
Reasoning with trees: interpreting CNNs using hierarchies

Caroline Mazini Rodrigues*
Laboratoire de Recherche de l'EPITA – LRE
Univ Gustave Eiffel, LIGM
caroline.mazinirodrigues@esiee.fr

Nicolas Boutry
Laboratoire de Recherche de l'EPITA – LRE
Le Kremlin-Bicêtre, 94270, France
nicolas.boutry@lrde.epita.fr

Laurent Najman
Univ Gustave Eiffel, CNRS, LIGM
Marne-la-Vallée, 77454, France
laurent.najman@esiee.fr

Abstract

Challenges persist in providing interpretable explanations for neural network reasoning in explainable AI (xAI). Existing methods like Integrated Gradients produce noisy maps, and LIME, while intuitive, may deviate from the model's reasoning. We introduce a framework that uses hierarchical segmentation techniques for faithful and interpretable explanations of Convolutional Neural Networks (CNNs). Our method constructs model-based hierarchical segmentations that maintain the model's reasoning fidelity and allows both human-centric and model-centric segmentation. This approach offers multiscale explanations, aiding bias identification and enhancing understanding of neural network decision-making. Experiments show that our framework, *xAiTrees*, delivers highly interpretable and faithful model explanations, not only surpassing traditional xAI methods but shedding new light on a novel approach to enhancing xAI interpretability. Code at: https://github.com/CarolMazini/reasoning_with_trees.

1 Introduction

In modern deep learning applications, especially in healthcare and finance, there is a growing need for transparency and explanation. Understanding a model's rationale is crucial before relying on its predictions. This need arises from biases present at various stages of model development and deployment. While some biases help in learning data distribution [12], others may indicate data imbalance, incorrect correlations, or prejudices in data collection.

To meet the demand for explanation, Explainable Artificial Intelligence (xAI) provides methods that clarify models' decision-making processes. In healthcare, tools like GradCAM [26], which shows heatmaps of important image regions, and LRP [2], which attributes importance to features (pixels), help in understanding deep learning models across various applications [4, 9, 5], including ultrasound [3] and X-ray [1] imaging. These techniques were crucial during the recent Covid-19 outbreaks, aiding in the diagnosis process [19, 13]. However, these methods are approximations of model behavior. Different techniques prioritize either faithfulness to the model's behavior or human interpretability, posing a challenge in balancing the two.

Object-structure-based visualizations enhance human interpretation by decomposing images in ways that mimic human perception, grouping objects by attributes like color, texture, and edges [15].

*<https://carolmazini.github.io/>

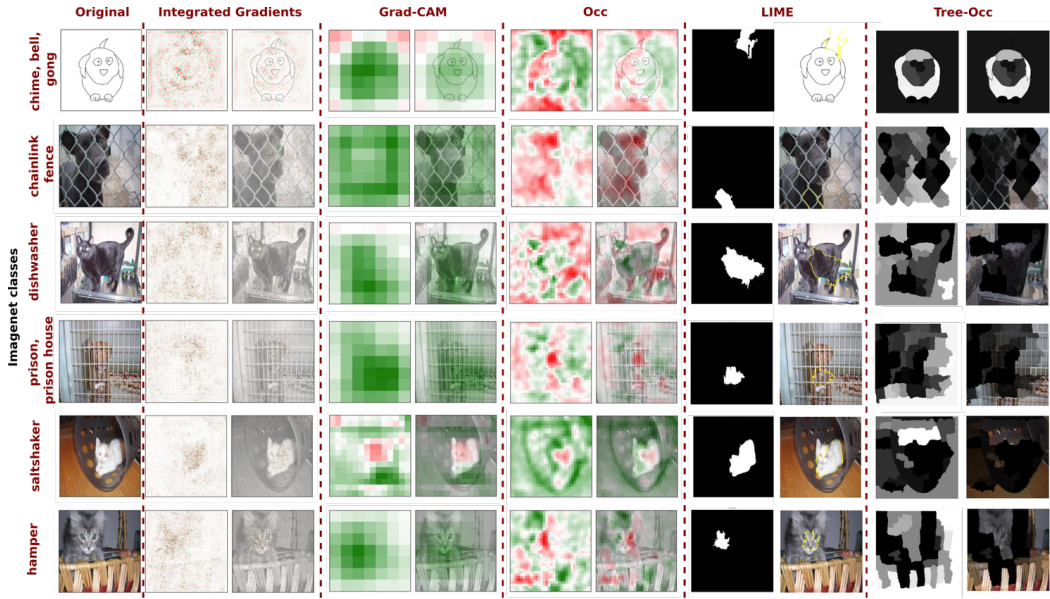


Figure 1: Explanations of six image classes predicted by VGG-16 and Resnet18 models trained on the Imagenet dataset. We compare four well-known xAI method explanations with one configuration of *xAiTrees*: Tree-Occ. Methods such as Integrated Gradients are quite noisy and difficult to interpret. Shapes such as the grades and the fence seem to be better highlighted by Tree-Occ, which is helpful for interpretation. When compared to highly interpretable methods like LIME, Tree-Occ avoids the mistake of highlighting the cat when the models predict classes such as dishwasher, saltshaker, and hamper.

Techniques such as LIME [24] and KernelSHAP [20] have used this approach effectively, segmenting images into meaningful parts to improve interpretability. However, the size of segmented regions affects the information extracted: small regions can be hard to interpret, while large regions may miss fine details. Additionally, using a segmentation framework introduces human bias, which aids comprehension but may reduce fidelity to the model’s actual behavior.

In this paper, we explore the trade-off between explaining Convolutional Neural Networks (CNNs) with model faithfulness and human interpretability. We introduce an innovative framework that combines hierarchical segmentation with region-based explanation methods, creating a human-friendly multiscale visualization, inspired by the Multiscale Interpretable Visualization (Ms-IV) technique [25]. Unlike traditional region-based xAI techniques that segment images into a fixed number of levels, *xAiTrees* leverages hierarchical segmentation to maintain varying degrees of abstraction within object structures. Additionally, to enhance our understanding of human-based segmentation and its relationship with model knowledge, we propose a *model-based segmentation* using pixel-wise xAI methods to reveal the model’s “vision”.

Through the experiments with these two approaches — human-based and model-based hierarchical segmentation explanations — we assess several aspects of explainability: the fidelity to the model’s behavior, the effectiveness in detecting bias within the models, and the ease of interpretability. The key contributions of this paper include:

1. A hierarchical segmentation explanation framework aimed at integrating the importance of multiscale regions in the model’s predictions, *xAiTrees*;
2. An integrated model-based segmentation approach within the framework *xAiTrees*, offering more faithful explanations to the model;
3. A quantitative comparison with established xAI techniques and a qualitative assessment against human-analysis for bias identification.

In this work, we demonstrate through extensive experimentation and analysis, both quantitative and qualitative, that our proposed framework significantly enhances explainability and interpretability. By conducting a comprehensive evaluation and comparison with state-of-the-art xAI visualization

methods, we provide robust evidence that our framework offers superior performance. The results highlight the efficacy of our approach in making complex models more transparent and understandable, addressing key challenges in the field of explainable artificial intelligence. We organize the paper as follows: in Section 2, we present some prior research on xAI. Section 3 outlines the preliminary concepts used in our framework, while in Section 4 we provide a detailed explanation of our methodology. In Section 5 we present and discuss our experimental results. Finally, we conclude and discuss possible future research directions in Section 7.

2 Related work

Classification problems and xAI: One fundamental task in machine learning is classification. The basic concept involves working with a training dataset, denoted as $\mathcal{DS} = (\mathcal{I}_i, GT_i)_{i \in [1, NbIm]}$, which consists of pairs of images \mathcal{I}_i and their associated labels GT_i . Each label belongs to one of a set of classes represented by $c \in [1, NbClasses]$. The goal is to train a model, denoted as Ξ , to effectively distinguish between different classes within the dataset.

In this configuration, we express Ξ as $\Xi = \Xi^{classif} \circ \Xi^{enc}$, the combination of two elements: an Ξ^{enc} , responsible for converting each input image \mathcal{I}_i into a feature vector, and a $\Xi^{classif}$, which analyzes these features to classify the images. The outcome of this process, referred to as the “logit” for image \mathcal{I}_i , is a vector $\mathbf{out}_i \in \mathbb{R}^{NbClasses}$ that signifies the activation levels across various classes. Typically, we apply a *Softmax* layer to \mathbf{out}_i to determine the class with the highest activation, ideally aligning with the ground truth label GT_i for perfect classification.

Pixel-wise explanations: In neural networks, optimizing the model Ξ involves the backpropagation process. Exploiting this process, certain explainable Artificial Intelligence (xAI) methods like Integrated Gradients [30], Guided-Backpropagation [29], and Deconvolution [34] utilize it to identify input features that enhance the response of a specific class, aiming to maximize the value of a particular position in the output vector \mathbf{out}_i . Consequently, attribution maps are generated, illustrating pixel-level explanations, as depicted in Fig. 1 for Integrated Gradients (IG).

Region-based explanations: Additional techniques like Sensitivity Analysis [34], LIME [24], and SHAP [20] utilize occlusions of image regions to assess the network’s sensitivity to each region within an image. These methods provide explanations at a region level rather than a pixel level, as illustrated in Fig. 1 for LIME.

Concept-based explanations: However, many of these techniques focus on explaining individual samples separately, which limits our understanding of how the model behaves globally across various scenarios. That is why methods like TCAV [16], ACE [11], Explanatory graphs [35], LGNN [31], and Ms-IV [25] aim to comprehend the overall behavior of the model. In particular, Ms-IV also considers the impact of occlusions, not on individual predictions, but on the model’s output space.

3 Preliminaries

To ensure a thorough understanding of the sequel, we provide in this section the general techniques and metrics employed during this work. In subsection **A**, we provide a brief overview of the selected hierarchical segmentation techniques, highlighting their significance. In subsection **B**, we shortly present the occlusion-based metrics used in the construction of our methodology.

A. Segmentation techniques: As an important step for our framework, we employ image segmentation algorithms that decompose images into more interpretable structures, enabling better human understanding and interpretation. We specifically employ hierarchical segmentation techniques due to their capability to decompose images into multiple levels of detail, from fine to coarse, mirroring how humans naturally perceive objects: initially observing the overall structure before delving into the finer details. A hierarchical segmentation algorithm produces a merging tree, that indicates how two given regions merge. In this paper, we use the tree structures available in the Hgra package [22, 21] : Binary Partition Tree (BPT) and Hierarchical watershed. See details in Appendix A.1.

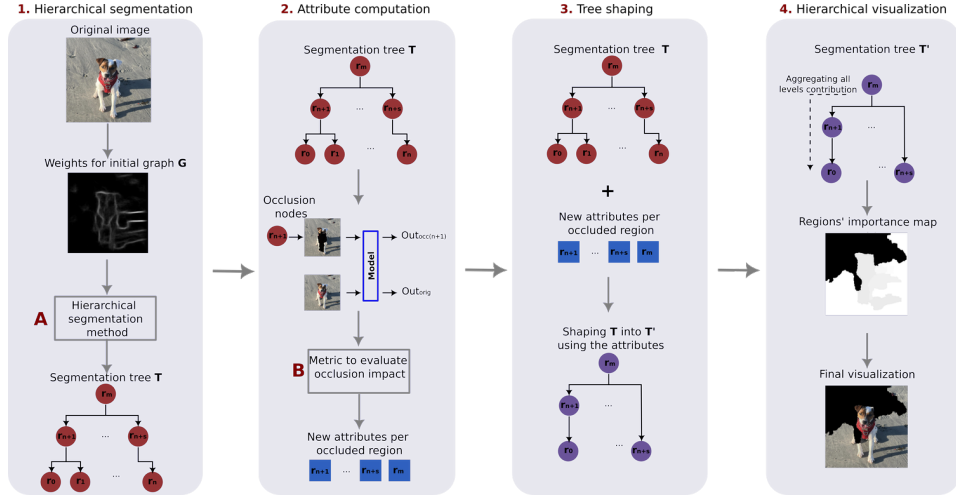


Figure 2: Our framework *xAiTrees* operates through four key steps: **1.** Generate a segmentation hierarchy using either the image’s edge map for human-based segmentation or pixel-wise importance based on xAI techniques for model-based segmentation. **2.** Systematically occlude each region of the segmentation to evaluate its impact on the model’s decision, obtaining an occlusion attribute for each region. **3.** Assess the persistence of the occlusion attribute using a shaping approach [33, 32]. **4.** Aggregate the contributions of each region from the highest to the lowest level of the tree to create a comprehensive multiscale visualization.

B. Occlusion-based metrics: In this work, we use two metrics to generate our segmentation based on the model explainability: (i) *Occlusion*, which is the impact of occluding an image region on its classification output, and (ii) CaOC which is the intra-class impact of occluding an image region. For (i), we assess how the output of a model changes when an image region is occluded. For (ii), we employ a sliding metric that ranks images based on the highest activations for a given class. We then measure the movement in this ranking after occluding a region of the image, determining the intra-class impact of the occlusion (detailed in Appendix A.2).

4 Methodology

In this section, we outline our four-step methodology (Figure 2): (1) hierarchical segmentation, (2) attribute computation, (3) tree shaping, and (4) hierarchical visualization. In step 1, we convert the data into a hierarchical representation, creating various regions at different scales in the image. In step 2, we evaluate some xAI-based attributes (B) on the regions. In step 3, we assess the importance of the region attributes. Finally, in step 4, we explain how to generate a visualization map from the importance of the attributes.

1. Hierarchical segmentation: Intuitively, any hierarchical segmentation algorithm works by iteratively merging first the pixels, then the regions, according to a similarity criterion. In this paper, we test two ways for measuring the similarity: human-based and model-based.

- The human-based approach relies on the Structured Edge Detection (SED) algorithm [10], which captures complex edge patterns and produces precise edge maps, in accordance with human intuition.
- The model-based approach uses a visual representation of the image’s pixels most influential in a model’s decision. Although less intuitive for humans, this approach helps to understand how the model reasons. We test pixel-wise explainable AI methods: Integrated Gradients (IG) [30], Guided-Backpropagation [29], Input x Gradient [27], and Saliency [28] (all from Captum framework). The methods were chosen for their state-of-the-art, pixel-wise importance attribution.

Using such a similarity criterion, we obtain a hierarchical segmentation, which can be represented as a tree T , completing the first step of our pipeline. See Figure 2, first column.

2. Attribute computation: The segmentation tree generated in the previous step provides many segments. We assess the model’s response on each segmented region in the tree, for all regions large enough. We apply a metric to evaluate the occlusion impact caused by each region. These occlusion scores reveal the influence of each segmented regions on the model’s output. The metric employed to assess the impact of regions can be any occlusion-based metric. See Fig. 2, second column.

3. Tree shaping: To assess the importance of the nodes’ attributes, it is not enough to simply take the regions with the highest attributes: there are too many of them. Instead, we rely on a process called *shaping* [33, 32]. The main idea is to look at the undirected, vertices-weighted graph G , whose vertices are the node of \mathbf{T} , whose edges are formed by the parent-children relationship in \mathbf{T} , and whose weights are the attributes of the nodes. We now look at the level-sets of G . A vertex of G (a node of \mathbf{T}) is important according to its persistence in the level sets of G . More precisely, a connected component is born when a local maximum of the attribute appear; when two connected components merge, one of the two maxima disappear, and the time of life of this maximum is its persistence. We can compute such persistence by building a new tree \mathbf{T}' on G , \mathbf{T}' is the tree of all the connected component of the upper-level sets of G . The persistence of a node of \mathbf{T} is easily computed on \mathbf{T}' by computing the length of the branch it belongs to. We refer to [33, 32] for more details. See Fig. 2, third column.

4. Construction of the hierarchical visualization: With \mathbf{T}' from the previous step, we now produce a visualization of the important regions. Using the persistence of a node directly for visualization can yield conflicting results for interpretation. Consider an example where we want to generate explanations for a model that classifies images of dogs. The persistence might indicate that eyes are the primary features for correct classification. If the image under scrutiny shows a dog with its owner, the persistence might erroneously highlight the eyes of both the human and the dog as relevant, which is misleading since only the dog’s eyes should matter (in an ideal, unbiased model). To avoid such effect, we recursively sum the persistence of each node from the root to the leaves of \mathbf{T}' . This ensures that smaller segments inherit the importance of their parent nodes. In our example, if the parent segment of the eyes is the entire face, the dog’s face carries importance for the model’s decision, while the human face does not. By adding the dog’s facial region information to the eye segments, we ensure the dog’s eyes are prioritized over the human eyes and, therefore, become more prominent in the explanation. This process aggregates the importance of various scales of the image into the pixels, resulting in a hierarchical, multi-scale, visualization. We use this aggregated persistence as the final score for each region of the hierarchical segmentation. We select a minimum importance score, and retain the regions accordingly. We superpose the retained region on the original image to generate the **Final visualization** (Fig. 2, fourth column).

5 Experiments and results

We evaluated the methods using two architectures, VGG-16 [28] and ResNet18 [14], trained on two datasets: Cat vs. Dog [7] (RGB images with a size of 224x224) and CIFAR-10 [17, 18] (RGB images with a size of 32x32). Explanations were generated for 512 images from the Cat vs. Dog dataset and 10,000 images from the CIFAR-10 dataset. A detailed description of the methods’ parameters and datasets is provided in the Appendix A.4 and A.5. We organize our experiments and results into two categories: quantitative and qualitative analysis. In the quantitative analysis, we conduct a series of experiments utilizing the metrics discussed in Section 3 to assess the impact of image region occlusion of various explainable frameworks. During our qualitative analysis, we delve into a more subjective examination, evaluating the human interpretability of the explanations generated by the models. The experiments were conducted on GPU (NVIDIA Quadro RTX 8000 48GB).

5.1 Quantitative evaluations

We selected state-of-the-art region-based methods as baseline (**B**) to be compared: Occlusion, Grad-CAM, LIME, and Ms-IV. Although ACE presents good concept-based explanations, we only use it in the human evaluation experiments because, as a global explanation method, it is not directly comparable to the local ones in these quantitative experiments (more details in Appendix A.6). We compare the baseline methods with two configurations of our proposed methodology: **C1** and **C2**. This is done to explore variations in the configuration of our methodology, including variations in the

Table 1: Percentage of images with the original class changed after the **exclusion** of selected explanation regions (a), and **inclusion** (exclusively) of this same regions. We test two configurations of our methodology (**C1** and **C2** – other configurations in Supplementary Materials) against four region-based baseline methods, Occlusion, Grad-CAM, LIME and Ms-IV, in two architectures, VGG-16 and ResNet18, and datasets, Cat vs. Dog and CIFAR10. We expect higher percentage of class change (**Ch.**) when the region is excluded (a) and lower when the region is included (b). **Same** column shows images maintaining the original class when the output was reduced, and **Total** is the sum of class change (**Ch.**) and class reduction (**Same**).

| % of images | | Cat vs. Dog | | | | | | Cifar10 | | | | | |
|-------------|---------------|-------------|------|-------|--------|------|-------|---------|------|-------|--------|------|-------|
| | | VGG | | | ResNet | | | VGG | | | ResNet | | |
| | | Ch. | Same | Total | Ch. | Same | Total | Ch. | Same | Total | Ch. | Same | Total |
| B | Occlusion | 0.05 | 0.93 | 0.98 | 0.06 | 0.89 | 0.95 | 0.13 | 0.50 | 0.63 | 0.17 | 0.44 | 0.61 |
| | Grad-CAM | 0.07 | 0.82 | 0.89 | 0.13 | 0.83 | 0.96 | 0.08 | 0.45 | 0.53 | 0.47 | 0.14 | 0.61 |
| | LIME | 0.07 | 0.83 | 0.90 | 0.07 | 0.76 | 0.83 | 0.31 | 0.38 | 0.69 | 0.29 | 0.37 | 0.66 |
| | Ms-IV | 0.06 | 0.76 | 0.82 | 0.07 | 0.66 | 0.73 | 0.19 | 0.44 | 0.63 | 0.21 | 0.43 | 0.64 |
| C1 | Tree-CaOC | 0.16 | 0.48 | 0.64 | 0.22 | 0.41 | 0.63 | 0.08 | 0.23 | 0.31 | 0.11 | 0.22 | 0.33 |
| | Tree-Occ | 0.31 | 0.63 | 0.94 | 0.35 | 0.60 | 0.95 | 0.32 | 0.26 | 0.58 | 0.28 | 0.22 | 0.50 |
| | IG-Tree-CaOC | 0.29 | 0.46 | 0.75 | 0.21 | 0.45 | 0.66 | 0.13 | 0.36 | 0.49 | 0.20 | 0.42 | 0.62 |
| | IG-Tree-Occ | 0.43 | 0.54 | 0.97 | 0.32 | 0.61 | 0.93 | 0.39 | 0.32 | 0.71 | 0.35 | 0.30 | 0.65 |
| C2 | TreeB-CaOC | 0.35 | 0.39 | 0.74 | 0.27 | 0.42 | 0.69 | 0.11 | 0.34 | 0.55 | 0.15 | 0.32 | 0.47 |
| | TreeB-Occ | 0.51 | 0.44 | 0.95 | 0.41 | 0.52 | 0.93 | 0.44 | 0.28 | 0.72 | 0.39 | 0.25 | 0.64 |
| | BP-TreeB-CaOC | 0.56 | 0.32 | 0.88 | 0.39 | 0.39 | 0.78 | 0.08 | 0.44 | 0.52 | 0.14 | 0.42 | 0.56 |
| | BP-TreeB-Occ | 0.63 | 0.35 | 0.98 | 0.55 | 0.37 | 0.92 | 0.49 | 0.27 | 0.76 | 0.41 | 0.25 | 0.66 |

(a)

| % of images | | Cat vs. Dog | | Cifar10 | |
|-------------|---------------|-------------|--------|---------|--------|
| | | VGG | ResNet | VGG | ResNet |
| B | Occlusion | 0.47 | 0.50 | 0.48 | 0.47 |
| | Grad-CAM | 0.51 | 0.30 | 0.47 | 0.0 |
| | LIME | 0.16 | 0.30 | 0.25 | 0.30 |
| | Ms-IV | 0.20 | 0.54 | 0.41 | 0.43 |
| C1 | Tree-CaOC | 0.26 | 0.32 | 0.43 | 0.43 |
| | Tree-Occ | 0.17 | 0.23 | 0.32 | 0.35 |
| | IG-Tree-CaOC | 0.41 | 0.43 | 0.43 | 0.43 |
| | IG-Tree-Occ | 0.19 | 0.42 | 0.36 | 0.38 |
| C2 | TreeB-CaOC | 0.16 | 0.22 | 0.42 | 0.41 |
| | TreeB-Occ | 0.11 | 0.19 | 0.26 | 0.29 |
| | BP-TreeB-CaOC | 0.38 | 0.28 | 0.46 | 0.44 |
| | BP-TreeB-Occ | 0.04 | 0.18 | 0.24 | 0.32 |

(b)

size of minimal regions in the visualizations, pixel weights for graph construction in segmentation, and methods for generating hierarchical segmentation. In **C1**, we present results for minimal regions of 500 pixels (Cat vs. Dog) and 64 pixels (CIFAR-10), utilizing edges and Integrated Gradients (IG) as pixel weights, and the watershed-by-area hierarchical segmentation. This configuration was selected for its stability across different minimal region sizes in the datasets, and its visualizations were used for human evaluation. While **C2** presents results for minimal regions of 200 pixels (Cat vs. Dog) and 4 pixels (CIFAR-10), using edges and Guided Backpropagation (BP) as pixel weights, and the BPT tree. This configuration yielded the highest performance. Comprehensive results for other configurations are provided in the Appendix A.5. Here, we propose two main quantitative evaluations: (i) Exclusion of important regions; and (ii) Inclusion of important regions. These experimental configurations are further discussed as follows:

Exclusion of important regions: Given that each region-based explainable AI (xAI) method identifies important regions that explain the prediction of a model, we performed occlusion of these regions, in order to measure the impact of each selection. For methods that assign scores to regions, we masked the 25% highest scores (this excludes LIME, which inherently provides information to directly mask each region, without the need for additional threshold of image region importance.).

In Table 1(a), we present results (**Ch.**, **Same**, **Total**) for each explainable technique (**B** – Baseline, **C1**, and **C2** – our proposition) applied to a network (VGG or ResNet) classifying images from a dataset (Cat vs. Dogs or Cifar10). **Ch.** is the percentage of images that changed class when the important region is concealed. **Same** is the percentage of images that remained in the same class after occlusion but with reduced classification certainty. **Total** is the percentage of all images with the class negatively impacted by the removal of important regions (sum of **Ch.** and **Same**). Higher **Ch.** values indicate that the identified regions are more class-representative. High **Same** values complement **Ch.**, suggesting that the best results are shown by higher **Ch.** and **Same** values. Thus, while **Total** sums **Ch.** and **Same**, the optimal result is reflected by initially higher **Ch.** and then higher **Same** values.

We can observe from the experiments in Table 1(a) that baseline methods such as **Occlusion** achieved **Total** values above 80%. However, the best results, based on **Ch.** being the most critical factor, were achieved using our methodology, specifically the **C2** configurations. Our techniques had the highest percentages of class changes in Cat vs. Dog images, with over 60% indicated by **Ch.**. For smaller images (CIFAR-10), the class change exceeded 40%. Among the baseline methods, LIME had the best results, with 7% class change for high-dimensional images and 31% for smaller images. This experiment demonstrates the superiority of our **C1** and **C2** configurations over state-of-the-art baselines in identifying the most impactful regions within an image. By achieving notably higher percentages of class changes (**Ch.**) followed by classification certainty (**Same**), in scenarios such as Cat vs. Dog images and CIFAR-10 datasets, our methodologies exhibit robustness and effectiveness

across various image classification tasks. These findings underscore the significance of our approach in providing more accurate insights into the interpretability of deep neural networks.

We present the results of a second experiment in Table 2. To address the issue of unhelpful explanations resulting from methods selecting the entire image as important, potentially leading to class changes upon occlusion, we introduce a novel metric, termed **Pixel Impact Rate (PIR)**. This metric quantifies the impact on class activation per occluded pixel. Complementing the percentage of class change, PIR distinguishes whether changes are primarily caused by complete or near-complete occlusion of the image. Higher PIR values indicate that each occluded pixel has a significant average impact, suggesting that concealing larger portions or the entire image leads to lower PIR, indicating less precision in the concealed area. Table 2 displays for each network, explainable technique, and dataset the average (**avg**) and standard deviation (**std**) of PIR.

Regarding the results of the Pixel Impact Rate (PIR) experiments displayed in Table 2, our configuration **C2**, particularly **BP-TreeB-CaOC**, demonstrated the best average PIR values for the Cats vs. Dog dataset overall. Compared to the baseline methods, **C2** showcased superior performance, with **Occlusion** demonstrating competitive results, followed by Grad-CAM. Considering the CIFAR-10 dataset, most methods, except for Grad-CAM on ResNet, presented similar magnitudes of PIR, indicating that the size of the regions was proportional to their impact. Grad-CAM on CIFAR-10 with ResNet occluded almost the entire image in most cases (smaller PIR values). Based on these results, we can highlight the distinct effectiveness of **C1** and **C2** in preserving region specificity and therefore increasing occluded pixel impact.

Inclusion of important regions: Additional experimentation was conducted to demonstrate a method’s capability to identify an image region with sufficient information for the original class. The goal of this experiment is to determine whether the selected important region, when the only one left unoccluded in the image, can maintain the classification in its expected class. This experiment elucidates the critical role of these identified regions, providing strong evidence that they indeed contain essential information for accurate classification. We occluded **all regions** in the images except for the one selected by each method. We then calculated the percentage of images that changed class. The results are presented in Table 1(b). Lower percentages indicate better performance, as they mean that a smaller percentage of images changed class, demonstrating that the chosen regions were sufficient to preserve the class for most of the images.

The metric presented in Table 1(b) highlights the capability of both LIME and our methodology **C1** and **C2** to identify regions that can sufficiently describe the class. However, our configuration, **BP-TreeB-Occ**, is still able to outperform LIME results, with, in some cases, less than half the number of images changing class. This shows that our configuration produces more essential information for class attribution. Some additional insights we obtained from these experiments include the following: **Occlusion** combined with our methodology appears to achieve superior results for local explanations (explaining individual images). Generally, using “model”-based segmentation leads to more faithful

Table 2: Pixel Impact Rate (PIR) of the chosen regions. The metric is the rate of the impact under occlusion (difference between the original class output and the output under occlusion) by the number of pixels of the occlusion mask. We test two configurations of our methodology (**C1** and **C2** – other configurations in Supplementary Materials) against four region-based baseline methods, Occlusion, Grad-CAM, LIME and Ms-IV, in two architectures, VGG-16 and ResNet18, and datasets, Cat vs. Dog and CIFAR10. We expect higher values, on average, for PIR, meaning each occluded pixel has a high impact.

| PIR | | Cat vs. Dog | | | | Cifar10 | | | |
|-----|---------------|-------------|----------|----------|----------|----------|----------|----------|----------|
| | | VGG | | ResNet | | VGG | | ResNet | |
| | | avg | std | avg | std | avg | std | avg | std |
| B | Occlusion | 4.60e-03 | 4.05e-03 | 1.50e-03 | 1.28e-03 | 1.09e-02 | 6.29e-02 | 9.47e-03 | 5.37e-02 |
| | Grad-CAM | 1.12e-03 | 1.02e-03 | 2.76e-04 | 2.07e-04 | 7.05e-04 | 4.64e-03 | 3.83e-04 | 1.19e-03 |
| | LIME | 9.03e-04 | 1.10e-03 | 3.47e-04 | 3.89e-04 | 1.38e-03 | 3.00e-02 | 1.19e-03 | 2.81e-02 |
| | Ms-IV | 4.30e-04 | 4.74e-04 | 1.83e-04 | 2.45e-04 | 1.34e-03 | 6.60e-03 | 1.25e-03 | 6.41e-03 |
| C1 | Tree-CaOC | 3.61e-04 | 4.70e-04 | 1.92e-04 | 2.86e-04 | 3.77e-02 | 3.31e-01 | 4.07e-02 | 3.38e-01 |
| | Tree-Occ | 3.66e-04 | 5.30e-04 | 1.69e-04 | 2.52e-04 | 6.05e-02 | 5.05e-01 | 5.83e-02 | 4.87e-01 |
| | IG-Tree-CaOC | 3.04e-04 | 3.48e-04 | 2.26e-04 | 3.09e-04 | 2.12e-02 | 2.41e-01 | 2.29e-02 | 2.60e-01 |
| | IG-Tree-Occ | 3.10e-04 | 3.61e-04 | 2.11e-04 | 3.05e-04 | 7.01e-03 | 1.55e-01 | 7.53e-03 | 1.58e-01 |
| C2 | TreeB-CaOC | 2.16e-04 | 2.91e-04 | 1.26e-04 | 2.26e-04 | 2.28e-02 | 2.39e-01 | 2.80e-02 | 2.82e-01 |
| | TreeB-Occ | 2.26e-04 | 3.32e-04 | 1.03e-04 | 1.81e-04 | 1.46e-02 | 2.42e-01 | 1.71e-02 | 2.66e-01 |
| | BP-TreeB-CaOC | 5.23e-03 | 3.59e-02 | 2.58e-03 | 1.81e-02 | 2.55e-02 | 1.68e-01 | 3.04e-02 | 2.07e-01 |
| | BP-TreeB-Occ | 8.64e-04 | 1.60e-02 | 1.18e-03 | 8.90e-03 | 1.61e-02 | 1.78e-01 | 2.47e-02 | 2.31e-01 |

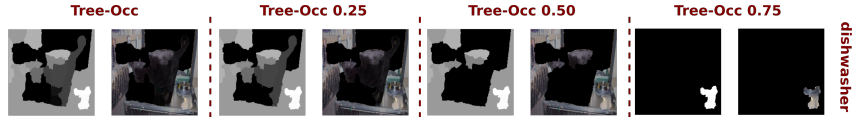


Figure 3: Different visualization levels on the explanation hierarchy. We illustrate a deeper analysis of the explanations of an image from Figure 1 using Tree-Occ (minimal region size of 500 pixels). We can note the evolution of the importance in the image’s shapes: the initial explanations show the sink as important but at the most selective level, the cat’s dish is the only one remaining. This analysis can be helpful to understand the reasoning behind predictions.

explanations. Our methodology outperformed traditional xAI methods used as baselines, including LIME. However, LIME showed consistently good results across all tests.

5.2 Qualitative analysis

As qualitative experiments, we want to visually evaluate the explanations for different interpretability tasks. In this section, we perform experiments to (i) identify reasons for misclassification of images, and (ii) evaluate explanations through the human interpretation of biased-trained networks. Our findings are reported in the following paragraphs:

Comparison of misclassified images: We searched for examples that were misclassified by models (VGG-16 and Resnet18) trained on Imagenet [8]. Figure 1 shows the explanations generated by Integrated Gradients, Grad-CAM, Occlusion, LIME, and Tree-Occ (500 pixels minimal region) of six images incorrectly classified. In Figure 1, the first column displays classes (such as chime, fence, dishwasher, among others) alongside examples of misclassified images. These images should have been classified as cat or dog. We then apply methods used in previous quantitative comparisons to generate visual explanations for why these images were misclassified. The figure illustrates that methods like Integrated Gradients, Grad-CAM, and Occlusion (Occ) may cause confusion in precisely identifying what caused the misclassification and may lead to poor human interpretation (we properly evaluate this in next experiment **Human evaluation in bias analysis**). Although LIME and our proposed Tree-Occ method can pinpoint interesting regions, the Tree-Occ method better illustrates the motivation behind misclassified results, as evident in the last column. For instance, in the fence example, it highlights the diamond pattern found on fences, while in the dishwasher example, it focuses solely on the sink region, disregarding the cat. Considering the hierarchical characteristic of our methodology, we can perform a deeper analysis of the explanations by selecting regions by the percentage of importance to be visualized, as shown in Figure 3 (more examples in Appendix A.6). In the last level of the dishwasher example, the model seems to focus on the cat’s dish after having focused on the sink (in the previous level).

Human evaluation in bias analysis: As previously mentioned, we used the configuration **C1** for human-interpretation evaluation. We trained three Resnet18 models subjected to data bias: (a) **Bias 1** – a model trained with dogs and only cats on cushions; (b) **Bias 2** – a model with cats and only dogs with grids; (c) **Bias 3** – and a model with dogs and only cats with humans (details of validation accuracy and visualizations in Supplementary Material). We presented the same five image visualizations (from corrected classified images by the biased class) for the baseline methods and the methods from **C1**. We intended to verify if: (i) humans can detect the wrong focus given based on a class prediction (**Detection**); and (ii) humans can recognize which was the cause of the bias (**Identification**).

To test (i) and (ii), for each **Bias** (a,b, or c) type we produce for each of the xAI methods an explanation image. By presenting five image explanations (the same images) for each of the xAI methodologies, we asked volunteers, based on the explanations provided, what they think the highlighted regions referred to (generated explanations and extra experiments in Appendix A.6).

Table 3 presents the results of evaluating 34 individuals from diverse continents (South America, Europe, and Asia), fields (Human, Biological, and Exact sciences), and levels of AI expertise (ranging from no knowledge to expert, with over half being non-experts). The experiment aims to identify effective methods for revealing trained-with biases. For each xAI method (IG, Grad-CAM, Occ, LIME, Ms-IV, ACE, Tree-Occ, Tree-CaOC, IG-Tree-Occ, IG-Tree-CaOC) used to explain biases (1,

Table 3: Human evaluation results for the tasks of bias (i) **Detection** and (ii) **Identification**. We proposed five image explanations (from the biased class) for each method and model trained with dataset bias: (1) dogs and only cats on cushions, (2) cats and only dogs with grids, and (3) dogs and only cats with humans. We present the percentage of volunteers that were able to: detect the bias **Detection (i)** by indicating the bias or the focus on the background; and identify the bias **Identification (ii)** by indicating the bias. We also present the percentage of people that did not understand the explanations (**Not identified**) and that found the explanation focusing on the **Animal**. We expect higher results for detection and identification, and lower for not identified and animal. Methods such as Ms-IV, ACE and Tree-CaOC (that are concept-aware methods) perform better. However, our method, Tree-CaOC, using human-based segmentation presented the best results for the three biased-datasets detection and identification.

| | | IG | Grad-CAM | Occ | LIME | Ms-IV | ACE | Tree-Occ | Tree-CaOC | IG-Tree-Occ | IG-Tree-CaOC |
|--------|---------------------|------|----------|------|------|-------|------|----------|-----------|-------------|--------------|
| Bias 1 | Detection (i) | 23.6 | 0.0 | 26.5 | 14.7 | 29.4 | 0.0 | 20.5 | 46.9 | 15.2 | 0.0 |
| | Identification (ii) | 11.8 | 0.0 | 5.9 | 5.9 | 14.7 | 0.0 | 2.9 | 18.8 | 0.0 | 0.0 |
| | Not identified | 35.3 | 2.9 | 41.2 | 38.2 | 26.5 | 30.3 | 20.6 | 25.0 | 36.4 | 12.1 |
| | Animal | 41.1 | 97.1 | 32.3 | 47.1 | 44.1 | 69.7 | 58.9 | 28.1 | 48.4 | 87.9 |
| Bias 2 | Detection (i) | 5.8 | 0.0 | 0.0 | 5.9 | 23.5 | 61.7 | 57.6 | 59.4 | 40.7 | 29.4 |
| | Identification (ii) | 2.9 | 0.0 | 0.0 | 5.9 | 17.6 | 44.1 | 39.4 | 46.9 | 31.3 | 26.5 |
| | Not identified | 11.8 | 0.0 | 14.7 | 35.3 | 23.5 | 35.3 | 27.3 | 31.3 | 37.5 | 38.2 |
| | Animal | 82.4 | 100.0 | 85.3 | 58.8 | 53.0 | 3.0 | 15.1 | 9.3 | 21.8 | 32.4 |
| Bias 3 | Detection (i) | 35.3 | 14.7 | 41.2 | 50.0 | 29.4 | 11.8 | 42.4 | 57.6 | 51.5 | 57.6 |
| | Identification (ii) | 20.6 | 11.8 | 32.4 | 35.3 | 17.6 | 5.9 | 12.1 | 36.4 | 18.2 | 39.4 |
| | Not identified | 38.2 | 14.7 | 29.4 | 47.1 | 47.1 | 61.8 | 48.5 | 36.4 | 45.5 | 39.4 |
| | Animal | 26.5 | 70.6 | 29.4 | 2.9 | 23.5 | 26.4 | 9.1 | 6.0 | 3.0 | 3.0 |

2, and 3), we show the percentage of participants who detected, identified, or did not identify the bias in the explanation. **Detection** indicates perceiving the xAI explanation as either background or reflecting the bias, while **Identification** denotes successful interpretation of the explanation as the induced bias. **Not Identification** refers to being unable to interpret the explanation. Higher percentages in the Identification row are desirable. If not, we prioritize high values in the Detection row. Lower values in the Not Identification or **Animal** rows indicate clearer human interpretation of our trained-with bias.

We can observe that the results of Table 3 demonstrate that IG and Grad-CAM explanations had some difficulties during interpretation. Their results obtained a lot of **Not identified** and/or **Animal**, meaning that the highlighted explanations were not clear to be our imposed biases. We remarked that the best results of detection and identification were found by methods that were linked to contextual information (or global explanations) such as Ms-IV, ACE, Tree-CaOC, and IG-Tree-CaOC. This occurs due to the nature of the method, which reflects, more globally, the model’s knowledge. However, this seems not always to be enough for humans to provide a complete interpretation of the model’s knowledge. Once again Tree-CaOC, one of our configurations, presented the highest results for all three **Bias** for detecting and identifying, by combining global-aware metric (CaOC) and a human-based segmentation (edge detection). In these experiments, we demonstrate that our method excels compared to other studies in a crucial aspect of explainable AI: human interpretability.

6 Limitations

The time of computation of hierarchies varies, it depends on factors such as the image size, and the smallest region’s size of the segmentation. The selection of the importance threshold for the final visualization of regions depends on experimentation. We found a relatively small number of volunteers for our qualitative experiments.

7 Conclusion

In this paper, we present a framework, *xAiTrees*, aimed at integrating multiscale region importance in model predictions, providing more faithful and interpretable explanations. Our approach outperforms traditional xAI methods like LIME, especially in identifying impactful regions, in datasets such as Cat vs. Dog and CIFAR-10. Quantitative evaluations highlight the superiority of configurations like BP-TreeB-CaOC, achieving the best average PIR values for the Cats vs. Dogs dataset and consistently higher percentages of class changes and classification certainty.

Our methods, such as BP-TreeB-Occ, offer crucial class attribution information with fewer changes compared to LIME, providing superior local explanations when combined with Occlusion. Qualitative analysis demonstrates that our Tree-Occ method better elucidates misclassification motivations and provides clearer, hierarchical interpretations of model predictions. Techniques like Tree-CaOC, merging global-aware metrics with human-based segmentation, excel in detection and identification tasks, achieving superior results in human interpretability. In summary, our framework delivers highly interpretable and faithful model explanations, significantly aiding in bias detection and identification, and demonstrating its effectiveness in the field of explainable AI. Therefore, potentially aiding to reduce the societal negative impact that could be generated by deep learning models in high-risk decision-making process.

References

- [1] Sudil Hasitha Piyath Abeyagunasekera, Yuvim Perera, Kenneth Chamara, Udari Kaushalya, Prasanna Sumathipala, and Oshada Senaweera. Lisa : Enhance the explainability of medical images unifying current XAI techniques. In *International conference for Convergence in Technology (I2CT)*, pages 1–9, 2022. 1
- [2] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Muller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS One*, 10(7):1–46, 2015. 1
- [3] Jannis Born, Nina Wiedemann, Manuel Cossio, Charlotte Buhre, Gabriel Brändle, Konstantin Leidermann, Julie Goulet, Avinash Aujayeb, Michael Moor, Bastian Rieck, and Karsten Borgwardt. Accelerating detection of lung pathologies with explainable ultrasound image analysis. *Applied Sciences*, 11(2), 2021. 1
- [4] Katarzyna Borys, Yasmin Alyssa Schmitt, Meike Nauta, Christin Seifert, Nicole Krämer, Christoph M Friedrich, and Felix Nensa. Explainable AI in medical imaging: An overview for clinical practitioners - beyond saliency-based XAI approaches. *European Journal of Radiology*, 2023. 1
- [5] Ahmad Chaddad, Jihao Peng, Jian Xu, and Ahmed Bouridane. Survey of explainable ai techniques in healthcare. *Sensors*, 23(2), 2023. 1
- [6] Jean Cousty, Gilles Bertrand, Laurent Najman, and Michel Couprie. Watershed cuts: Minimum spanning forests and the drop of water principle. *IEEE transactions on pattern analysis and machine intelligence*, 31(8):1362–1374, 2008. 13
- [7] Will Cukierski. Dogs vs. cats redux: Kernels edition. <https://www.kaggle.com/competitions/dogs-vs-cats-redux-kernels-edition/data>, 2016. Accessed: 2024-05-20. 5
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255. Ieee, 2009. 8
- [9] S. Priya Dharshini, K. Ram Kumar, S. Venkatesh, K. Narasimhan, and K. Adalarasu. An overview of interpretability techniques for explainable artificial intelligence (XAI) in deep learning-based medical image analysis. In *International Conference on Advanced Computing and Communication Systems (ICACCS)*, volume 1, pages 175–182, 2023. 1
- [10] Piotr Dollár and C Lawrence Zitnick. Fast edge detection using structured forests. *IEEE transactions on pattern analysis and machine intelligence*, 37(8):1558–1570, 2014. 4
- [11] Amirata Ghorbani, James Wexler, James Zou Y, and Been Kim. Towards automatic concept-based explanations. In *Advances in Neural Information Processing Systems*, volume 32, pages 1–10, 2019. 3, 21
- [12] Anirudh Goyal and Yoshua Bengio. Inductive biases for deep learning of higher-level cognition. In *Proceedings of the Royal Society A*, volume 478, pages 1–49, 2022. 1
- [13] Arman Haghanifar, Mahdiyar Molahasani Majdabadi, Younhee Choi, S. Deivalakshmi, and Seokbum Ko. Covid-cxnet: Detecting covid-19 in frontal chest x-ray images using deep learning. *Multimedia Tools and Applications*, 81:30615–30645, 2022. 1

- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *29th IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 5
- [15] David H. Hubel and Torsten N. Wiesel. Receptive fields of single neurons in the cat’s striate cortex. *Journal of Physiology*, 148:574–591, 1959. 1
- [16] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie J. Cai, James Wexler, Fernanda B. Viégas, and Rory Sayres. Interpretability beyond feature attribution: Quantitative testing with Concept Activation Vectors (TCAV). In *35th International Conference on Machine Learning (ICML)*, pages 2668–2677, 2018. 3
- [17] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. 5
- [18] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). Technical report, Canadian Institute for Advanced Research, 2009. 5
- [19] Siyuan Lu, Ziquan Zhu, Juan Manuel Gorriz, Shui-Hua Wang, and Yu-Dong Zhang. Nagnn: Classification of covid-19 based on neighboring aware representation from deep graph neural network. *International Journal of Intelligent Systems*, 37:1572–159, 2022. 1
- [20] Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *31st International Conference on Neural Information Processing Systems (NeurIPS)*, pages 4768–4777, 2017. 2, 3
- [21] Benjamin Perret, Giovanni Chierchia, Jean Cousty, Silvio Jamil Ferzoli Guimaraes, Yukiko Kenmochi, and Laurent Najman. Higma (hierarchical graph analysis) documentation. <https://higma.readthedocs.io/>, 2018. Accessed = 2024-03-07. 3
- [22] Benjamin Perret, Giovanni Chierchia, Jean Cousty, Silvio Jamil Ferzoli Guimaraes, Yukiko Kenmochi, and Laurent Najman. Higma: Hierarchical graph analysis. *SoftwareX*, 10:100335, 2019. 3
- [23] Huy Phan. Pytorch models trained on cifar-10 dataset. https://github.com/huyvnphan/PyTorch_CIFAR10, 2021. Accessed: 2024-05-20. 14
- [24] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?": Explaining the predictions of any classifier. In *22nd International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 1135–1144, 2016. 2, 3
- [25] Caroline Mazini Rodrigues, Nicolas Boutry, and Laurent Najman. Unsupervised discovery of interpretable visual concepts. *Information Sciences*, 661:1–26, 2024. 2, 3, 13
- [26] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *16th International Conference on Computer Vision (ICCV)*, pages 618–626, 2017. 1
- [27] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *International Conference on Machine Learning (ICML)*, volume 70, pages 3145–3153, 2017. 4
- [28] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *3rd International Conference on Learning Representations (ICLR)*, pages 1–14, 2015. 4, 5
- [29] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. In *3rd International Conference on Learning Representations (ICLR)*, pages 1–14, 2015. 3, 4
- [30] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *34th International Conference on Machine Learning (ICML)*, pages 3319–3328, 2017. 3, 4
- [31] Randy Tan, Lei Gao, Naimul Khan, and Ling Guan. Interpretable artificial intelligence through locality guided neural networks. *Neural Networks*, 155:58–73, 2022. 3
- [32] Yongchao Xu, Edwin Carlinet, Thierry Géraud, and Laurent Najman. Hierarchical segmentation using tree-based shape space. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39:1–14, 04 2016. 4, 5

- [33] Yongchao Xu, Thierry Géraud, and Laurent Najman. Connected filtering on tree-based shape-spaces. *IEEE transactions on pattern analysis and machine intelligence*, 38(6):1126–1140, 2015. [4](#), [5](#)
- [34] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *13th European Conference on Computer Vision (ECCV)*, pages 818–833, 2014. [3](#)
- [35] Quanshi Zhang, Ruiming Cao, Feng Shi, Ying Nian Wu, and Song-Chun Zhu. Interpreting CNN knowledge via an explanatory graph. In *Conference on Artificial Intelligence (AAAI)*, pages 4454–4463, 2018. [3](#)

A Appendix

A.1 A. Segmentation techniques

As an important step for our framework, we employ segmentation techniques so we can decompose images, based on specific attributes, into more interpretable structures, enabling better human understanding and interpretation. We specifically employ hierarchical segmentation techniques due to their capability to decompose images into multiple levels of detail, mirroring how humans naturally perceive objects: initially observing the overall structure before delving into the finer details.

Trees: A tree is an acyclic graph, consisting of nodes that connect to zero or more other nodes. It starts with a “root” node that branches out to other nodes, ending in “leaves” with no children. In image representation, the root node represents the entire image, and each leaf represents a pixel, resulting in as many leaves as pixels. The structure between the root and leaves groups pixels into clusters at each level based on similarity metrics, with each level abstracting the one below. Using a segmentation tree, we can make *cuts* at various levels to obtain different numbers and sizes of segmented regions.

Binary Partition Tree (BPT): A Binary Partition Tree (BPT) is a data structure in which each node represents a region of the image. Similarly, the tree starts with a root node representing the entire image and branches out through a series of binary splits until reaching the leaf nodes, representing the individual pixels. Different from the tree, in which a node could have multiple splits, in the BPT each split, divides a region into two smaller sub-regions based on a criterion.

Watershed: This algorithm [6] constructs a hierarchical segmentation tree based on a minimum-spanning forest rooted in the local minima of an edge-weighted graph. In this context, local minima are points in the graph where the surrounding edge weights are higher, representing the lowest values in their neighborhood. These minima serve as starting points for the segmentation. The algorithm iteratively merges regions beginning from these local minima, guided by the edge weights that indicate dissimilarity between adjacent pixels. By progressively combining these regions, the algorithm builds the segmentation tree, effectively capturing the hierarchical structure of the image.

A.2 B. Occlusion-based metrics:

Here, we discuss the metrics used to generate our segmentation based on the model explainability (block **B** in Figure 2). We present two metrics: (i) *Occlusion*, which is the impact of occluding an image region on its classification output, and (ii) CaOC which is the intra-class impact of occluding an image region. For (i), we assess how the output of a model changes when an image region is occluded. For (ii), we employ a sliding metric that ranks images based on the highest activations for a given class. We then measure the movement in this ranking after occluding a region of the image, determining the intra-class impact of the occlusion.

Occlusion: Let us say we have a model Ξ producing an output \mathbf{out}_i for an image \mathcal{I}_i . By concealing portions of this image, creating a new image $\mathcal{I}_i^\blacksquare$, we obtain a different model output $\mathbf{out}_i^\blacksquare$. The significance of the occluded area concerning a particular class c is assessed by comparing the outputs:

$$\left| \mathbf{out}_{i,c} - \mathbf{out}_{i,c}^\blacksquare \right|. \quad (1)$$

If there is a significant difference, it indicates that the model strongly relies on this region for class activation, meaning that these regions have a high *impact* on the model’s decision.

CAOC: In the Ms-IV method, introduced by Rodrigues *et al.* [25], CaOC employs rankings to assess how occlusions affect the model’s output space. A *ranking* is a sequence of objects ordered according to a specific criterion, from the object most aligned with it to the least aligned. Suppose the criterion is to maximize class c . In that case, the first index i in this sequence represents the object (in our case, the image \mathcal{I}_i) with the highest activation for class c in the output \mathbf{out}_i . If we define a function *argsort* to obtain the indices of an ordered sequence of objects, we can derive the sequence of image indices that maximize class c : $Seq_c = \mathit{argsort}(\mathbf{out}_{\cdot,c}, \mathit{decreasing})$, with $\mathbf{out}_{\cdot,c}$ the vector of outputs for class c of a set of input images $(\mathcal{I}_i)_{i \in [1, NbIm]}$.

CaOC computes an initial ranking Seq_c for a subset of images $\mathcal{DS}' \subset \mathcal{DS}$, and then a subsequent ranking Seq'_c after occluding one region of image $\mathcal{I}_i \in \mathcal{DS}'$. The significance of this occluded image region for the model is determined by the difference in the positions of this image in the rankings given by $|position(Seq_{i,c}, \mathcal{I}_i) - position(Seq'_{i,c}, \mathcal{I}_i)|$.

This metric aims to assess the impact of occluding image regions not only against the original output \mathbf{out}_i but also against the outputs of a range of images. Incorporating the model’s output space into the analysis ensures that explanations consider the broader context (global model’s behavior). Hence, we can characterize it as globally aware, even when explaining a single sample.

A.3 Tested framework’s configuration

We tested four different sizes of minimal region for filtering the initial segmentation. For Cat vs. Dog dataset: 200, 300, 400, and 500 pixels. For CIFAR10: 4, 16, 32 and 64 pixels.

For the model-based segmentation we tested four xAI techniques to generation the initial graph \mathbf{G} on Figure 2: Integrated Gradients (IG), Guided-Backpropagation (BP), Input X Gradient (I X G), and Saliency (S).

We tested three algorithms to construct the hierarchical segmentation: Binary Partition Tree (BPT), Watershed with Area, and Watershed with Volume.

We tested two different occlusion based metrics to obtain the impact of regions used to shape the hierarchical tree: CaOC and OCC.

When we refer to **Tree-CaOC** or **TreeW-CaOC**, we mean the human-based segmentation (edges’ map) using Watershed area and CaOC as occlusion metric. When we refer to **IG-Tree-Occ** or **IG-TreeW-Occ**, we mean the model-based segmentation (using Integrated Gradients (IG) attributions) using Watershed area and Occ (simple occlusion – Equation (1)) as occlusion metric. When we refer to **BP-TreeB-Occ**, we mean the model-based segmentation (using Guided Backpropagation (BP) attributions) using BPT and Occ as occlusion metric.

A.4 Parameters of the baseline methods

For Grad-CAM method, we used the last convolutional layer of each architecture to generate the visualizations. For Occlusion (from Captum framework) we used, for Cat vs. Dog the stride of $3 \times 7 \times 7$ and sliding window of $3 \times 14 \times 14$, for CIFAR10 the stride of $3 \times 2 \times 2$ and sliding window of $3 \times 4 \times 4$. For LIME, we used the standard configuration for Cat vs. Dog (Quickshift kernel size of 4) and, Quickshift kernel size of 2 for CIFAR10. All the other methods followed the standard configuration.

A.5 Quantitative evaluations

Models’ description: Table 4 shows the number of images in train and validation sets for Cat vs. Dog and CIFAR10 datasets. We also include the train and validation accuracies for the models ResNet18 and VGG-16 used in the quantitative evaluations.

Cat vs. Dog models were trained with initial weights from Imagenet, learning rate $1e - 7$, cross-entropy loss, the Adam optimizer, and early stop in 20 epochs of non-improving validation loss.

CIFAR10 models were adapted to receive 32×32 input images, and they were trained with initial weights from Imagenet, learning rate $1e - 2$, cross-entropy loss and the stochastic gradient descent optimizer (code from [23]).

Table 4: Number of images and accuracy on train and validation sets for ResNet18 and VGG-16 models. We train the models with two different dataset: Cat vs. Dog and CIFAR10.

| | | Train | | Val. | |
|-------------|----------|-------------|----------|-------------|----------|
| | | Num. images | Acc. (%) | Num. images | Acc. (%) |
| Cat vs. Dog | ResNet18 | 19,891 | 98.21 | 5,109 | 97.86 |
| | VGG-16 | | 99.04 | | 98.61 |
| CIFAR10 | ResNet18 | 50,000 | 99.56 | 10,000 | 92.53 |
| | VGG-16 | | 99.84 | | 93.54 |

Table 5: Percentage of images with the original class changed after the **exclusion** of selected explanation regions for Cat vs. Dog dataset. Highlighted in blue are the configurations presented in the main paper. We tested hierarchies constructed by filtering out smaller regions than 200, 300, 400 and 500 pixels, segmentation based on **Edges**, Integrated Gradients (**IG**), Guided-Backpropagation (**BP**), Input X Gradients (**I X G**) and **Saliency**. We tested three different strategies to for the first hierarchical segmentation: BPT, watershed with area attribute, and watershed with volume attribute. **Same** column shows images maintaining the original class when the output was reduced, and **Total** is the sum of class change (**Ch.**) and class reduction (**Same**).

| % of images | | Cat vs. Dog | | | | | | | | | | | | | | | | | | | | |
|-------------|------------------|-------------|------|------|------|------|------|-------|------|----------|------|--------|------|------|------|------|------|-------|------|----------|------|------|
| | | VGG | | | | | | | | | | ResNet | | | | | | | | | | |
| | | Edges | | IG | | BP | | I X G | | Saliency | | Edges | | IG | | BP | | I X G | | Saliency | | |
| Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | | | |
| 200 | BPT | CaOC | 0.35 | 0.39 | 0.32 | 0.15 | 0.56 | 0.32 | 0.26 | 0.12 | 0.46 | 0.38 | 0.27 | 0.42 | 0.12 | 0.29 | 0.39 | 0.10 | 0.21 | 0.28 | 0.43 | |
| | | Occ | 0.51 | 0.44 | 0.33 | 0.23 | 0.63 | 0.35 | 0.27 | 0.19 | 0.64 | 0.35 | 0.41 | 0.52 | 0.17 | 0.34 | 0.55 | 0.37 | 0.13 | 0.30 | 0.39 | 0.55 |
| | Watershed area | CaOC | 0.15 | 0.48 | 0.23 | 0.45 | 0.21 | 0.46 | 0.22 | 0.46 | 0.15 | 0.49 | 0.20 | 0.44 | 0.18 | 0.45 | 0.16 | 0.50 | 0.17 | 0.46 | 0.17 | 0.46 |
| | | Occ | 0.30 | 0.66 | 0.42 | 0.55 | 0.41 | 0.56 | 0.42 | 0.55 | 0.31 | 0.62 | 0.34 | 0.61 | 0.31 | 0.63 | 0.32 | 0.64 | 0.30 | 0.63 | 0.31 | 0.63 |
| | Watershed volume | CaOC | 0.18 | 0.48 | 0.16 | 0.49 | 0.12 | 0.46 | 0.18 | 0.47 | 0.13 | 0.50 | 0.18 | 0.44 | 0.17 | 0.46 | 0.12 | 0.49 | 0.17 | 0.45 | 0.18 | 0.45 |
| | | Occ | 0.33 | 0.64 | 0.33 | 0.61 | 0.31 | 0.64 | 0.37 | 0.57 | 0.32 | 0.63 | 0.33 | 0.62 | 0.29 | 0.65 | 0.25 | 0.69 | 0.30 | 0.65 | 0.32 | 0.64 |
| 300 | BPT | CaOC | 0.36 | 0.38 | 0.22 | 0.07 | 0.56 | 0.30 | 0.17 | 0.07 | 0.46 | 0.38 | 0.28 | 0.41 | 0.07 | 0.19 | 0.39 | 0.37 | 0.06 | 0.13 | 0.30 | 0.41 |
| | | Occ | 0.50 | 0.46 | 0.23 | 0.13 | 0.61 | 0.34 | 0.17 | 0.11 | 0.58 | 0.41 | 0.41 | 0.53 | 0.09 | 0.27 | 0.48 | 0.39 | 0.07 | 0.20 | 0.38 | 0.53 |
| | Watershed area | CaOC | 0.16 | 0.49 | 0.25 | 0.43 | 0.22 | 0.45 | 0.24 | 0.46 | 0.15 | 0.50 | 0.20 | 0.42 | 0.20 | 0.44 | 0.17 | 0.47 | 0.18 | 0.47 | 0.19 | 0.45 |
| | | Occ | 0.30 | 0.65 | 0.43 | 0.53 | 0.42 | 0.55 | 0.42 | 0.55 | 0.31 | 0.63 | 0.35 | 0.59 | 0.30 | 0.63 | 0.32 | 0.62 | 0.30 | 0.63 | 0.31 | 0.62 |
| | Watershed volume | CaOC | 0.17 | 0.51 | 0.18 | 0.49 | 0.15 | 0.43 | 0.19 | 0.48 | 0.14 | 0.51 | 0.20 | 0.41 | 0.18 | 0.44 | 0.12 | 0.48 | 0.18 | 0.45 | 0.18 | 0.44 |
| | | Occ | 0.32 | 0.63 | 0.34 | 0.60 | 0.30 | 0.64 | 0.36 | 0.57 | 0.32 | 0.61 | 0.34 | 0.61 | 0.29 | 0.64 | 0.26 | 0.69 | 0.30 | 0.64 | 0.32 | 0.63 |
| 400 | BPT | CaOC | 0.36 | 0.39 | 0.15 | 0.05 | 0.51 | 0.29 | 0.10 | 0.05 | 0.43 | 0.40 | 0.29 | 0.42 | 0.04 | 0.16 | 0.35 | 0.37 | 0.04 | 0.11 | 0.29 | 0.41 |
| | | Occ | 0.47 | 0.48 | 0.15 | 0.09 | 0.54 | 0.37 | 0.10 | 0.08 | 0.51 | 0.47 | 0.41 | 0.52 | 0.06 | 0.23 | 0.42 | 0.39 | 0.05 | 0.14 | 0.36 | 0.55 |
| | Watershed area | CaOC | 0.16 | 0.49 | 0.26 | 0.46 | 0.25 | 0.43 | 0.24 | 0.46 | 0.17 | 0.50 | 0.21 | 0.41 | 0.20 | 0.44 | 0.18 | 0.46 | 0.20 | 0.46 | 0.21 | 0.42 |
| | | Occ | 0.30 | 0.64 | 0.43 | 0.53 | 0.42 | 0.54 | 0.42 | 0.55 | 0.31 | 0.62 | 0.34 | 0.61 | 0.32 | 0.61 | 0.32 | 0.61 | 0.30 | 0.63 | 0.31 | 0.63 |
| | Watershed volume | CaOC | 0.18 | 0.49 | 0.20 | 0.47 | 0.16 | 0.45 | 0.22 | 0.45 | 0.16 | 0.50 | 0.21 | 0.40 | 0.19 | 0.45 | 0.13 | 0.47 | 0.19 | 0.42 | 0.20 | 0.44 |
| | | Occ | 0.33 | 0.63 | 0.35 | 0.60 | 0.30 | 0.64 | 0.38 | 0.57 | 0.31 | 0.62 | 0.34 | 0.61 | 0.30 | 0.64 | 0.25 | 0.69 | 0.29 | 0.63 | 0.33 | 0.63 |
| 500 | BPT | CaOC | 0.35 | 0.40 | 0.11 | 0.04 | 0.45 | 0.29 | 0.07 | 0.03 | 0.40 | 0.42 | 0.30 | 0.41 | 0.04 | 0.13 | 0.31 | 0.34 | 0.04 | 0.11 | 0.29 | 0.42 |
| | | Occ | 0.45 | 0.49 | 0.12 | 0.06 | 0.49 | 0.36 | 0.07 | 0.04 | 0.47 | 0.49 | 0.41 | 0.51 | 0.04 | 0.17 | 0.38 | 0.38 | 0.04 | 0.11 | 0.37 | 0.54 |
| | Watershed area | CaOC | 0.16 | 0.48 | 0.29 | 0.46 | 0.26 | 0.42 | 0.25 | 0.47 | 0.18 | 0.48 | 0.22 | 0.41 | 0.21 | 0.45 | 0.20 | 0.46 | 0.20 | 0.46 | 0.21 | 0.42 |
| | | Occ | 0.31 | 0.63 | 0.43 | 0.54 | 0.41 | 0.55 | 0.41 | 0.54 | 0.31 | 0.63 | 0.35 | 0.60 | 0.32 | 0.61 | 0.34 | 0.61 | 0.30 | 0.62 | 0.30 | 0.66 |
| | Watershed volume | CaOC | 0.19 | 0.48 | 0.20 | 0.47 | 0.16 | 0.44 | 0.22 | 0.45 | 0.18 | 0.51 | 0.22 | 0.39 | 0.20 | 0.46 | 0.14 | 0.47 | 0.18 | 0.43 | 0.22 | 0.43 |
| | | Occ | 0.33 | 0.63 | 0.34 | 0.61 | 0.29 | 0.66 | 0.38 | 0.56 | 0.32 | 0.60 | 0.33 | 0.61 | 0.30 | 0.64 | 0.25 | 0.68 | 0.29 | 0.64 | 0.32 | 0.63 |

Table 6: Percentage of images with the original class changed after the **exclusion** of selected explanation regions for CIFAR10 dataset. Highlighted in blue are the configurations presented in the main paper. We tested hierarchies constructed by filtering out smaller regions than 200, 300, 400 and 500 pixels, segmentation based on **Edges**, Integrated Gradients (**IG**), Guided-Backpropagation (**BP**), Input X Gradients (**I X G**) and **Saliency**. We tested three different strategies to for the first hierarchical segmentation: BPT, watershed with area attribute, and watershed with volume attribute. **Same** column shows images maintaining the original class when the output was reduced, and **Total** is the sum of class change (**Ch.**) and class reduction (**Same**).

| % of images | | CIFAR10 | | | | | | | | | | | | | | | | | | | | |
|-------------|------------------|---------|------|------|------|------|------|-------|------|----------|------|--------|------|------|------|------|------|-------|------|----------|------|------|
| | | VGG | | | | | | | | | | ResNet | | | | | | | | | | |
| | | Edges | | IG | | BP | | I X G | | Saliency | | Edges | | IG | | BP | | I X G | | Saliency | | |
| Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | Ch. | Same | | | |
| 4 | BPT | CaOC | 0.11 | 0.34 | 0.07 | 0.44 | 0.08 | 0.44 | 0.07 | 0.43 | 0.14 | 0.44 | 0.15 | 0.32 | 0.12 | 0.43 | 0.14 | 0.42 | 0.12 | 0.43 | 0.19 | 0.43 |
| | | Occ | 0.44 | 0.28 | 0.50 | 0.26 | 0.49 | 0.27 | 0.51 | 0.25 | 0.46 | 0.30 | 0.39 | 0.25 | 0.41 | 0.24 | 0.41 | 0.25 | 0.41 | 0.24 | 0.39 | 0.28 |
| | Watershed area | CaOC | 0.12 | 0.41 | 0.18 | 0.42 | 0.18 | 0.42 | 0.18 | 0.42 | 0.18 | 0.42 | 0.17 | 0.40 | 0.21 | 0.42 | 0.21 | 0.42 | 0.21 | 0.42 | 0.21 | 0.41 |
| | | Occ | 0.38 | 0.33 | 0.39 | 0.34 | 0.39 | 0.34 | 0.40 | 0.34 | 0.38 | 0.35 | 0.34 | 0.30 | 0.33 | 0.33 | 0.34 | 0.33 | 0.34 | 0.33 | 0.33 | 0.34 |
| | Watershed volume | CaOC | 0.12 | 0.41 | 0.18 | 0.42 | 0.18 | 0.42 | 0.18 | 0.42 | 0.18 | 0.42 | 0.17 | 0.40 | 0.21 | 0.42 | 0.21 | 0.42 | 0.21 | 0.42 | 0.21 | 0.42 |
| | | Occ | 0.39 | 0.33 | 0.39 | 0.35 | 0.38 | 0.35 | 0.39 | 0.34 | 0.38 | 0.35 | 0.35 | 0.30 | 0.33 | 0.33 | 0.33 | 0.33 | 0.34 | 0.33 | 0.33 | 0.33 |
| 16 | BPT | CaOC | 0.06 | 0.16 | 0.05 | 0.20 | 0.06 | 0.26 | 0.05 | 0.18 | 0.13 | 0.39 | 0.09 | 0.15 | 0.08 | 0.21 | 0.11 | 0.25 | 0.07 | 0.18 | 0.18 | 0.39 |
| | | Occ | 0.29 | 0.20 | 0.37 | 0.22 | 0.42 | 0.25 | 0.33 | 0.19 | 0.42 | 0.31 | 0.26 | 0.17 | 0.30 | 0.18 | 0.34 | 0.22 | 0.27 | 0.16 | 0.37 | 0.29 |
| | Watershed area | CaOC | 0.12 | 0.40 | 0.15 | 0.43 | 0.15 | 0.43 | 0.16 | 0.43 | 0.16 | 0.43 | 0.16 | 0.39 | 0.20 | 0.42 | 0.20 | 0.42 | 0.21 | 0.42 | 0.21 | 0.42 |
| | | Occ | 0.38 | 0.33 | 0.40 | 0.33 | 0.39 | 0.33 | 0.40 | 0.33 | 0.38 | 0.35 | 0.34 | 0.30 | 0.35 | 0.32 | 0.35 | 0.32 | 0.35 | 0.32 | 0.34 | 0.33 |
| | Watershed volume | CaOC | 0.12 | 0.40 | 0.15 | 0.43 | 0.15 | 0.43 | 0.16 | 0.43 | 0.16 | 0.43 | 0.16 | 0.39 | 0.21 | 0.42 | 0.20 | 0.42 | 0.21 | 0.42 | 0.21 | 0.42 |
| | | Occ | 0.39 | 0.32 | 0.39 | 0.34 | 0.37 | 0.35 | 0.40 | 0.34 | 0.38 | 0.35 | 0.35 | 0.30 | 0.34 | 0.32 | 0.33 | 0.33 | 0.34 | 0.32 | 0.34 | 0.33 |
| 32 | BPT | CaOC | 0.03 | 0.06 | 0.03 | 0.10 | 0.05 | 0.16 | 0.03 | 0.10 | 0.10 | 0.25 | 0.05 | 0.07 | 0.06 | 0.11 | 0.08 | 0.16 | 0.05 | 0.10 | 0.14 | 0.25 |
| | | Occ | 0.20 | 0.12 | 0.24 | 0.14 | 0.33 | 0.21 | 0.21 | 0.11 | 0.36 | 0.28 | 0.19 | 0.10 | 0.20 | 0.12 | 0.27 | 0.17 | 0.18 | 0.10 | 0.33 | 0.26 |
| | Watershed area | CaOC | 0.11 | 0.36 | 0.14 | 0.43 | 0.13 | 0.43 | 0.14 | 0.43 | 0.14 | 0.43 | 0.15 | 0.35 | 0.20 | 0.42 | 0.19 | 0.42 | 0.20 | 0.42 | 0.20 | 0.42 |
| | | Occ | 0.37 | 0.32 | 0.40 | 0.33 | 0.39 | 0.33 | 0.40 | 0.33 | 0.38 | 0.35 | 0.33 | 0.29 | 0.35 | 0.31 | 0.35 | 0.32 | 0.35 | 0.32 | 0.33 | 0.33 |
| | Watershed volume | CaOC | 0.11 | 0.36 | 0.14 | 0.44 | 0.13 | 0.44 | 0.14 | 0.43 | 0.14 | 0.43 | 0.15 | 0.35 | 0.20 | 0.42 | 0.19 | 0.42 | 0.20 | 0.42 | 0.20 | 0.42 |
| | | Occ | 0.37 | 0.32 | 0.38 | 0.34 | 0.36 | 0.35 | 0.39 | 0.34 | 0.38 | 0.34 | 0.34 | 0.28 | 0.34 | 0.32 | 0.33 | 0.33 | 0.35 | 0.32 | 0.34 | 0.33 |
| 64 | BPT | CaOC | 0.01 | 0.03 | 0.01 | 0.05 | 0.02 | 0.08 | 0.01 | 0.05 | 0.05 | 0.11 | 0.03 | 0.03 | 0.03 | 0.06 | 0.06 | 0.08 | 0.03 | 0.05 | 0.07 | 0.11 |
| | | Occ | 0.13 | 0.07 | 0.12 | 0.07 | 0.22 | 0.13 | 0.10 | 0.05 | 0.25 | 0.17 | 0.12 | 0.06 | 0.12 | 0.06 | 0.19 | 0.11 | 0.10 | 0.05 | 0.24 | 0.16 |
| | Watershed area | CaOC | 0.08 | 0.23 | 0.13 | 0.36 | 0.12 | 0.35 | 0.13 | 0.37 | 0.14 | 0.38 | 0.11 | 0.22 | 0.20 | 0.42 | 0.18 | 0.37 | 0.18 | 0.36 | 0.18 | 0.38 |
| | | Occ | 0.32 | 0.26 | 0.39 | 0.32 | 0.38 | 0.33 | 0.39 | 0.32 | 0.37 | 0.34 | 0.28 | 0.22 | 0.35 | 0.30 | 0.35 | 0.30 | 0.35 | 0.31 | 0.33 | 0.32 |
| | Watershed volume | CaOC | 0.08 | 0.22 | 0.13 | 0.38 | 0.12 | 0.37 | 0.13 | 0.37 | 0.14 | 0.36 | 0.11 | 0.21 | 0.18 | 0.36 | 0.18 | 0.37 | 0.18 | 0.36 | 0.17 | 0.36 |
| | | Occ | 0.32 | 0.26 | 0.37 | 0.33 | 0.35 | 0.34 | 0.38 | 0.33 | 0.36 | 0.34 | 0.29 | 0.22 | 0.34 | 0.31 | 0.33 | 0.32 | 0.35 | 0.31 | 0.33 | 0.32 |

Exclusion of important regions: Given that each region-based explainable AI (xAI) method identifies important regions in an image to explain the prediction of a model, we performed occlusion of these regions to measure the impact of each selection and evaluate the methods. For methods that assign scores to regions, we masked the 25% highest scores.

We present, in Tables 5 and 6, the complete experiments of different configurations of our framework for the datasets Cat vs. Dog and CIFAR10 respectively.

PIR values: To address the issue of unhelpful explanations resulting from methods selecting the entire image as important, potentially leading to class changes upon occlusion, we introduce a novel metric termed **Pixel Impact Rate (PIR)**. This metric quantifies the impact on class activation per occluded pixel. Complementing the percentage of class change, PIR distinguishes whether changes

are primarily caused by complete or near-complete occlusion of the image. Higher PIR values indicate that each occluded pixel has a significant average impact, suggesting that concealing larger portions or the entire image leads to lower PIR, indicating less precision in the concealed area.

Tables 7 and 8 display for each network, and tested configurations of our framework, the average (**avg**) and standard deviation (**std**) of PIR, for the datasets Cat vs. Dog and CIFAR10 respectively.

Inclusion of important regions: Additional experimentation was conducted to demonstrate a method’s capability to identify an image region with sufficient information for the original class. The goal of this experiment is to determine whether the selected important region, when the only one left unoccluded in the image, can maintain the classification in its expected class. This experiment elucidates the critical role of these identified regions, providing strong evidence that they indeed contain essential information for accurate classification. We occluded **all regions** in the images except for the one selected by each method. We then calculated the percentage of images that changed class. We present the results from both datasets, Cat vs. Dog and CIFAR10, and all the tested framework configurations in Tables 9 and 10, respectively. Lower percentages indicate better performance, as they mean that a smaller percentage of images changed class, demonstrating that the chosen regions were sufficient to preserve the class for most of the images.

Table 8: Pixel Impact Rate (PIR) of selected explanation regions for CIFAR10 dataset. Highlighted in blue are the configurations presented in the main paper. The metric is the rate of impact under occlusion (difference between the original class output and the output under occlusion) by the number of pixels of the occlusion mask. We tested hierarchies constructed by filtering out smaller regions than 200, 300, 400 and 500 pixels, segmentation based on **Edges**, Integrated Gradients (**IG**), Guided-Backpropagation (**BP**), Input X Gradients (**IXG**) and **Saliency**. We tested three different strategies to for the first hierarchical segmentation: BPT, watershed with area attribute, and watershed with volume attribute. We expect higher values, on average (**avg**), for PIR, meaning each occluded pixel has a high impact.

| PIR | | CIFAR10 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|-----|------------------|---------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|-----|-----|-----|-----|-----|-----|----------|--|--|--|--|--|--|--|--|--|--|--|--|--|
| | | VGG | | | | | | | | | ResNet | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | | Edges | | | IG | | | BP | | | IXG | | | Saliency | | | Edges | | | IG | | | BP | | | IXG | | | Saliency | | | | | | | | | | | | | |
| avg | std | std | avg | std | std | avg | std | std | avg | std | std | avg | std | std | avg | std | std | avg | std | std | avg | std | std | avg | std | std | avg | std | std | | | | | | | | | | | | | |
| 4 | BPT | CatOC | 2.23e-02 | 2.39e-01 | 2.64e-02 | 1.73e-01 | 1.68e-01 | 2.61e-02 | 1.69e-01 | 3.11e-03 | 2.26e-02 | 2.80e-02 | 2.82e-01 | 3.02e-02 | 1.84e-01 | 3.04e-02 | 2.07e-01 | 3.09e-02 | 1.94e-01 | 3.06e-03 | 2.05e-02 | | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 1.46e-02 | 2.42e-01 | 2.01e-02 | 2.04e-01 | 1.61e-02 | 1.78e-01 | 2.03e-02 | 2.09e-01 | 1.72e-03 | 1.29e-02 | 1.71e-02 | 2.66e-01 | 2.85e-02 | 2.35e-01 | 2.47e-02 | 2.31e-01 | 2.94e-02 | 2.40e-01 | 1.79e-03 | 1.30e-02 | | | | | | | | | | | | | | | | | | | | |
| | Watershed area | CatOC | 1.08e-02 | 1.60e-01 | 7.49e-03 | 4.60e-02 | 8.37e-02 | 5.20e-02 | 7.69e-03 | 4.86e-02 | 7.11e-03 | 4.60e-02 | 1.24e-02 | 1.72e-01 | 6.09e-03 | 3.99e-02 | 5.87e-03 | 3.75e-02 | 5.92e-03 | 3.72e-02 | 5.65e-03 | 3.77e-02 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 8.71e-03 | 1.65e-01 | 9.50e-03 | 5.63e-02 | 9.85e-03 | 6.05e-02 | 9.09e-03 | 5.52e-02 | 8.73e-03 | 5.21e-02 | 8.82e-03 | 1.61e-01 | 6.09e-03 | 4.04e-02 | 5.75e-03 | 3.85e-02 | 5.94e-03 | 4.05e-02 | 5.89e-03 | 3.91e-02 | | | | | | | | | | | | | | | | | | | | |
| | Watershed volume | CatOC | 1.08e-02 | 1.57e-01 | 6.94e-03 | 4.41e-02 | 7.47e-03 | 4.72e-02 | 7.58e-03 | 4.88e-02 | 6.46e-03 | 4.14e-02 | 1.18e-02 | 1.68e-01 | 5.53e-03 | 3.67e-02 | 5.52e-03 | 3.71e-02 | 5.94e-03 | 3.95e-02 | 5.32e-03 | 3.59e-02 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 7.60e-03 | 1.51e-01 | 9.22e-03 | 5.85e-02 | 9.56e-03 | 5.83e-02 | 8.82e-03 | 5.33e-02 | 7.97e-03 | 4.96e-02 | 7.71e-03 | 1.46e-01 | 5.77e-03 | 3.71e-02 | 5.76e-03 | 3.84e-02 | 5.82e-03 | 4.02e-02 | 5.23e-03 | 3.51e-02 | | | | | | | | | | | | | | | | | | | | |
| 16 | BPT | CatOC | 3.54e-02 | 3.20e-01 | 3.62e-02 | 2.79e-01 | 3.56e-02 | 2.80e-01 | 3.38e-02 | 2.78e-01 | 1.34e-02 | 1.89e-01 | 4.12e-02 | 3.48e-01 | 4.06e-02 | 4.21e-02 | 3.08e-01 | 3.90e-02 | 2.94e-01 | 1.17e-02 | 1.65e-01 | | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 8.41e-02 | 6.02e-01 | 8.33e-02 | 5.84e-01 | 5.23e-02 | 4.61e-01 | 9.64e-02 | 6.37e-01 | 3.27e-03 | 8.07e-02 | 7.71e-02 | 5.69e-01 | 8.00e-02 | 5.55e-01 | 6.14e-02 | 4.91e-01 | 8.62e-02 | 5.75e-01 | 2.69e-03 | 7.94e-02 | | | | | | | | | | | | | | | | | | | | |
| | Watershed area | CatOC | 1.18e-02 | 1.68e-01 | 3.01e-03 | 1.65e-02 | 3.04e-03 | 1.66e-02 | 2.72e-03 | 1.47e-02 | 2.91e-03 | 1.64e-02 | 1.32e-02 | 1.77e-01 | 2.92e-03 | 1.57e-02 | 2.70e-03 | 1.44e-02 | 2.79e-03 | 1.49e-02 | 2.78e-03 | 2.61e-02 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 9.00e-03 | 1.71e-01 | 3.86e-03 | 1.89e-02 | 3.77e-03 | 1.92e-02 | 3.91e-03 | 1.99e-02 | 4.01e-03 | 2.00e-02 | 9.88e-03 | 1.77e-01 | 3.33e-03 | 1.72e-02 | 3.13e-03 | 1.67e-02 | 3.31e-03 | 1.70e-02 | 3.38e-03 | 1.76e-02 | | | | | | | | | | | | | | | | | | | | |
| | Watershed volume | CatOC | 1.15e-02 | 1.62e-01 | 2.79e-03 | 1.95e-02 | 2.67e-03 | 1.53e-02 | 2.78e-03 | 1.59e-02 | 2.68e-03 | 1.57e-02 | 1.25e-02 | 1.73e-01 | 2.63e-03 | 1.45e-02 | 2.53e-03 | 1.41e-02 | 2.79e-03 | 1.54e-02 | 2.68e-03 | 2.59e-02 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 7.92e-03 | 1.57e-01 | 3.68e-03 | 1.87e-02 | 3.74e-03 | 1.95e-02 | 3.82e-03 | 1.95e-02 | 3.59e-03 | 1.84e-02 | 8.36e-03 | 1.58e-01 | 3.22e-03 | 1.69e-02 | 3.19e-03 | 1.73e-02 | 3.30e-03 | 1.70e-02 | 3.11e-03 | 1.65e-02 | | | | | | | | | | | | | | | | | | | | |
| 32 | BPT | CatOC | 1.76e-02 | 2.18e-01 | 2.56e-02 | 2.48e-01 | 3.62e-02 | 2.97e-01 | 2.33e-02 | 2.30e-01 | 4.33e-02 | 3.67e-01 | 3.04e-02 | 2.99e-01 | 3.37e-02 | 2.88e-01 | 4.48e-02 | 3.26e-01 | 3.26e-02 | 2.76e-01 | 4.30e-02 | 3.58e-01 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 1.07e-01 | 6.95e-01 | 1.19e-01 | 7.24e-01 | 9.65e-02 | 6.51e-01 | 1.11e-01 | 6.99e-01 | 4.22e-02 | 4.31e-01 | 9.23e-02 | 6.24e-01 | 9.76e-02 | 6.56e-01 | 9.12e-02 | 6.16e-01 | 8.89e-02 | 6.00e-01 | 3.43e-02 | 3.84e-01 | | | | | | | | | | | | | | | | | | | | |
| | Watershed area | CatOC | 1.99e-02 | 2.30e-01 | 3.11e-03 | 6.00e-02 | 3.40e-03 | 6.91e-02 | 2.76e-03 | 5.99e-02 | 2.17e-03 | 2.89e-02 | 2.15e-02 | 2.35e-01 | 3.13e-03 | 5.16e-02 | 3.24e-03 | 7.13e-02 | 3.13e-03 | 6.64e-02 | 2.51e-02 | 5.20e-02 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 1.58e-02 | 2.53e-01 | 2.50e-03 | 1.15e-02 | 2.44e-03 | 1.12e-02 | 2.42e-03 | 1.11e-02 | 2.60e-03 | 1.17e-02 | 1.96e-02 | 2.76e-01 | 2.28e-03 | 1.11e-02 | 2.10e-03 | 1.04e-02 | 2.31e-03 | 1.11e-02 | 2.37e-03 | 1.13e-02 | | | | | | | | | | | | | | | | | | | | |
| | Watershed volume | CatOC | 1.98e-02 | 2.29e-01 | 2.49e-03 | 4.73e-02 | 2.80e-03 | 6.48e-02 | 2.54e-03 | 5.29e-02 | 2.70e-03 | 6.54e-02 | 2.24e-02 | 2.42e-01 | 2.95e-03 | 5.87e-02 | 2.46e-03 | 5.74e-02 | 2.88e-03 | 6.65e-02 | 2.75e-03 | 5.77e-02 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 1.58e-02 | 2.54e-01 | 2.48e-03 | 1.13e-02 | 2.49e-03 | 1.19e-02 | 2.51e-03 | 1.14e-02 | 2.43e-03 | 1.12e-02 | 1.78e-02 | 2.60e-01 | 2.22e-03 | 1.07e-02 | 2.10e-03 | 1.03e-02 | 2.29e-03 | 1.10e-02 | 2.18e-03 | 1.05e-02 | | | | | | | | | | | | | | | | | | | | |
| 64 | BPT | CatOC | 6.47e-03 | 1.26e-01 | 1.25e-02 | 1.70e-01 | 2.49e-02 | 2.53e-01 | 1.28e-02 | 1.63e-01 | 3.40e-02 | 3.17e-01 | 1.86e-02 | 2.31e-01 | 2.38e-02 | 2.53e-01 | 3.53e-02 | 3.18e-01 | 2.11e-02 | 2.27e-01 | 4.60e-02 | 3.78e-01 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 8.85e-02 | 6.37e-01 | 8.74e-02 | 6.31e-01 | 1.13e-01 | 7.19e-01 | 7.68e-02 | 5.96e-01 | 1.04e-01 | 6.85e-01 | 7.92e-02 | 5.86e-01 | 7.71e-02 | 5.74e-01 | 9.51e-02 | 6.34e-01 | 6.25e-02 | 5.09e-01 | 9.35e-02 | 6.37e-01 | | | | | | | | | | | | | | | | | | | | |
| | Watershed area | CatOC | 3.77e-02 | 3.31e-01 | 2.12e-02 | 2.41e-01 | 2.34e-02 | 2.53e-01 | 2.03e-02 | 2.45e-01 | 1.47e-02 | 2.01e-01 | 4.07e-02 | 3.38e-01 | 2.29e-02 | 2.50e-01 | 2.32e-02 | 2.55e-01 | 1.79e-02 | 2.21e-01 | 1.56e-02 | 2.05e-01 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 6.05e-02 | 5.05e-01 | 7.01e-03 | 1.55e-01 | 7.12e-03 | 1.49e-01 | 5.97e-03 | 1.45e-01 | 5.55e-03 | 1.24e-01 | 5.83e-02 | 4.87e-01 | 7.21e-03 | 1.52e-01 | 1.52e-01 | 1.58e-01 | 5.50e-03 | 1.36e-01 | 7.09e-03 | 1.59e-01 | | | | | | | | | | | | | | | | | | | | |
| | Watershed volume | CatOC | 3.71e-02 | 3.24e-01 | 1.79e-02 | 2.25e-01 | 1.83e-02 | 2.25e-01 | 2.04e-02 | 2.49e-01 | 1.98e-02 | 2.35e-01 | 4.12e-02 | 3.40e-01 | 1.79e-02 | 2.20e-01 | 1.68e-02 | 2.10e-01 | 1.74e-02 | 2.16e-01 | 2.03e-02 | 2.37e-01 | | | | | | | | | | | | | | | | | | | | |
| | | Occ | 5.99e-02 | 4.97e-01 | 7.14e-03 | 1.53e-01 | 7.28e-03 | 1.57e-01 | 5.97e-03 | 1.34e-01 | 9.57e-03 | 1.86e-01 | 5.84e-02 | 4.90e-01 | 6.11e-03 | 1.40e-01 | 6.54e-03 | 1.47e-01 | 5.02e-03 | 1.23e-01 | 9.20e-03 | 1.90e-01 | | | | | | | | | | | | | | | | | | | | |

Table 9: Percentage of images with the original class changed after the **inclusion** (exclusively) of selected explanation regions for Cat vs. Dog dataset. Highlighted in blue are the configurations presented in the main paper. We tested hierarchies constructed by filtering out smaller regions than 200, 300, 400 and 500 pixels, segmentation based on **Edges**, Integrated Gradients (**IG**), Guided-Backpropagation (**BP**), Input X Gradients (**I X G**) and **Saliency**. We tested three different strategies to for the first hierarchical segmentation: BPT, watershed with area attribute, and watershed with volume attribute. We expect smaller rate values of class change.

| % of images | | | Cat vs. Dog | | | | | | | | | |
|-------------|------------------|------|-------------|------|------|-------|----------|--------|------|------|-------|----------|
| | | | VGG | | | | | ResNet | | | | |
| | | | Edges | IG | BP | I X G | Saliency | Edges | IG | BP | I X G | Saliency |
| 200 | BPT | CaOC | 0.16 | 0.49 | 0.38 | 0.49 | 0.37 | 0.22 | 0.45 | 0.28 | 0.47 | 0.35 |
| | | Occ | 0.11 | 0.18 | 0.04 | 0.24 | 0.11 | 0.19 | 0.40 | 0.18 | 0.43 | 0.31 |
| | Watershed area | CaOC | 0.30 | 0.45 | 0.39 | 0.46 | 0.42 | 0.34 | 0.47 | 0.43 | 0.45 | 0.49 |
| | | Occ | 0.22 | 0.22 | 0.26 | 0.24 | 0.31 | 0.30 | 0.49 | 0.41 | 0.49 | 0.52 |
| | Watershed volume | CaOC | 0.27 | 0.44 | 0.43 | 0.50 | 0.45 | 0.34 | 0.47 | 0.49 | 0.46 | 0.50 |
| | | Occ | 0.24 | 0.30 | 0.36 | 0.28 | 0.29 | 0.30 | 0.50 | 0.46 | 0.47 | 0.50 |
| 300 | BPT | CaOC | 0.19 | 0.49 | 0.37 | 0.49 | 0.36 | 0.22 | 0.49 | 0.28 | 0.49 | 0.30 |
| | | Occ | 0.10 | 0.27 | 0.05 | 0.35 | 0.10 | 0.17 | 0.46 | 0.22 | 0.48 | 0.25 |
| | Watershed area | CaOC | 0.30 | 0.43 | 0.39 | 0.44 | 0.43 | 0.33 | 0.46 | 0.40 | 0.44 | 0.47 |
| | | Occ | 0.18 | 0.20 | 0.25 | 0.24 | 0.29 | 0.24 | 0.46 | 0.40 | 0.43 | 0.47 |
| | Watershed volume | CaOC | 0.28 | 0.45 | 0.43 | 0.45 | 0.45 | 0.31 | 0.46 | 0.46 | 0.45 | 0.46 |
| | | Occ | 0.21 | 0.29 | 0.34 | 0.27 | 0.29 | 0.25 | 0.49 | 0.45 | 0.45 | 0.45 |
| 400 | BPT | CaOC | 0.21 | 0.49 | 0.39 | 0.49 | 0.36 | 0.22 | 0.50 | 0.29 | 0.50 | 0.29 |
| | | Occ | 0.10 | 0.34 | 0.07 | 0.40 | 0.11 | 0.18 | 0.48 | 0.25 | 0.49 | 0.27 |
| | Watershed area | CaOC | 0.26 | 0.41 | 0.38 | 0.44 | 0.42 | 0.31 | 0.47 | 0.38 | 0.44 | 0.46 |
| | | Occ | 0.17 | 0.19 | 0.21 | 0.22 | 0.29 | 0.24 | 0.43 | 0.36 | 0.40 | 0.45 |
| | Watershed volume | CaOC | 0.27 | 0.44 | 0.43 | 0.46 | 0.45 | 0.33 | 0.47 | 0.43 | 0.47 | 0.46 |
| | | Occ | 0.19 | 0.26 | 0.30 | 0.28 | 0.28 | 0.22 | 0.47 | 0.41 | 0.44 | 0.43 |
| 500 | BPT | CaOC | 0.20 | 0.49 | 0.41 | 0.49 | 0.35 | 0.22 | 0.50 | 0.31 | 0.50 | 0.29 |
| | | Occ | 0.08 | 0.38 | 0.11 | 0.43 | 0.12 | 0.16 | 0.49 | 0.29 | 0.50 | 0.26 |
| | Watershed area | CaOC | 0.26 | 0.41 | 0.37 | 0.44 | 0.41 | 0.32 | 0.43 | 0.35 | 0.41 | 0.46 |
| | | Occ | 0.17 | 0.19 | 0.22 | 0.22 | 0.28 | 0.23 | 0.42 | 0.32 | 0.40 | 0.42 |
| | Watershed volume | CaOC | 0.25 | 0.45 | 0.42 | 0.44 | 0.42 | 0.32 | 0.44 | 0.42 | 0.45 | 0.45 |
| | | Occ | 0.19 | 0.28 | 0.31 | 0.26 | 0.29 | 0.22 | 0.45 | 0.37 | 0.42 | 0.41 |

A.6 Qualitative analysis

Models’ description: Table 11 shows the number of images in the train set for three Cat vs. Dog ResNet18 biased models. For **Bias 1** the biased class is composed of only cats on top of cushions. For **Bias 2** the biased class is composed of only dogs next to grades. For **Bias 3** the biased class is composed of only cats with humans. We also include the accuracy percentage per class when predicting a non-biased validation set composed by 5,109 images.

The biased models were trained with initial weights from Imagenet, learning rate $5e-7$, cross-entropy loss, the Adam optimizer, and early stop in 20 epochs of non-improving validation loss.

Comparison of misclassified images: Considering the hierarchical characteristic of our methodology, we can perform a deeper analysis of the explanations by selecting regions by the percentage of importance to be visualized, as shown in Figure 4. In the last level of the dishwasher example, the model seems to focus on the cat’s dish after having focused on the sink (in the previous level).

Human evaluation in bias analysis: As mentioned on the paper, we used the configuration **C1** for human-interpretation evaluation compared to baseline techniques: IG, Grad-CAM, OCC, LIME, Ms-IV, ACE. We presented the same five image visualizations (from corrected classified images by the biased class) for the baseline methods and the methods from **C1**. We intended to verify if: (i) humans can detect the wrong focus given based on a class prediction (**Detection**); and (ii) humans can recognize which was the cause of the bias (**Identification**).

To test (i) and (ii), for each **Bias** type we produce for each of the xAI methods an explanation image. By presenting five image explanations (the same images) for each of the xAI methodologies, we asked volunteers, based on the explanations provided, what did they think the highlighted regions referred to. The five image explanations are presented in Figure 5 for each **Bias** type (1-(a), 2-(b), and 3-(c)).

Table 10: Percentage of images with the original class changed after the **inclusion** (exclusively) of selected explanation regions for CIFAR10 dataset. Highlighted in blue are the configurations presented in the main paper. We tested hierarchies constructed by filtering out smaller regions than 200, 300, 400 and 500 pixels, segmentation based on **Edges**, Integrated Gradients (**IG**), Guided-Backpropagation (**BP**), Input X Gradients (**I X G**) and **Saliency**. We tested three different strategies to for the first hierarchical segmentation: BPT, watershed with area attribute, and watershed with volume attribute. We expect smaller rate values of class change.

| % of images | | | CIFAR10 | | | | | | | | | |
|-------------|------------------|------|---------|------|------|-------|----------|--------|------|------|-------|----------|
| | | | VGG | | | | | ResNet | | | | |
| | | | Edges | IG | BP | I X G | Saliency | Edges | IG | BP | I X G | Saliency |
| 4 | BPT | CaOC | 0.42 | 0.45 | 0.46 | 0.45 | 0.47 | 0.41 | 0.45 | 0.44 | 0.44 | 0.45 |
| | | Occ | 0.26 | 0.21 | 0.24 | 0.17 | 0.36 | 0.29 | 0.29 | 0.32 | 0.28 | 0.38 |
| | Watershed area | CaOC | 0.43 | 0.47 | 0.47 | 0.47 | 0.47 | 0.43 | 0.47 | 0.46 | 0.46 | 0.46 |
| | | Occ | 0.31 | 0.42 | 0.42 | 0.42 | 0.43 | 0.35 | 0.44 | 0.43 | 0.44 | 0.44 |
| | Watershed volume | CaOC | 0.43 | 0.47 | 0.47 | 0.48 | 0.47 | 0.43 | 0.46 | 0.47 | 0.46 | 0.46 |
| | | Occ | 0.30 | 0.43 | 0.43 | 0.43 | 0.43 | 0.34 | 0.44 | 0.44 | 0.44 | 0.44 |
| 16 | BPT | CaOC | 0.44 | 0.45 | 0.45 | 0.45 | 0.43 | 0.43 | 0.44 | 0.44 | 0.44 | 0.43 |
| | | Occ | 0.32 | 0.23 | 0.24 | 0.25 | 0.31 | 0.34 | 0.30 | 0.32 | 0.32 | 0.35 |
| | Watershed area | CaOC | 0.43 | 0.46 | 0.46 | 0.46 | 0.46 | 0.43 | 0.45 | 0.46 | 0.45 | 0.45 |
| | | Occ | 0.31 | 0.41 | 0.39 | 0.40 | 0.41 | 0.34 | 0.42 | 0.42 | 0.42 | 0.42 |
| | Watershed volume | CaOC | 0.43 | 0.46 | 0.46 | 0.46 | 0.46 | 0.43 | 0.45 | 0.46 | 0.46 | 0.45 |
| | | Occ | 0.30 | 0.41 | 0.41 | 0.40 | 0.41 | 0.34 | 0.42 | 0.42 | 0.43 | 0.42 |
| 32 | BPT | CaOC | 0.46 | 0.46 | 0.45 | 0.46 | 0.42 | 0.44 | 0.45 | 0.44 | 0.45 | 0.42 |
| | | Occ | 0.36 | 0.33 | 0.29 | 0.34 | 0.33 | 0.37 | 0.36 | 0.34 | 0.37 | 0.35 |
| | Watershed area | CaOC | 0.43 | 0.45 | 0.45 | 0.45 | 0.45 | 0.43 | 0.44 | 0.45 | 0.45 | 0.45 |
| | | Occ | 0.31 | 0.38 | 0.39 | 0.38 | 0.39 | 0.34 | 0.40 | 0.41 | 0.41 | 0.41 |
| | Watershed volume | CaOC | 0.43 | 0.45 | 0.45 | 0.45 | 0.45 | 0.43 | 0.44 | 0.45 | 0.45 | 0.44 |
| | | Occ | 0.30 | 0.39 | 0.39 | 0.39 | 0.38 | 0.34 | 0.41 | 0.41 | 0.41 | 0.40 |
| 64 | BPT | CaOC | 0.47 | 0.46 | 0.46 | 0.47 | 0.44 | 0.45 | 0.46 | 0.45 | 0.46 | 0.43 |
| | | Occ | 0.40 | 0.41 | 0.37 | 0.42 | 0.38 | 0.40 | 0.41 | 0.39 | 0.42 | 0.38 |
| | Watershed area | CaOC | 0.43 | 0.43 | 0.43 | 0.44 | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 |
| | | Occ | 0.32 | 0.36 | 0.35 | 0.35 | 0.37 | 0.35 | 0.38 | 0.38 | 0.38 | 0.39 |
| | Watershed volume | CaOC | 0.43 | 0.43 | 0.43 | 0.44 | 0.43 | 0.43 | 0.43 | 0.44 | 0.43 | 0.43 |
| | | Occ | 0.32 | 0.37 | 0.38 | 0.36 | 0.37 | 0.35 | 0.39 | 0.40 | 0.39 | 0.38 |

Table 11: Number of images (for a normal and an induced biased class) for training three biased ResNet18 models. We also present the accuracy of the models when predicting each class image from a non-biased validation dataset (5,109 images).

| | Normal class | Acc. orig. val normal (%) | Bias class | Acc. orig. val bias (%) |
|--------|--------------|---------------------------|------------|-------------------------|
| Bias 1 | 138 | 86.91 | 69 | 84.82 |
| Bias 2 | 85 | 97.97 | 56 | 37.81 |
| Bias 3 | 161 | 86.28 | 46 | 80.96 |

Here, we display the text provided to the volunteers for this experiment:

| |
|---|
| <p>[FORM] Part I - Determining the focus of the images: For each question, we provide two rows of images:</p> <ul style="list-style-type: none"> • The first row displays the original images, each representing a specific class. • The second row showcases an image for each image from the first row, highlighting the important parts for the class. <p>[IMPORTANT] What is a class?</p> <p>A class refers to a category or type of object, animal, or characteristic depicted in the images. For instance, a class of cat images would include images featuring cats, while a class of dog images would comprise images featuring dogs. Similarly, a class of cartoon images would include images characterized by cartoon-like features. In essence, a class represents a distinct category used to classify and organize images based on their content or characteristics.</p> <p>Throughout the questions, our objective is to identify the common important parts present in the images of the first row, as indicated by the corresponding images in the second row.</p> <p>If no common important parts are identified for most of the images, the answer should be Not identified.</p> |
|---|

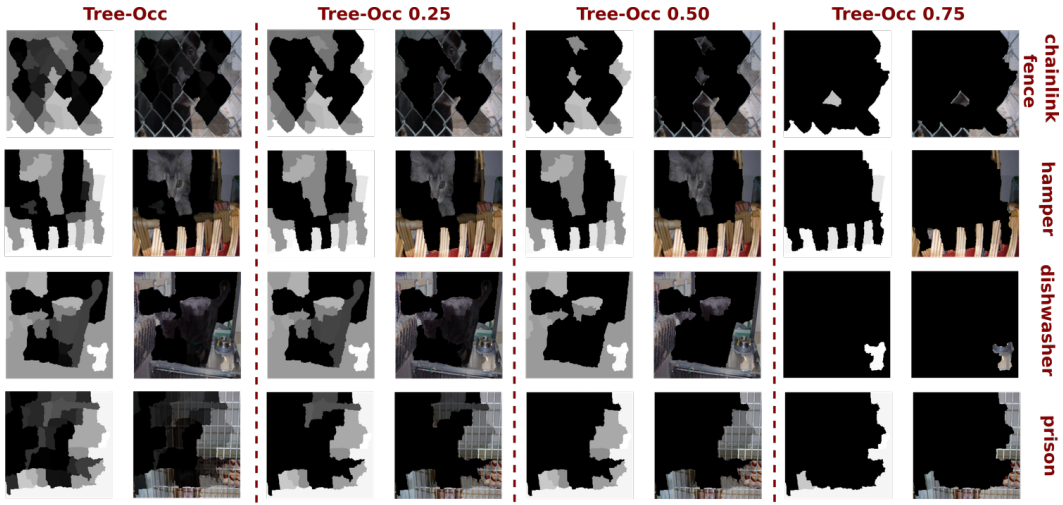


Figure 4: Different visualization levels on the explanation hierarchy. We illustrate a deeper analysis of the explanations of four images from Figure 1 using Tree-Occ (minimal region size of 500 pixels). We can note the evolution of the importance in the images’ shapes, for examples: in the hamper image, although the hamper is the most important, the cat has also an important that disappears at the more selective level (Tree-Occ 0.75); in the dishwasher the initial explanations show the sink as important but at the most selective level, the cat’s dish is the only one remaining. This analysis can be helpful to understand the reasoning behind predictions.

And for each method visualization:

For the following three questions, the second row of images displays significant image components to the class of the animal.
 What are the significant components of the images highlighted, as depicted in the second row of images?

To test ACE similarly as we did with the other methods, we highlight the top five concepts found (described as sufficient in the original ACE paper [11]) in the same five selected images. However, we also show the visualizations of the ten most activated images for the top five found concepts in Figure 6.

In our final qualitative experiment, using the same methods as the previous human evaluation, we presented four image explanation visualizations for non-biased models to determine xAI model preferences. The images are presented in Figure 7. We presented the following explanation and question:

[FORM] Part II: Choosing the best representation:
 For the next questions, you will be asked to answer which image number do you prefer to describe the class we indicate.
 You should choose the image that seems to highlight class features in an easier way to understand.
 Which image do you think better shows representative parts of the animal?

For the two first images (Figures 7 (a) and (b)), over 70% preferred Tree-Occ and Tree-CaOC over others. For the third image (Figure 7 (c)), IG was preferred by 26.5%, followed by Tree-OCC and Occlusion with 20.6%. Grad-CAM was preferred in the fourth image (Figure 7 (d)), with 60.6%, followed by Tree-CaOC with 18.2%. The visualizations suggest a preference for explanations that highlight the complete concept (cat or dog) rather than focusing on specific small animals’ regions.

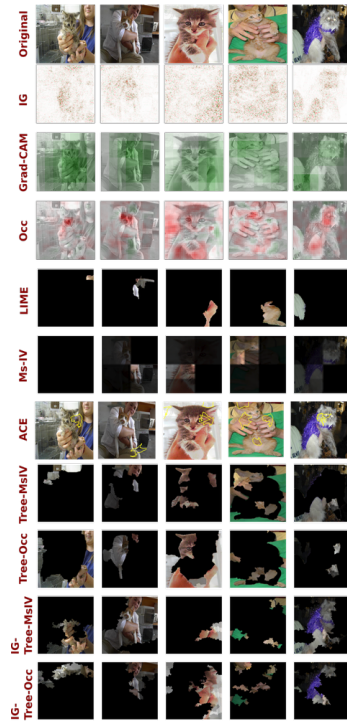
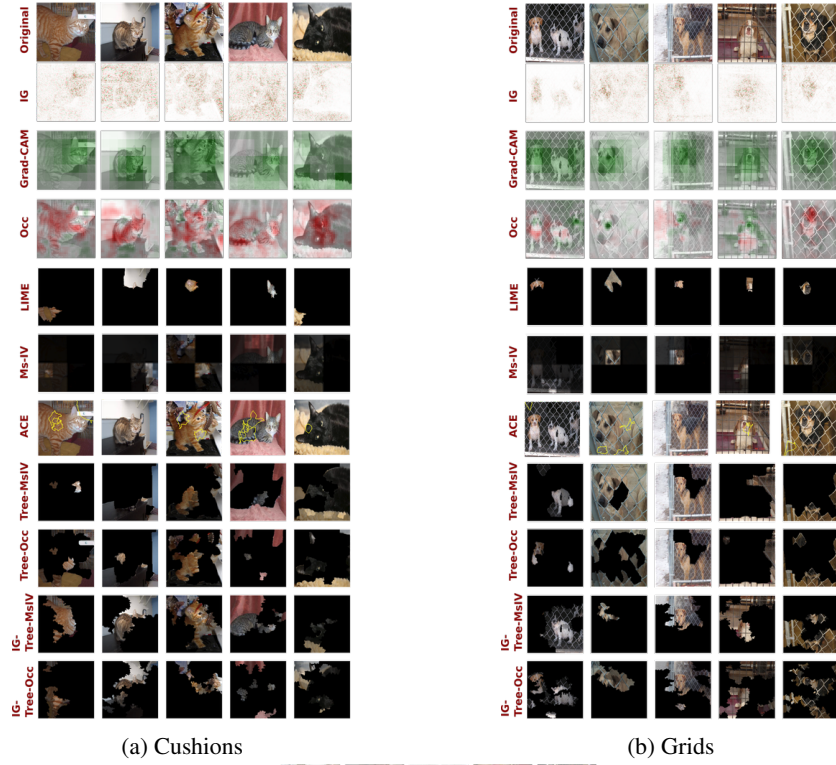


Figure 5: Explanations of visualizations used on our human-based evaluations for bias detection and identification, of all the ten compared methods: IG, Grad-CAM, OCC, LIME, Ms-IV, ACE, Tree-MsIV, Tree-Occ, IG-Tree-Msiv, and IG-Tree-Occ. We showed the same five image explanations for all the methods.

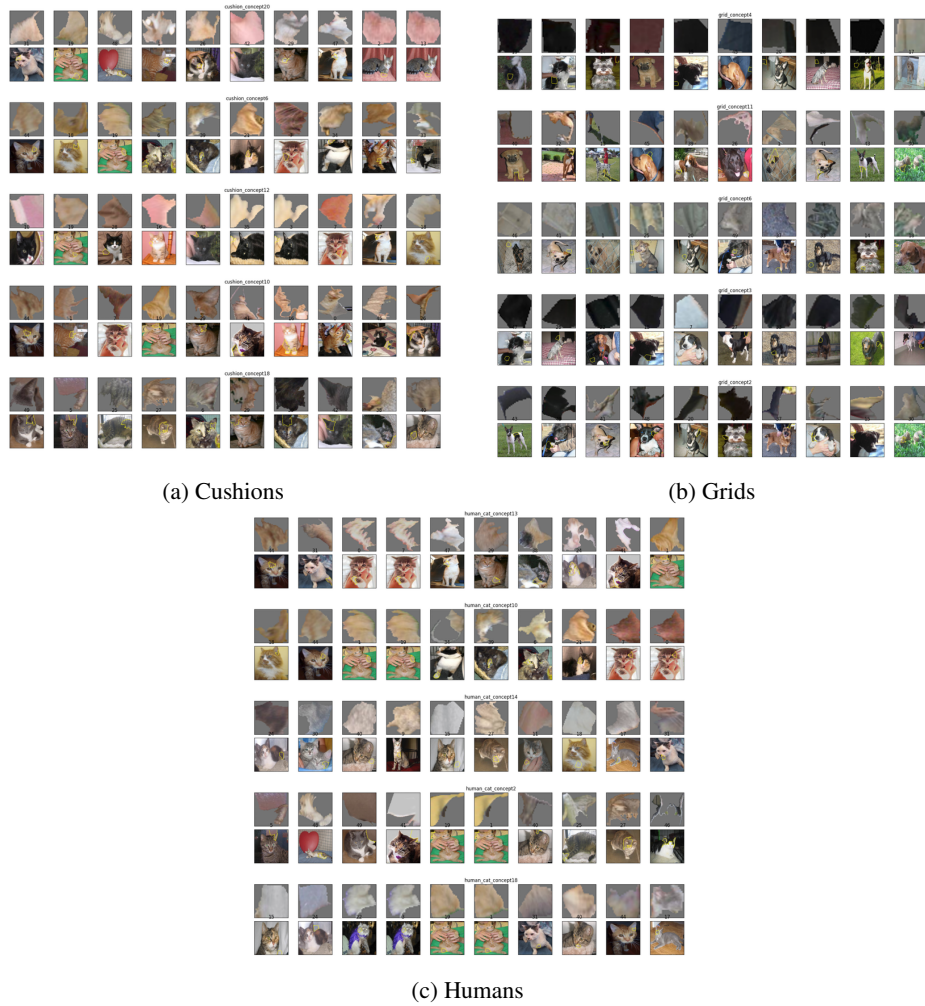


Figure 6: Original explanations of the top 5 concepts generated by ACE for the three biased models. Instead of showing the 10 most concepts’ activated images we draw these five top concepts on the 5 selected images from Figure 5 to have a fairer comparison with the other methods.

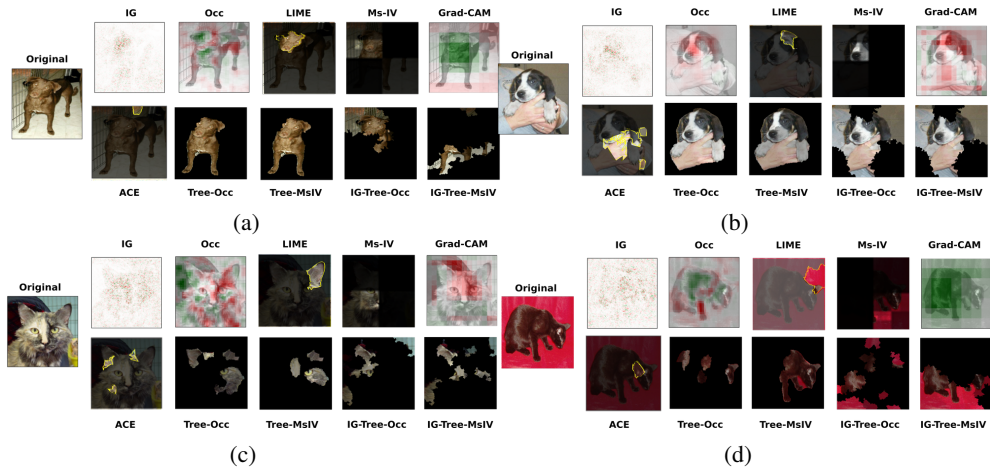


Figure 7: Explanations of visualizations used on our human-based evaluations for preference analysis, of all the ten compared methods: IG, Grad-CAM, OCC, LIME, Ms-IV, ACE, Tree-MsIV, Tree-Occ, IG-Tree-Msiv, and IG-Tree-Occ. We showed the same five image explanations for all the methods.