapacity of current approaches to adequately explain their decisions and actions to users limits their effectiveness. Therefore, CNNs must be equipped with the ability to explain their reasoning, characterise their strengths and shortcomings, and convey an understanding of how they will behave in the future. In this study, we have introduced ADVISE, a new explainability method that could quantify and leverage the relevance of each unit of the feature map to provide better visual explanations in CNNs. To this end, we have proposed a method to estimate the kernel density of gradients with an adaptive bandwidth for each unit in the feature map in order to calculate the number of peaks as the unit's relevance score. The cumulative gradient of units with the same relevance score for the class of interest was then calculated to visualise the latent representations in CNNs. We have also proposed a protocol for evaluating the visual explainability of CNN models quantitatively.

In our experiments, we used AlexNet, VGG16, ResNet50, and Xception pretrained on ILSVRC. We have compared ADVISE with the state-of-the-art visual explainable methods and showed that our proposed method outperformed competing approaches in quantifying feature-relevance and visual explainability while maintaining competitive time complexity. Our experiments further demonstrated that ADVISE meets the *sensitivity* and *implementation independence* axioms while passing the sanity checks.

It is worth mentioning that different metrics have been proposed to evaluate interpretability methods, each with its own set of pros and cons. This lack of consensus on evaluating interpretability methods is related to the fact that we do not know how exactly and transparently our model works and have no specific ground truth against which to compare it. As a result, more experiments on various computer vision tasks and other applications that benefit from the use of deep neural network architectures are required to demonstrate that ADVISE can meet a range of metrics for evaluating interpretability, as we intend to do in the future work.

## REFERENCES

[1] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 1–11. Curran Associates, Inc., 2018. 3, 6

[2] Gereziher Adhane, Mohammad Mahdi Dehshibi, and David Masip. A deep convolutional neural network for classification of aedes albopictus mosquitoes. *IEEE Access*, 9:72681–72690, 2021. 1
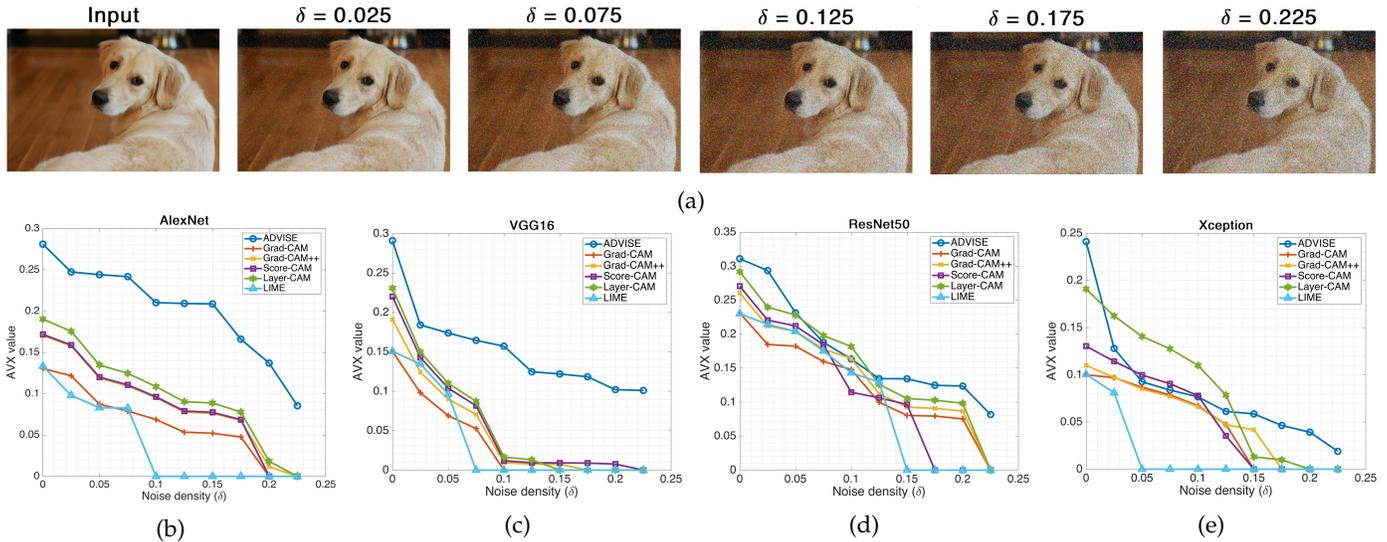
Fig. 5: (a) An ablated image by randomly replacing pixels with the salt and pepper noise with the noise density of $\delta = [0.025, 0.075, 0.125, 0.175, 0.225]$. (b-e) Changes in the performance of the ADVISE and five additional visual explanation methods in AlexNet, VGG16, ResNet50, and Xception pretrained models on ILSVRC as a function of (AVX, $\delta$).

[3] Gereziher Adhane, Mohammad Mahdi Dehshibi, and David Masip. On the use of uncertainty in classifying aedes albopictus mosquitoes. *IEEE Journal of Selected Topics in Signal Processing*, pages 1–11, 2021. 1

[4] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115, 2020. 1, 5

[5] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3319–3327, 2017. 1, 3, 4

[6] Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, Klaus-Robert Müller, and Wojciech Samek. Layer-wise relevance propagation for neural networks with local renormalization layers. In *Artificial Neural Networks and Machine Learning – ICANN 2016*, pages 63–71. Springer International Publishing, 2016. 1

[7] Ali Borji and James Tanner. Reconciling Saliency and Object Center-Bias Hypotheses in Explaining Free-Viewing Fixations. *IEEE Transactions on Neural Networks and Learning Systems*, 27(6):1214–1226, 2015. 6

[8] Adrian W. Bowman. An alternative method of cross-validation for the smoothing of density estimates. *Biometrika*, 71(2):353–360, 1984. 3

[9] Chunshui Cao, Xianming Liu, Yi Yang, Yinan Yu, Jiang Wang, Zilei Wang, Yongzhen Huang, Liang Wang, Chang Huang, Wei Xu, Deva Ramanan, and Thomas S. Huang. Look and Think Twice: Capturing Top-Down Visual Attention with Feedback Convolutional Neural Networks. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 2956–2964, 2015. 3

[10] Aditya Chattopadhay, Anirban Sarkar, Prantik Howlader, and Vineeth N. Balasubramanian. Grad-CAM++: Generalized Gradient-Based Visual Explanations for Deep Convolutional Networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 839–847, 2018. 2, 6

[11] Long Chen, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, Wei Liu, and Tat-Seng Chua. Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6298–6306, 2017. 1

[12] Weijie Chen, Di Xie, Yuan Zhang, and Shiliang Pu. All you need is a few shifts: Designing efficient convolutional neural networks for image classification. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7234–7243, 2019. 1

[13] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1800–1807, 2017. 2, 4, 5, 6, 7

[14] Mohammad Mahdi Dehshibi, Bita Baiani, Gerard Pons, and David Masip. A deep multimodal learning approach to perceive basic needs of humans from instagram profile. *IEEE Transactions on Affective Computing*, pages 1–13, 2021. 1

[15] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):677–691, 2017. 1

[16] Ruth C. Fong and Andrea Vedaldi. Interpretable Explanations of Black Boxes by Meaningful Perturbation. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 3429–3437, 2017. 3

[17] Abel Gonzalez-Garcia, Davide Modolo, and Vittorio Ferrari. Do semantic parts emerge in convolutional neural networks? *International Journal of Computer Vision*, 126(5):476–494, 2018. 1

[18] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)*, 51(5):93:1–93:42, 2018. 1

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 2, 4, 5, 6, 7

[20] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A Benchmark for Interpretability Methods in Deep Neural Networks. In *Advances in Neural Information Processing Systems*, pages 1–12. Curran Associates, Inc., 2019. 7

[21] Dominik Janzing, Lenon Minorics, and Patrick Blöbaum. Feature relevance quantification in explainable AI: A causal problem. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, pages 2907–2916. PMLR, 2020. 7

[22] Peng-Tao Jiang, Chang-Bin Zhang, Qibin Hou, Ming-Ming Cheng, and Yunchao Wei. Layercam: Exploring hierarchical class activation maps for localization. *IEEE Transactions on Image Processing*, 30:5875–5888, 2021. 1, 3, 6

[23] Sangwon Kim, Mira Jeong, and Byoung Chul Ko. Lightweight surrogate random forest support for model simplification and feature relevance. *Applied Intelligence*, pages 1–11, 2021. 1

[24] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105. Curran Associates, Inc., 2012. 2, 4, 5, 6, 7

[25] Xiao-Hui Li, Yuhan Shi, Haoyang Li, Wei Bai, Caleb Chen Cao, and Lei Chen. An Experimental Study of Quantitative Evaluations on Saliency Methods. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 3200––3208. Association for Computing Machinery, 2021. 5

[26] Yixuan Li, Jason Yosinski, Jeff Clune, Hod Lipson, and John E. Hopcroft. Convergent learning: Do different neural networks

[26] learn the same representations? In *Forth International Conference on Learning Representations, ICLR*, pages 196–212, 2016. 1

[27] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. In *Second International Conference on Learning Representations, ICLR*, pages 1——10, 2014. 1

[28] Zachary C. Lipton. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31——57, 2018. 1

[29] Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 4768–4777. Curran Associates Inc., 2017. 1

[30] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5188–5196, 2015. 1

[31] Sachin Mehta, Mohammad Rastegari, Anat Caspi, Linda Shapiro, and Hannaneh Hajishirzi. Espnet: Efficient spatial pyramid of dilated convolutions for semantic segmentation. In *Computer Vision – ECCV 2018*, pages 561–580. Springer International Publishing, 2018. 1

[32] Elizbar A. Nadaraya. On Estimating Regression. *Theory of Probability & Its Applications*, 9(1):141–142, 1964. 4

[33] Anh Nguyen, Alexey Dosovitskiy, Jason Yosinski, Thomas Brox, and Jeff Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 3395——3403. Curran Associates Inc., 2016. 1

[34] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2(11):e7, 2017. 1

[35] Emanuel Parzen. On Estimation of a Probability Density Function and Mode. *The Annals of Mathematical Statistics*, 33(3):1065–1076, 1962. 3

[36] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. Association for Computing Machinery, 2016. 1, 2

[37] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 2, 4, 5, 6

[38] Wojciech Samek, Thomas Wiegand, and Klaus-Robert Müller. Explainable artificial intelligence: understanding, visualizing and interpreting deep learning models. *ITU JOURNAL: ICT DISCOVERIES*, 1(S1):39–48, 2017. 1

[39] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Gradcam: Visual explanations from deep networks via gradient-based localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 618–626, 2017. 1, 2, 5, 6

[40] Sungyong Seo, Jing Huang, Hao Yang, and Yan Liu. Interpretable convolutional neural networks with dual local and global attention for review rating prediction. In *Proceedings of the eleventh ACM conference on recommender systems*, pages 297–305. Association for Computing Machinery, 2017. 1

[41] Hideaki Shimazaki and Shigeru Shinomoto. Kernel bandwidth optimization in spike rate estimation. *Journal of Computational Neuroscience*, 29(1):171–182, 2010. 3

[42] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3145–3153. PMLR, 2017. 1

[43] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *Third International Conference on Learning Representations, ICLR*, pages 1–14, 2015. 2, 4, 5, 6, 7

[44] Leon Sixt, Maximilian Granz, and Tim Landgraf. When explanations lie: Why many modified bp attributions fail. In *Proceedings of the 37th International Conference on Machine Learning*, pages 9046–9057. PMLR, 2020. 3

[45] Pascal Sturmfels, Scott Lundberg, and Su-In Lee. Visualizing the Impact of Feature Attribution Baselines. *Distill*, 5(1):e22, 2020. 3, 6

[46] Mukund Sundararajan and Amir Najmi. The Many Shapley Values for Model Explanation. In *Proceedings of the 37th International Conference on Machine Learning*, pages 9269–9278. PMLR, 2020. 7

[47] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3319–3328. PMLR, 2017. 3, 6, 7

[48] Ahmed Taha, Xitong Yang, Abhinav Shrivastava, and Larry Davis. A generic visualization approach for convolutional neural networks. In *Computer Vision – ECCV 2020*, pages 734–750. Springer International Publishing, 2020. 1

[49] Mingxing Tan, Ruoming Pang, and Quoc V. Le. Efficientdet: Scalable and efficient object detection. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10778–10787, 2020. 1

[50] Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and Xia Hu. Score-cam: Score-weighted visual explanations for convolutional neural networks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 111–119, 2020. 2, 6

[51] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. 6

[52] Christian Wolf and Markus Lappe. Salient objects dominate the central fixation bias when orienting toward images. *Journal of vision*, 21(8):23–23, 2021. 6

[53] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Computer Vision – ECCV 2014*, pages 818–833. Springer International Publishing, 2014. 1

[54] Jianming Zhang, Sarah Adel Bargal, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. Top-Down Neural Attention by Excitation Backprop. *International Journal of Computer Vision*, 126(10):1084–1102, 2018. 3

[55] Lin Zhang, Lei Zhang, Xuanqin Mou, and David Zhang. FSIM: A Feature Similarity Index for Image Quality Assessment. *IEEE Transactions on Image Processing*, 20(8):2378–2386, 2011. 6

[56] Quanshi Zhang, Yu Yang, Haotian Ma, and Ying Nian Wu. Interpreting cnns via decision trees. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6254–6263, 2019. 1

[57] Bolei Zhou, Aditya Khosla, Àgata Lapedriza, Aude Oliva, and Antonio Torralba. Object detectors emerge in deep scene cnns. In *Third International Conference on Learning Representations, ICLR*, pages 1–12, 2015. 1

[58] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning Deep Features for Discriminative Localization. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2921–2929, 2016. 2

[59] Yanzhao Zhou, Yi Zhu, Qixiang Ye, Qiang Qiu, and Jianbin Jiao. Weakly Supervised Instance Segmentation Using Class Peak Response. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3791–3800, 2018. 3

**Mohammad Mahdi Dehshibi** received his PhD in Computer Science in 2017 from IAU, Iran. He is currently a research postdoctoral fellow at Universitat Oberta de Catalunya, Spain. He was also a visiting researcher at Unconventional Computing Lab, UWE, Bristol, UK. He has contributed to more than 60 papers published in scientific journals and international conferences. His research interests include Affective Computing, Unconventional Computing, Cellular Automata and Deep Learning.

**Mona Ashtari-Majlan** received her Master's degree in Health Systems Engineering from Amirkabir University of Technology, Tehran, in 2021. She is a PhD candidate in computer science at Universitat Oberta de Catalunya, Spain. Her area of interest includes Biomedical Image Processing, Computer Vision, and Deep Learning.

**Gereziher Adhane** is currently a PhD student at Universitat Oberta de Catalunya, Spain. He obtained his MSc from Osmania University (India) in 2013/14. His research interests includes deep learning, computer vision and fairness in AI.

**David Masip** is a professor in the Department of Computer Science, Multimedia and Telecommunications at Universitat Oberta de Catalunya (UOC) since February 2007, and the director of the UOC Doctoral School since 2015. He leads the SUNAI (Scene Understanding and Artificial Intelligence) research group. He studied computer science at the Universitat Autonoma de Barcelona (UAB) and received his PhD in September 2005, receiving the UAB's best thesis award in computer science.