# Use HiResCAM instead of Grad-CAM for faithful explanations of convolutional neural networks

Rachel Lea Draelos Ph.D.[a] and Lawrence Carin Ph.D.[b]

[a]Corresponding author. rlb61@duke.edu. Duke University Computer Science
[b]lawrence.carin@kaust.edu.sa. Duke University Electrical and Computer Engineering

## Abstract

Explanation methods facilitate the development of models that learn meaningful concepts and avoid exploiting spurious correlations. We illustrate a previously unrecognized limitation of the popular neural network explanation method Grad-CAM: as a side effect of the gradient averaging step, Grad-CAM sometimes highlights locations the model did not actually use. To solve this problem, we propose HiResCAM, a novel class-specific explanation method that is guaranteed to highlight only the locations the model used to make each prediction. We prove that HiResCAM is a generalization of CAM and explore the relationships between HiResCAM and other gradient-based explanation methods. Experiments on PASCAL VOC 2012, including crowd-sourced evaluations, illustrate that while HiResCAM's explanations faithfully reflect the model, Grad-CAM often expands the attention to create bigger and smoother visualizations. Overall, this work advances convolutional neural network explanation approaches and may aid in the development of trustworthy models for sensitive applications.

## 1 Introduction

Machine learning models sometimes rely on spurious correlations, predicting horses from copyright watermarks [1] or pneumonia from metal tokens [2]. Even more alarming, models absorb prejudices present in training data and exhibit racial and gender bias [3, 4]. Explanation methods are one tool to identify and combat these dangerous properties. By improving understanding of how a model makes predictions, explanation methods can facilitate selection of more reasonable, fair models. For convolutional neural networks (CNNs) that process images, visual explanation methods are often used. These methods highlight the parts of an input image that contributed most to a particular prediction.

Gradient-based visual explanation methods for CNNs are popular due to their computational efficiency. Input-level gradient-based explanation methods [5, 6, 7, 8, 9] are quick, but produce explanations that are not class-specific in practice due to a "white noise" appearance. Output-level gradient-based methods such as Class Activation Mapping (CAM) [10] and Grad-CAM [11] are class-specific, and thus have been deployed in numerous settings including sensitive applications like medical imaging [12, 13, 14, 15, 16, 17]. Grad-CAM in particular is often used because its explanations can be calculated for any CNN architecture, while CAM is restricted to a subset of CNN architectures.

We identify a previously unreported limitation of Grad-CAM: due to its gradient averaging step, Grad-CAM is not guaranteed to reflect locations the model used for prediction, and therefore can produce misleading explanations. To solve this problem, we develop HiResCAM, an explanation method in the CAM family that is faithful to the model. Grad-CAM and HiResCAM explanations often differ noticeably as exemplified in Figure 1. In this work, we offer the following contributions:

Figure 1: Grad-CAM and HiResCAM produce different explanations for the same model, image, and class. HiResCAM provably reflects the locations the model used for computation, while Grad-CAM does not. (A) HiResCAM explanations are often more focal than Grad-CAM explanations. (B) Sometimes HiResCAM highlights the correct object while Grad-CAM does not; (C) other times, HiResCAM highlights more parts of the image than Grad-CAM. In this example, the Grad-CAM explanation gives the impression that the model predicted "potted plant" from the plant alone, but the HiResCAM explanation reveals that the model also attended to other parts of the image. Best viewed in color.

- We develop HiResCAM, a novel explanation method, and prove that HiResCAM explanations are guaranteed to reflect the locations the model used for any CNN ending in one fully connected layer, while the same guarantee does not hold for Grad-CAM;
- We prove that HiResCAM is a generalization of CAM and conceptually connect HiResCAM to other gradient-based explanation methods;
- In experiments on natural images, we quantify how Grad-CAM explanations deviate from the model's calculations, and show through examples and crowd-sourced assessment that Grad-CAM's explanations are often bigger and rounder than the faithful HiResCAM explanations.

## 2  Related Work

Our proposed attention mechanism, HiResCAM, is part of the family of gradient-based neural network explanation methods [18].

### 2.1  Input-level approaches

Saliency mapping [5], DeconvNets [6], and Guided Backpropagation [7] are gradient-based explanation methods that compute the gradient of the class score with respect to the input image to visualize

important image regions. These methods are identical except for handling of ReLU nonlinearities [19]. Gradient $*$ Input [8] is a related method in which the gradient of the class score with respect to the input is multiplied element-wise against the input itself. Layer-wise relevance propagation (LRP) [9] proceeds layer-by-layer, starting with the output, to redistribute the final score across the pixels of the input. While not originally formulated as a gradient-based explanation method, $\epsilon$-LRP is in fact equivalent to Gradient $*$ Input where the gradient is calculated in a modified manner using the ratio between the output and input at each nonlinearity [18]. A limitation of the aforementioned approaches is the "white noise" appearance of the final explanation caused by shattered gradients [20, 19], which prevents the explanation from being class-specific in practice.

## 2.2 Output-level approaches

Class Activation Mapping (CAM) [10] is an explanation method for a particular class of neural networks that consist of convolutions followed by global average pooling of feature maps and one final fully connected layer. CAM explanations are produced by multiplying the class-specific weights of the final layer by the corresponding feature maps prior to pooling. CAM may be considered a gradient-based method, as the final class-specific weights are the gradient of the score with respect to the feature maps. Grad-CAM [11] is a generalization of CAM, in which gradients are averaged over the spatial dimensions to produce importance weights. The Grad-CAM explanation is a sum of feature maps weighted by the importance weights. Grad-CAM was proposed with the intention of extending CAM explanations to a broader class of CNN architectures.

Unlike the input-level approaches, which rely on propagation through all layers back to the level of the original image, CAM and Grad-CAM produce explanations at a layer of the network closer to the output. The low-dimensional explanation is then upsampled for superimposition over the input image, an acceptable step because typical CNNs preserve spatial relationships. Guided Grad-CAM [11] is a Grad-CAM variant obtained via an element-wise product of the Guided Backpropagation [7] and Grad-CAM explanations.

Recent work has called into question some gradient-based explanation methods. Nie *et al.* [19] demonstrate that Guided Backpropagation and DeconvNets perform partial image recovery due to their handling of ReLU nonlinearities and max pooling. Adebayo *et al.* [21] reveal that Guided Backpropagation, Guided Grad-CAM, and Gradient $*$ Input produce convincing explanations even when model parameters have been randomized or when a model has been trained on randomly labeled data; however, saliency mapping and Grad-CAM pass their sanity checks.

Because they are are class-specific and produce less noisy explanations than input-level methods, CAM and Grad-CAM form the foundation of numerous weakly supervised localization methods [22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34]. CAM and Grad-CAM have also been widely used to explain model behavior, including in sensitive medical applications where protecting patients from biased, unreasonable models is particularly important [12, 13, 14, 15, 16, 17]. Unfortunately, as we will demonstrate, Grad-CAM is not a reliable explanation method and sometimes highlights locations the model did not use. We thus propose a new explanation method, HiResCAM, and prove that HiResCAM explanations faithfully reflect the model's computations.

## 3 Methods

### 3.1 Problem setup

Consider a CNN that takes in an input image $\mathbf{X}$. The goal of class-specific visual model explanation is to produce an attention map for a class $m$ with the same shape as input image $\mathbf{X}$ and values in $[0, 1]$ such that higher values indicate regions of the input image that increase the model's score for class $m$.

An output-level CNN explanation approach accomplishes this goal by applying the trained CNN with fixed parameters to an input image $\mathbf{X}$, in order to obtain an explanation at the level of some convolutional feature maps. The explanation is then upsampled and superimposed over $\mathbf{X}$ for visualization.
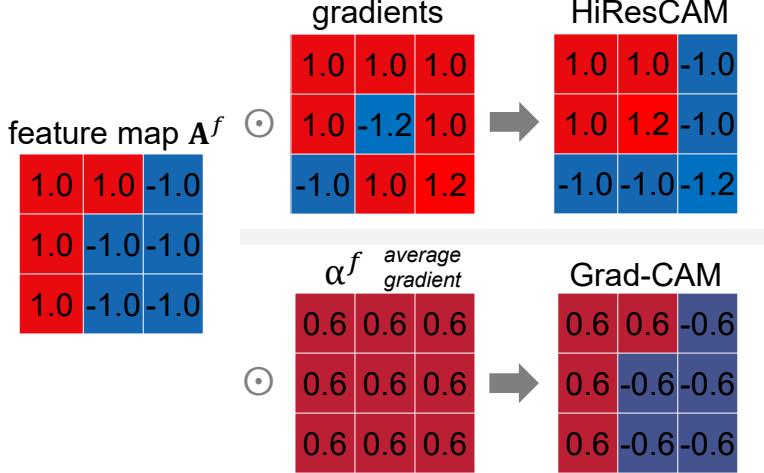
Figure 2: 2D example of how HiResCAM addresses the limitation of the gradient averaging step in Grad-CAM. The Grad-CAM explanation (equation 2) matches the relative magnitudes and positive-negative pattern of the original feature map (the "inverted red L shape" here), even though the gradients suggest that some elements should be re-scaled and/or change sign. HiResCAM (equation 3) does not average over the gradients and instead element-wise multiplies the feature map with the gradients directly, thereby producing attention that reflects the model's computations and emphasizes the most important locations for a particular prediction. Best viewed in color.

## 3.2   Motivation

We propose a new output-level, gradient-based CNN explanation approach, High-Resolution Class Activation Mapping (HiResCAM), which produces for every class $m = 1, ..., M$ a class-specific attention map $\tilde{\mathcal{A}}_m^{\text{HiResCAM}} \in \mathbb{R}^{D_1 \times D_2}$, $i.e.$ attention over the spatial dimensions $D_1$ and $D_2$ excluding the feature dimension $F$.

HiResCAM is inspired by the popular attention mechanism Grad-CAM [11], and is designed to address a limitation of the Grad-CAM averaging step. In Grad-CAM, feature map importance weights $\alpha^f$ are calculated by averaging gradients over the spatial dimensions of the low-dimensional CT representation. Such averaging is likely motivated by the global average pooling step built in to the architecture of Class Activation Mapping (CAM) [10]. However, the averaging limits the extent to which the final visualization depicts the locations within the image that the model is using to make predictions.

Figure 2 illustrates the fundamental problem: each $\alpha^f \mathbf{A}^f$ subcomponent of the final Grad-CAM explanation must always match the relative magnitudes of the feature map $\mathbf{A}^f$, and either (a) exactly match the positive-negative pattern of the feature map (when $\alpha^f$ is positive), or (b) invert the positive-negative pattern (when $\alpha^f$ is negative). Rescaling and sign changes of individual elements of the feature map are "blurred out." In HiResCAM, rescaling and sign changes are preserved, producing more high-resolution attention that reflects the model's computations.

## 3.3   Grad-CAM formulation

Before specifying the HiResCAM formulation, we review how Grad-CAM explanations are calculated. Define $s_m$ as a CNN's raw score for class $m$ before a sigmoid or softmax function is applied to produce predicted probabilities. To obtain a Grad-CAM explanation for class $m$, we first compute the gradient of $s_m$ with respect to a collection of feature maps $\mathbf{A} = \{\mathbf{A}^f\}_{f=1}^F$ produced by a convolutional layer. For 2D data, this gradient $\frac{\partial s_m}{\partial \mathbf{A}}$ is 3-dimensional, $[F, D_1, D_2]$, matching the shape of the collection of feature maps. Note that the collection of feature maps selected could be produced by the last convolutional layer, or an earlier convolutional layer, as the authors of Grad-CAM state that Grad-CAM explanations can be calculated at any convolutional layer [11].

After computing $\frac{\partial s_m}{\partial \mathbf{A}}$, we calculate a vector of importance weights [11] $\boldsymbol{\alpha}_m \in \mathbb{R}^F$, where each element $\alpha_m^f$ will be used to re-weight the corresponding feature map $\mathbf{A}^f$. The importance weights are obtained by global average pooling the gradient over the spatial dimensions:

$$\alpha_m^f = \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \frac{\partial s_m}{\partial \mathbf{A}_{d_1 d_2}^f}. \tag{1}$$

The importance weights indicate which features are most relevant to this particular class throughout the image overall. The final Grad-CAM explanation is produced as an importance-weighted combination of the feature maps:

$$\tilde{\mathcal{A}}_m^{\text{GradCAM}} = \sum_{f=1}^{F} \alpha_m^f \mathbf{A}^f. \tag{2}$$

Following standard practice for use of Grad-CAM [11], the attention map is then post-processed for better visualization by applying a ReLU and normalizing the attention values to the range $[0, 1]$. This step ensures that the regions positively associated with a class will be easily visible.

### 3.4 HiResCAM formulation

HiResCAM addresses the limitation of Grad-CAM illustrated in Figure 2. The first step of HiResCAM is the same as the first step of Grad-CAM: compute $\frac{\partial s_m}{\partial \mathbf{A}}$, the gradient of $s_m$ with respect to the feature maps $\mathbf{A}$.

In the second step of HiResCAM, the attention map is produced by element-wise multiplying the gradient and the feature maps before summing over the feature dimension:

$$\tilde{\mathcal{A}}_m^{\text{HiResCAM}} = \sum_{f=1}^{F} \frac{\partial s_m}{\partial \mathbf{A}^f} \odot \mathbf{A}^f. \tag{3}$$

The purpose of HiResCAM is to reflect the model's computations; therefore, when the gradients indicate that some elements of the feature map should be scaled or have their sign inverted, HiResCAM performs these operations. In contrast, Grad-CAM blurs the effect of the gradients across each feature map.

HiResCAM can be applied to 2D, 3D, or $n$-D CNNs; the formulation works for feature maps with any number of spatial dimensions.

### 3.5 CNNs ending in one fully connected layer

#### 3.5.1 HiResCAM highlights locations that increase the class score

HiResCAM can be calculated at any convolutional layer of a CNN. However, for an explanation with guaranteed properties, HiResCAM must be applied at the last convolutional layer of any CNN that ends in one fully connected layer, which includes many modern CNN architectures. In this section we prove that for such CNNs, HiResCAM has an intuitive interpretation: the resulting explanation is guaranteed to highlight all locations within the image that increase the class score. In contrast, Grad-CAM's explanations do not reflect calculation of the class score.

*Proof.* Consider a 2D CNN constructed from the following layers:

1. Convolutional layers: $\text{conv}(\mathbf{X}) = \{\mathbf{A}^f\}_{f=1}^F$. Previously we used $\{\mathbf{A}^f\}_{f=1}^F$ to denote convolutional feature maps from any layer; now we consider $\{\mathbf{A}^f\}_{f=1}^F$ as the feature maps of specifically the last convolutional layer.
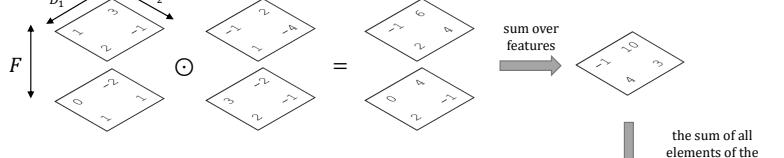
2. One final fully connected layer:

$$\mathbf{s} = \mathbf{W}\mathbf{A} + \mathbf{b}, \tag{4}$$

where $\mathbf{s} \in \mathbb{R}^M$ are the raw scores for the $M$ classes, the weight matrix is $\mathbf{W} \in \mathbb{R}^{M \times FD_1 D_2}$, and $\mathbf{A} \in \mathbb{R}^{FD_1 D_2 \times 1}$ is the convolutional output $\{\mathbf{A}^f\}_{f=1}^F$ flattened. If we consider only a

5

**HiResCAM calculation**

$$\tilde{\mathcal{A}}_m^{\text{HiResCAM}} = \sum_{f=1}^{F} \frac{\partial s_m}{\partial \mathbf{A}^f} \odot \mathbf{A}^f$$

$$\frac{\partial s_m}{\partial \mathbf{A}} = \mathbf{w}_m \in \mathbb{R}^{F \times D_1 \times D_2} \qquad \mathbf{A} \in \mathbb{R}^{F \times D_1 \times D_2} \qquad \mathbf{w}_m \odot \mathbf{A} \qquad \tilde{\mathcal{A}}_m^{\text{HiResCAM}} = \sum_{f=1}^{F} \mathbf{w}_m \odot \mathbf{A}$$

sum over features

the sum of all elements of the HiResCAM explanation is the input-dependent part of the abnormality score

**Abnormality score calculation**

$$s_m = \mathbf{w}_m \mathbf{A} + b_m$$

$$\mathbf{A} \in \mathbb{R}^{FD_1 D_2 \times 1}$$

$$\mathbf{w}_m \in \mathbb{R}^{1 \times FD_1 D_2}$$

$$\begin{bmatrix} 1 & 3 & 2 & -1 & 0 & -2 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} -1 \\ 2 \\ 1 \\ -4 \\ 3 \\ -2 \\ 2 \\ -1 \end{bmatrix} = \begin{matrix} (1)(-1)+(3)(2)+(2)(1)+(-1)(-4) \\ +(0)(3)+(-2)(-2)+(1)(2)+(1)(-1) \end{matrix} = \mathbf{w}_m \mathbf{A} = 16$$

$$s_m = \boxed{16} + b_m$$

**Grad-CAM calculation**

$$\tilde{\mathcal{A}}_m^{\text{GradCAM}} = \sum_{f=1}^{F} \alpha_m^f \mathbf{A}^f$$

$$\alpha_m^f = \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \mathbf{w}_m \qquad \mathbf{A} \in \mathbb{R}^{F \times D_1 \times D_2} \qquad \alpha_m^f \mathbf{A}^f \qquad \tilde{\mathcal{A}}_m^{\text{GradCAM}}$$

$$\alpha_m^1 = 1.25$$
$$\alpha_m^2 = 0$$

sum over features

the sum of all elements of the Grad-CAM explanation is NOT a part of the abnormality score

Figure 3: Specific example demonstrating that for CNNs ending in a single fully connected layer, HiResCAM explanations directly reflect the calculation of the class score while Grad-CAM explanations do not. Integer input values were chosen for simplicity; actual weights and activations are not integers.

single class $m$, and extract out the relevant row of weights $\mathbf{w}_m$ and bias $b_m$, the expression for the $m^{th}$ class score is

$$s_m = \mathbf{w}_m \mathbf{A} + b_m, \tag{5}$$

where $\mathbf{w}_m \in \mathbb{R}^{1 \times FD_1 D_2}$.

Then $\frac{\partial s_m}{\partial \mathbf{A}} = \mathbf{w}_m$ and the HiResCAM explanation for class $m$ is

$$\begin{aligned} \tilde{\mathcal{A}}_m^{\text{HiResCAM}} &= \sum_{f=1}^{F} \frac{\partial s_m}{\partial \mathbf{A}} \odot \mathbf{A} \\ &= \sum_{f=1}^{F} \mathbf{w}_m \odot \mathbf{A}, \end{aligned} \tag{6}$$

where here we consider $\mathbf{w}_m$ and $\mathbf{A}$ represented with the following dimensions: $\mathbf{w}_m \in \mathbb{R}^{F \times D_1 \times D_2}$ and $\mathbf{A} \in \mathbb{R}^{F \times D_1 \times D_2}$. In other words, we exploit the fact that $\mathbf{A}$ has spatial information and we can thus infer spatial information for the corresponding elements of $\mathbf{w}_m$.

HiResCAM highlights only relevant locations because $\mathbf{w}_m \odot \mathbf{A}$ in the HiResCAM calculation (where $\mathbf{w}_m \in \mathbb{R}^{F \times D_1 \times D_2}$ and $\mathbf{A} \in \mathbb{R}^{F \times D_1 \times D_2}$) is the intermediate computation in the input-specific part of the class score $\mathbf{w}_m \mathbf{A}$ (where $\mathbf{w}_m \in \mathbb{R}^{1 \times FD_1 D_2}$ and $\mathbf{A} \in \mathbb{R}^{FD_1 D_2 \times 1}$). $\square$

Figure 3 provides a specific example demonstrating how HiResCAM explanations reflect the class score calculation, while Grad-CAM explanations do not.

### 3.5.2 Grad-CAM can highlight irrelevant locations that are not guaranteed to increase the class score

Grad-CAM does not directly visualize important locations, and Grad-CAM does not reflect the model's computations, even if Grad-CAM is applied at the last convolutional layer. Considering the 2D CNN described in the previous section, we begin calculating the Grad-CAM explanation by substituting in $\mathbf{w}_m$ for $\frac{\partial s_m}{\partial \mathbf{A}}$ in the importance weights equation (1):

$$\alpha_m^f = \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \mathbf{w}_m. \tag{7}$$

The final Grad-CAM explanation is then calculated by multiplying these importance weights $\alpha_m^f$ against the corresponding feature maps of $\mathbf{A}$. From equation 7 we see that all the model's fully connected layer weights corresponding to a given feature map are averaged together; averaging of the model's weights prevents the Grad-CAM explanation from reflecting the model's class score calculation.

### 3.5.3 HiResCAM connection to regression

A CNN ending in a single fully connected layer is a feature extractor (convolutions) followed by regression (the fully connected layer). The feature extractor produces a collection of numbers summarizing the input, and this collection of numbers is fed in to the regression to make the final prediction. Regression is fully interpretable. For a simple regression model $y = w_1 x_1 + w_2 x_2$, a *global explanation* consists of inspecting the numeric values of the learned coefficients $w_1$ and $w_2$, while a *local explanation* considers the values of $w_1 x_1$ and $w_2 x_2$ to understand the prediction for a particular input example. The global explanation is equivalent to the model gradient while the local explanation is equivalent to the gradient multiplied element-wise by the input [18]. HiResCAM is thus the local explanation that demonstrates for a particular example the most important locations for the prediction.

### 3.6 CNNs ending in global average pooling then one fully-connected layer: the CAM architecture

In the previous section, we considered CNNs that end in one fully-connected layer. In this section, we consider a narrower group of CNNs: those of the "CAM architecture" which end with global average pooling of feature maps followed by one fully-connected layer. By considering CAM architecture CNNs, we prove that HiResCAM is a generalization of CAM (Section 3.6.2), which then reveals that CAM, HiResCAM, and Grad-CAM applied at the last convolutional layer produce identical explanations for the CAM architecture (Section 3.6.3).

### 3.6.1 Class Activation Mapping (CAM)

Class Activation Mapping [10], or CAM, is a CNN explanation method that requires the CAM architecture. CAM was developed with the observation that for this particular architecture, it is straightforward to visualize the locations in the input image that contribute to an increased score for a particular class.

Consider the following 2D CNN:

1. Convolutional layers: $\text{conv}(\mathbf{X}) = \{\mathbf{A}^f\}_{f=1}^F$ where again here we consider $\{\mathbf{A}^f\}_{f=1}^F$ as the feature maps of specifically the last convolutional layer.

2. Global average pooling of each feature map to a scalar:

$$a^f = \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \mathbf{A}_{d_1 d_2}^f. \tag{8}$$

7

3. Final fully connected layer:

$$s_m = w_m^1 a^1 + w_m^2 a^2 + \cdots + w_m^F a^F \text{ for } m = 1, ..., M \tag{9}$$

where $s_m$ is the score for the $m^{th}$ class and there are $M$ classes total.

The CAM explanation for class $m$ is then defined as:

$$\tilde{\mathcal{A}}_m^{\text{CAM}} = w_m^1 \mathbf{A}^1 + w_m^2 \mathbf{A}^2 + \cdots + w_m^F \mathbf{A}^F, \tag{10}$$

or equivalently:

$$\tilde{\mathcal{A}}_m^{\text{CAM}} = \sum_{f=1}^{F} w_m^f \mathbf{A}^f. \tag{11}$$

### 3.6.2 HiResCAM is a generalization of CAM

In this section we demonstrate that HiResCAM (equation 3) is a generalization of CAM.

*Proof.* As before, to calculate the HiResCAM explanation for class $m$, we first calculate $\frac{\partial s_m}{\partial \mathbf{A}^f}$. To get a useful expression for $\frac{\partial s_m}{\partial \mathbf{A}^f}$ in a CAM architecture, we plug in the right-hand side of the global average pooling equation 8 into the final fully connected layer equation 9:

$$\begin{aligned}
s_m = {} & w_m^1 \left( \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \mathbf{A}_{d_1 d_2}^1 \right) \\
& + w_m^2 \left( \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \mathbf{A}_{d_1 d_2}^2 \right) \\
& + \cdots + w_m^F \left( \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \mathbf{A}_{d_1 d_2}^F \right).
\end{aligned} \tag{12}$$

From the above expression we can see that the gradient $\frac{\partial s_m}{\partial \mathbf{A}^f}$ is[1]

$$\frac{\partial s_m}{\partial \mathbf{A}^f} = \frac{1}{D_1 D_2} w_m^f, \tag{13}$$

which means that the HiResCAM explanation for a model with a CAM architecture is:

$$\tilde{\mathcal{A}}_m^{\text{HiResCAM}} = \frac{1}{D_1 D_2} \sum_{f=1}^{F} w_m^f \mathbf{A}^f. \tag{14}$$

This is identical to the CAM explanation (equation 11) except for a constant factor of $\frac{1}{D_1 D_2}$ which disappears in the subsequent normalization step.

Thus, HiResCAM is a generalization of CAM because both methods yield identical explanations for any CAM architecture, but HiResCAM is applicable to a broader class of architectures as shown previously in Section 3.5. □

---

[1]For another perspective on equation 13, consider that the gradient of the score for class $m$ with respect to the pooling outputs $a^1, a^2, ..., a^F$ are the respective $m$-specific weights $w_m^1, w_m^2, ..., w_m^F$. Each $w_m^f$ value is then propagated back through the global average pooling step by considering how many elements were pooled together, and distributing $w_m^f$ equally across all those elements - that is, dividing $w_m^f$ by $D_1 D_2$, which is the total number of elements in each feature map.

### 3.6.3 CAM, HiResCAM, and Grad-CAM produce identical explanations for CNNs of the CAM architecture

Previous work [11] has demonstrated that Grad-CAM is also a generalization of CAM. Thus, HiResCAM and Grad-CAM are alternative generalizations of the CAM method, and for the CAM architecture the following methods produce identical explanations: CAM, HiResCAM applied at the last convolutional layer, and Grad-CAM applied at the last convolutional layer. Examples of CNNs following the CAM architecture include ResNets [35], GoogLeNet [36], and DenseNets [37].

### 3.7 All other CNNs

Section 3.5 considered CNNs ending in one fully connected layer, while Section 3.6 considered the narrower class of CAM architecture CNNs. For all other CNNs, including those ending in multiple fully connected layers, CAM explanations cannot be calculated, and while HiResCAM and Grad-CAM explanations can be calculated, they are not provably guaranteed to highlight only relevant locations. To the best of our knowledge, no gradient-based neural network explanation method yet exists which can produce consistently class-specific explanations guaranteed to highlight only the locations the model is using, for any arbitrary CNN architecture or layer (see Related Work). For sensitive applications where trustworthy class-specific explanations are required, we recommend using CAM with a CAM architecture CNN, or HiResCAM with any CNN ending in only one fully connected layer.

### 3.8 HiResCAM connection to Gradient ∗ Input

The preceding sections demonstrated HiResCAM's relationship to CAM, Grad-CAM, and regression. HiResCAM is also related to Gradient ∗ Input. If HiResCAM were to be applied at the level of the input image, it would be equivalent to Gradient ∗ Input. While the "level of application" is a simple distinction, it has several important implications. First, HiResCAM explanations are clean and class-specific, whereas Gradient ∗ Input and other pixel-space explanations produce visualizations that are too noisy to be class-specific [11]. Second, HiResCAM is not susceptible to the same issue that caused Gradient ∗ Input to fail sanity checks [21]: namely, HiResCAM does not involve element-wise multiplication with the raw input image. Third, when HiResCAM is used as recommended, the HiResCAM explanations can be seamlessly integrated into model training in a computationally efficient manner, as the gradients necessary to compute the HiResCAM explanation correspond to particular model weights and can thus be accessed during the forward pass without any extra backward passes required.

## 4 Model explanation is not weakly-supervised segmentation

There are close ties between model explanation and weakly-supervised segmentation (WSS), but these tasks have different goals. This section provides context for the Experiments section by assessing the relationship between model explanation and WSS.

### 4.1 The goal of model explanation

The goal of model explanation is to demonstrate what locations in an image a model used to make a particular prediction - *even if that means highlighting areas outside of the relevant object*. For example, if a model has used tracks to identify a train, the explanation should highlight the tracks. If the model has used water to identify a boat, the explanation should highlight the water. Any performance metric to evaluate explanation correctness ("explanation quality") thus must be calculated against a ground truth of "locations the model used for each prediction" which in turn can only be uncovered through mathematical properties of a model and an explanation method - for example, by proving that HiResCAM exactly explains locations used by CNNs ending in one fully connected layer. Locations the model used may not have any relation to object segmentation maps, nor can they be created manually by a human, for if humans were able to understand models well enough to circumscribe regions used for each prediction then there would be no need for explanation methods in the first place.

## 4.2 The goal of weakly-supervised segmentation

The goal of weakly-supervised segmentation is to identify all pixels that are part of an object using only a classification model trained on whole image labels. Because classifiers tend to focus on small discriminative parts of objects [34, 33, 32, 22], a key goal of weakly-supervised segmentation methods is to expand the attention of the classifier beyond small discriminative areas so that the attention covers more of the object, such as by using conditional random fields [32].

Somewhat confusingly, many methods for weakly-supervised segmentation are based on the model explanation methods CAM or Grad-CAM and thus may be named similarly, including Mixup-CAM [31], Sub-Category CAM [32], Puzzle-CAM [38], and FickleNet [34]. However, it is critical to keep in mind that even though these WSS approaches are leveraging model explanations, they have a fundamentally different goal.

## 4.3 IoU must not be used to evaluate explanation correctness

One common performance metric for WSS is intersection-over-union (IoU), which is highest when the predicted segmentation for a class fully overlaps the ground truth segmentation for that class without spreading to other regions. While IoU is a reasonable metric for judging WSS performance, IoU should never be used to evaluate explanation correctness. Unfortunately, some prior work attempts to experimentally evaluate explanation correctness using IoU calculated against ground truth object segmentations, or using the closely related setup of asking humans to subjectively judge how well explanations correspond to an object (which in effect corresponds to humans estimating an IoU-like quantity). The unspoken assumption in these experiments is that the classification model is always using the relevant object to predict the class and so a "good" explanation method will achieve a high IoU. However, as mentioned previously, models are *not* guaranteed to always use the relevant object to predict the class, and indeed, the possibility for undesirable model behavior is a primary motivator behind development of explanation methods. Any time a model behaves unexpectedly or exploits spurious correlations, the IoU of a truthful model explanation will be low, but it would be false to then conclude that the low IoU means the explanation was of poor quality. The only way to know if an explanation method is faithful to a particular type of model is to prove it.

## 4.4 IoU of a faithful explanation method provides insight into a model

Although IoU cannot be used to evaluate the quality of an explanation method, IoU calculated based on an explanation method with guaranteed properties (*e.g.* CAM or HiResCAM) can be used to evaluate a particular *model*. In these cases, high IoU indicates that the model tends to make predictions using areas within the objects of interest, as desired, while low IoU indicates that the model tends to make predictions using areas outside the object of interest and thus is exploiting background or correlated objects.

## 4.5 HiResCAM is an explanation method

The goal of HiResCAM is not to expand the size of the attention maps or to yield high performance on weakly-supervised segmentation. HiResCAM is an explanation method that faithfully represents the locations the model has used to make a prediction, even if these locations are outside the object of interest.

# 5 Experiments

## 5.1 Datasets

To compare HiResCAM and Grad-CAM, in Sections 5.4 through 5.6 we conduct experiments on PASCAL VOC 2012 [39], a data set of 2D RGB natural images with ground truth segmentation maps of 20 classes. Following prior work [27], we combine PASCAL VOC 2012 with SBD [40] to create an augmented training set of 7,087 images. Results are reported on the PASCAL VOC 2012 validation set which contains 1,449 images annotated with 2,148 segmentation maps.

In Section 5.7 we provide a qualitative comparison of HiResCAM and Grad-CAM on the medical imaging dataset RAD-ChestCT, which includes 36,316 chest computed tomography scans annotated with 83 different abnormalities [41].

## 5.2 Models

For the experiments on natural images we report results for a ResNet-34 variant [35] and a DenseNet-121 variant [37] which both end in one fully connected layer and avoid global average pooling, so that HiResCAM and Grad-CAM do not collapse to CAM. These models are defined in the Appendix. Both models use convolutional layers pretrained on ImageNet [42]. All layers are refined using only whole-image labels from PASCAL VOC 2012. PASCAL VOC 2012 experiments were conducted on an NVIDIA Titan RTX GPU with 24 GB of memory. Results on the RAD-ChestCT dataset are from an AxialNet model [43] and were obtained using an NVIDIA Tesla V100 GPU with 32 GB of memory.

## 5.3 Code

All code to replicate all experiments will be made publicly available on GitHub following publication.

## 5.4 Explanation correctness: HiResCAM reflects model computations while Grad-CAM does not

We first quantify the explanation correctness of HiResCAM and Grad-CAM on PASCAL VOC 2012 by calculating the L2 distance between the explanation and the ground truth of the locations the model used. The models are CNNs ending in a single fully connected layer, meaning that the HiResCAM explanation is provably identical to the locations the model used; therefore, the calculated L2 distance between the HiResCAM explanation and the ground truth of the locations used is always 0. However, for Grad-CAM, the explanation is *not* equivalent to the locations used, and so the distance is always nonzero, demonstrating that Grad-CAM provides misleading explanations that do not illustrate the model's behavior (Table 1).

## 5.5 Explanation utility for weakly supervised segmentation

As detailed in Section 4, weakly supervised segmentation (WSS) cannot be used to evaluate an explanation's correctness. However, as CAM explanations are frequently used as a starting point for WSS methods, this section provides WSS performance of raw HiResCAM and Grad-CAM explanations.

To calculate IoU, an explanation with continuous values in the range $[0, 1]$ must be binarized into a proposed segmentation map with values in the set $\{0, 1\}$. Choosing a poor binarization threshold can lead to a low IoU that is not reflective of how well the explanation overlaps the object. To ensure a fair comparison, we therefore selected the best binarization threshold for each explanation method, object class, and model separately, considering thresholds between 0.02 and 0.98 in increments of 0.02. We did not apply any post-processing techniques such as region growing to avoid confounding the results.

The results are summarized in Table 2. Interestingly, Grad-CAM outperforms HiResCAM at WSS. We hypothesize that this may be because Grad-CAM tends to "expand" explanations beyond the regions the model used, an idea that we investigate further in the next section.

Because the HiResCAM explanation reflects the model's computations, we can additionally use the HiResCAM IoU to understand the extent to which these particular models predict each class based on the object itself versus background or other correlated objects. As is apparent from the relatively low IoUs, the models appear to be making some use of background and correlated objects in their predictions.

## 5.6 Humans perceive Grad-CAM and HiResCAM explanations differently

Figure 4 displays HiResCAM and Grad-CAM explanations for 72 randomly-selected image-class pairs. Inspection of this figure suggests that Grad-CAM may be expanding attention maps. In

Table 1: Explanation correctness evaluated by mean L2 distance between the explanation and the ground truth of the locations each model used to make predictions on the PASCAL VOC 2012 validation set. An L2 distance of 0 indicates that the explanation is correct and exactly reflects the locations the model used. A nonzero L2 distance indicates that the explanation is misleading and does not accurately reflect the model's computations.

**L2 Distance to Locations Model Used**

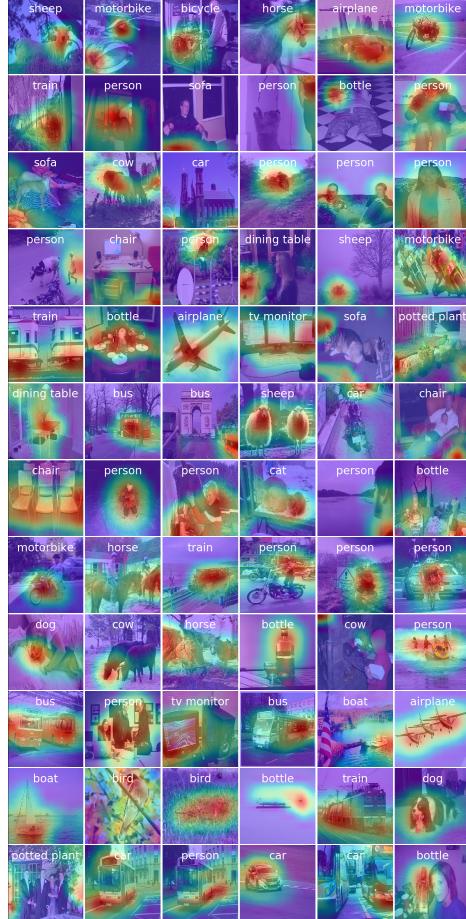| Label | DenseNet-121v | | ResNet-34v | |
| --- | --- | --- | --- | --- |
| | Grad-CAM | HiResCAM | Grad-CAM | HiResCAM |
| airplane | 2.80 | 0 | 3.03 | 0 |
| bicycle | 1.02 | 0 | 0.96 | 0 |
| bird | 2.76 | 0 | 2.69 | 0 |
| boat | 0.99 | 0 | 0.96 | 0 |
| bottle | 0.85 | 0 | 0.81 | 0 |
| bus | 0.91 | 0 | 0.90 | 0 |
| car | 0.72 | 0 | 1.05 | 0 |
| cat | 2.24 | 0 | 2.34 | 0 |
| chair | 0.72 | 0 | 0.74 | 0 |
| cow | 0.88 | 0 | 1.12 | 0 |
| dining table | 0.77 | 0 | 0.85 | 0 |
| dog | 1.47 | 0 | 2.57 | 0 |
| horse | 1.06 | 0 | 1.29 | 0 |
| motorbike | 1.11 | 0 | 1.28 | 0 |
| person | 0.95 | 0 | 1.08 | 0 |
| potted plant | 0.70 | 0 | 0.64 | 0 |
| sheep | 1.20 | 0 | 1.06 | 0 |
| sofa | 0.71 | 0 | 0.63 | 0 |
| train | 1.12 | 0 | 1.56 | 0 |
| tv monitor | 0.75 | 0 | 0.71 | 0 |

order to quantify how human perception of HiResCAM and Grad-CAM explanations differ, we ran Amazon Mechanical Turk (AMT) experiments. We selected 250 image-class pairs randomly from the PASCAL VOC 2012 validation set and plotted the ResNet-34v Grad-CAM explanation next to the HiResCAM explanation such that HiResCAM sometimes randomly appeared on the left as "Image A" and other times randomly appeared on the right as "Image B" (Figure 5). Workers were never informed which explanation method appeared on the right versus the left. Workers were asked three multiple-choice questions to compare the two explanations:

- A size question, with mutually exclusive options "Image A highlight is bigger;" "Image B highlight is bigger;" "Image A and B highlights are the same size."

- A shape question, with mutually exclusive options "Image A highlight is smoother/rounder;" "Image B highlight is smoother/rounder;" "Image A and B highlights have similar shapes"

- A focus question, with mutually exclusive options "Image A highlight is more focused on the *label*;" "Image B highlight is more focused on the *label*;" "Image A and B highlights have similar focus on the *label*" where *label* was replaced with the particular class being explained, e.g. "person" or "tv monitor." If neither Image A nor Image B was focused on the object, the workers were instructed to choose the "similar focus" option.

For each of the three questions for the 250 paired explanations, five unique workers provided an answer, for a total of 3,750 human evaluations. The total number of workers who participated was 42. Worker quality was assessed by calculating how often a worker's answer agreed with the mode(s) of the other 4 answers provided per question.

AMT results are shown in Figure 6. Although we demonstrated previously that Grad-CAM does not reflect the model's computations and therefore should not be used for model explanation, the results of the AMT task intriguingly suggest that Grad-CAM may have extra utility for WSS, as the Grad-CAM explanations are typically bigger, smoother/rounder, and more focused on the object than the HiResCAM explanations.
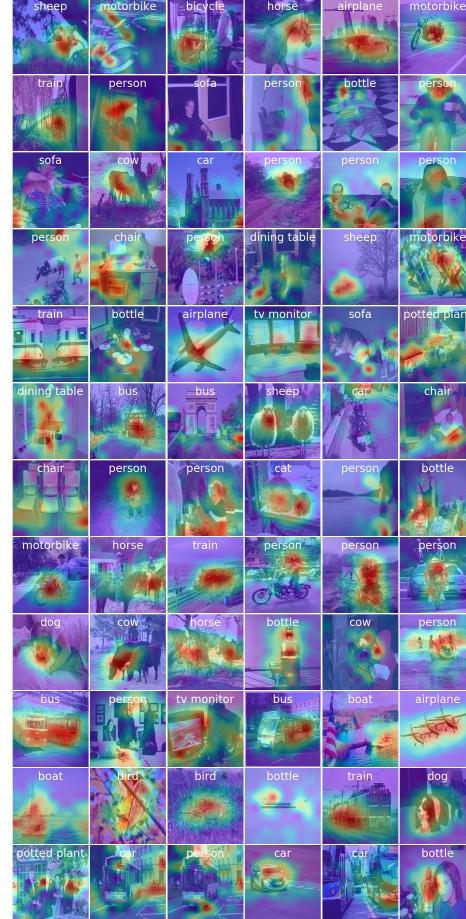
## Grad-CAM                    HiResCAM



Figure 4: HiResCAM and Grad-CAM explanations for 72 randomly-selected PASCAL VOC 2012 validation set image-class pairs for ResNet-34v. Each half of the figure includes the same images and classes. The two halves differ only in the calculation of the explanation. Best viewed in color.

## Image A        Image B



Figure 5: Example of a visualization used in the AMT human evaluation task comparing Grad-CAM and HiResCAM explanations produced using the same model on the same input image and class. HiResCAM appeared as Image A in 50% of the visualizations and as Image B in the remaining 50%. Best viewed in color.

|  | Bigger Size | | Smoother/ Rounder | | More Focused on Label | |
|---|---|---|---|---|---|---|
|  | Grad | HiRes | Grad | HiRes | Grad | HiRes |
| dog | 71 | 14 | 52 | 38 | 53 | 26 |
| horse | 68 | 16 | 44 | 32 | 44 | 32 |
| sofa | 67 | 27 | 70 | 20 | 80 | 13 |
| airplane | 63 | 14 | 46 | 26 | 60 | 14 |
| motorbike | 60 | 20 | 55 | 25 | 56 | 20 |
| train | 58 | 13 | 42 | 33 | 62 | 16 |
| chair | 57 | 37 | 59 | 37 | 70 | 23 |
| bird | 57 | 29 | 50 | 33 | 61 | 26 |
| person | 52 | 36 | 68 | 25 | 65 | 22 |
| cat | 52 | 23 | 50 | 27 | 49 | 25 |
| car | 51 | 36 | 63 | 24 | 52 | 38 |
| tv monitor | 50 | 30 | 80 | 20 | 60 | 20 |
| sheep | 49 | 40 | 54 | 34 | 49 | 40 |
| dining table | 48 | 40 | 52 | 25 | 48 | 37 |
| boat | 46 | 37 | 63 | 29 | 49 | 34 |
| cow | 44 | 44 | 53 | 36 | 49 | 31 |
| bus | 44 | 38 | 47 | 27 | 38 | 47 |
| bottle | 40 | 54 | 60 | 14 | 54 | 23 |
| bicycle | 31 | 42 | 42 | 31 | 40 | 38 |
| potted plant | 22 | 68 | 50 | 40 | 40 | 45 |

Figure 6: Amazon Mechanical Turk human evaluation results comparing ResNet-34v Grad-CAM and HiResCAM explanations on size, shape, and focus. Considering the "Bigger Size" comparison as an example, the "Grad" column indicates the percent of the time that workers judged the Grad-CAM explanation to be bigger in size than HiResCAM, while the "HiRes" column indicates the percent of the time workers judged HiResCAM to be bigger. The percents do not add up to 100 across the Grad and HiRes columns because for each characteristic workers were also allowed to indicate that Grad-CAM and HiResCAM were equivalent. For most classes, workers perceived Grad-CAM explanations as bigger, smoother/rounder, and more focused on the relevant object. Best viewed in color.

Table 2: Explanation utility for weakly supervised segmentation: mean IoU calculated for class-specific explanations versus the corresponding class-specific ground truth segmentation maps provided in the PASCAL VOC 2012 validation set. Note that the primary purpose of HiResCAM is accurate explanation (Table 1), not WSS. The "overall" row at the bottom was calculated by averaging IoU across all image-class pairs, so classes appearing more frequently are weighted proportionally more.

| | Mean IoU | | | |
| | DenseNet-121v | | ResNet-34v | |
| Label | Grad-CAM | HiResCAM | Grad-CAM | HiResCAM |
|---|---|---|---|---|
| airplane | 30.6 | 35.5 | 36.0 | 36.0 |
| bicycle | 13.9 | 14.3 | 15.3 | 15.7 |
| bird | 25.4 | 31.4 | 29.4 | 30.9 |
| boat | 20.8 | 19.8 | 24.1 | 23.4 |
| bottle | 19.4 | 19.0 | 20.2 | 19.0 |
| bus | 40.5 | 41.8 | 48.3 | 45.6 |
| car | 30.1 | 27.8 | 34.1 | 27.8 |
| cat | 39.6 | 39.5 | 42.5 | 41.0 |
| chair | 15.4 | 12.4 | 15.4 | 12.6 |
| cow | 37.8 | 34.9 | 36.9 | 37.5 |
| dining table | 31.5 | 28.4 | 31.5 | 28.4 |
| dog | 36.3 | 38.4 | 38.0 | 38.5 |
| horse | 27.7 | 30.6 | 35.3 | 35.3 |
| motorbike | 36.8 | 37.7 | 42.8 | 41.3 |
| person | 29.9 | 28.4 | 29.4 | 23.5 |
| potted plant | 19.2 | 14.8 | 19.5 | 15.0 |
| sheep | 36.2 | 37.1 | 35.2 | 34.6 |
| sofa | 32.0 | 26.3 | 36.3 | 28.3 |
| train | 32.6 | 34.1 | 41.5 | 43.0 |
| tv monitor | 35.2 | 31.1 | 38.4 | 33.4 |
| overall | **29.4** | 28.9 | **31.8** | 29.1 |

To better understand the manner in which Grad-CAM explanations expand beyond the locations the model used, we generated step-by-step examples showing the intermediate calculations for HiResCAM and Grad-CAM explanations. One such example is shown in Figure 7. This example suggests that Grad-CAM may expand attention maps by over-emphasizing top feature maps.

### 5.7 In medical images, Grad-CAM can create the incorrect impression that the model has focused on the wrong organ

Reliable model explanation is particularly important for sensitive applications, such as those in criminal justice or healthcare. Concurrent work [43] assessed HiResCAM and Grad-CAM for explainable multiple abnormality prediction in volumetric medical images. We leveraged that final AxialNet model to generate additional HiResCAM and Grad-CAM visualizations that reveal how Grad-CAM sometimes creates the false impression that the model focused on the wrong anatomical structure when predicting an abnormality (Figure 8). For a more detailed comparison of HiResCAM and Grad-CAM in medical imaging, see [43].

## 6   Conclusion

In this work, we propose the model explanation method HiResCAM, a new generalization of CAM that is provably guaranteed to highlight locations used by any CNN ending in one fully connected layer. We demonstrate that the related method Grad-CAM can create misleading explanations, but because it also tends to expand attention maps, Grad-CAM yields superior WSS performance. Overall, Grad-CAM's attention expansion properties may be useful for downstream segmentation tasks, while for model understanding HiResCAM provides faithful, class-specific explanations.
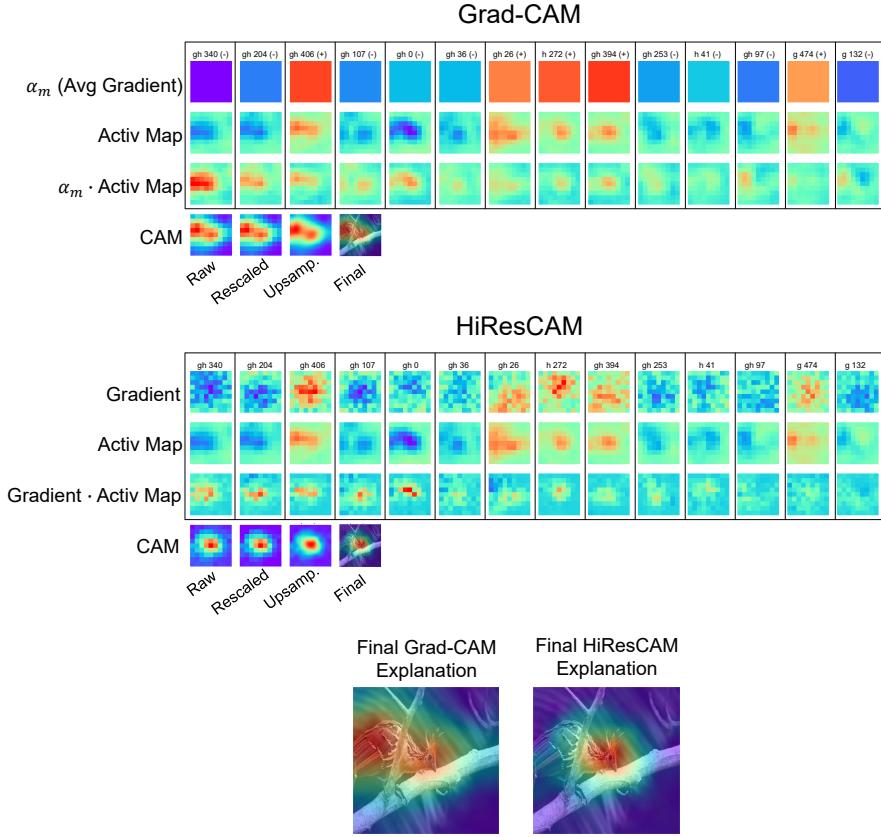
Figure 7: Step-by-step example of how Grad-CAM and HiResCAM explanations are calculated, using a real PASCAL VOC 2012 validation set image with gradients and activation maps from ResNet-34v. The Grad-CAM explanation has expanded the attention beyond the locations the model used by blurring out the gradient information and thus placing more emphasis on certain features that in this case happen to cover more of the object. Figure arrangement: Each boxed column corresponds to one feature dimension. As this model ends in 512 features, for human comprehensibility only the union of the top 12 features that contribute most to each final explanation are shown. Because there was overlap between the top 12 features for HiResCAM and the top 12 features for Grad-CAM for this image, a total of 14 unique features are included. Further explanation of this figure's construction is provided in the Appendix. Best viewed in color, with zoom.



Figure 8: Examples of Grad-CAM creating the incorrect impression that an AxialNet model focused on the wrong anatomical structure. The HiResCAM and Grad-CAM explanations were generated using exactly the same model on the same input CT volume. The only difference is the explanation method. Both of the abnormalities shown here are lung findings, and HiResCAM indicates that the model used the lung fields to predict these lung findings. However, Grad-CAM creates the impression that the model predicted these lung abnormalities based on the body wall and heart, which are irrelevant. Best viewed in color; text annotations added for clarity.

16

## Acknowledgements

## Funding Sources

## References

[1] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Unmasking clever hans predictors and assessing what machines really learn. *Nature communications*, 10(1):1–8, 2019.

[2] John R Zech, Marcus A Badgeley, Manway Liu, Anthony B Costa, Joseph J Titano, and Eric Karl Oermann. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study. *PLoS medicine*, 15(11):e1002683, 2018.

[3] Alex Najibi. Racial discrimination in face recognition technology. *Harvard University*, 2020.

[4] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018.

[5] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *ICLR*, 2014.

[6] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.

[7] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.

[8] Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*, 2016.

[9] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015.

[10] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016.

[11] 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 *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.

[12] Hyunkwang Lee, Sehyo Yune, Mohammad Mansouri, Myeongchan Kim, Shahein H Tajmir, Claude E Guerrier, Sarah A Ebert, Stuart R Pomerantz, Javier M Romero, Shahmir Kamalian, et al. An explainable deep-learning algorithm for the detection of acute intracranial haemorrhage from small datasets. *Nature biomedical engineering*, 3(3):173–182, 2019.

[13] Harsh Panwar, PK Gupta, Mohammad Khubeb Siddiqui, Ruben Morales-Menendez, Prakhar Bhardwaj, and Vaishnavi Singh. A deep learning and grad-cam based color visualization approach for fast detection of covid-19 cases using chest x-ray and ct-scan images. *Chaos, Solitons & Fractals*, 140:110190, 2020.

[14] Ivo M Baltruschat, Hannes Nickisch, Michael Grass, Tobias Knopp, and Axel Saalbach. Comparison of deep learning approaches for multi-label chest x-ray classification. *Scientific reports*, 9(1):1–10, 2019.

[15] Yan Shen and Mingchen Gao. Dynamic routing on deep neural network for thoracic disease classification and sensitive area localization. In *International Workshop on Machine Learning in Medical Imaging*, pages 389–397. Springer, 2018.

[16] F Pasa, V Golkov, F Pfeiffer, D Cremers, and D Pfeiffer. Efficient deep network architectures for fast chest x-ray tuberculosis screening and visualization. *Scientific reports*, 9(1):1–9, 2019.

[17] Michael T Lu, Alexander Ivanov, Thomas Mayrhofer, Ahmed Hosny, Hugo JWL Aerts, and Udo Hoffmann. Deep learning to assess long-term mortality from chest radiographs. *JAMA network open*, 2(7):e197416–e197416, 2019.

[18] Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Gradient-based attribution methods. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pages 169–191. Springer, 2019.

[19] Weili Nie, Yang Zhang, and Ankit Patel. A theoretical explanation for perplexing behaviors of backpropagation-based visualizations. In *International Conference on Machine Learning*, pages 3809–3818. PMLR, 2018.

[20] David Balduzzi, Marcus Frean, Lennox Leary, JP Lewis, Kurt Wan-Duo Ma, and Brian McWilliams. The shattered gradients problem: If resnets are the answer, then what is the question? *arXiv preprint arXiv:1702.08591*, 2017.

[21] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 9505–9515, 2018.

[22] Alexander Kolesnikov and Christoph H Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *European conference on computer vision*, pages 695–711. Springer, 2016.

[23] Xiang Wang, Shaodi You, Xi Li, and Huimin Ma. Weakly-supervised semantic segmentation by iteratively mining common object features. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1354–1362, 2018.

[24] Zilong Huang, Xinggang Wang, Jiasi Wang, Wenyu Liu, and Jingdong Wang. Weakly-supervised semantic segmentation network with deep seeded region growing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7014–7023, 2018.

[25] Yunchao Wei, Huaxin Xiao, Honghui Shi, Zequn Jie, Jiashi Feng, and Thomas S Huang. Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7268–7277, 2018.

[26] Peng-Tao Jiang, Qibin Hou, Yang Cao, Ming-Ming Cheng, Yunchao Wei, and Hong-Kai Xiong. Integral object mining via online attention accumulation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2070–2079, 2019.

[27] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12275–12284, 2020.

[28] Junsong Fan, Zhaoxiang Zhang, Chunfeng Song, and Tieniu Tan. Learning integral objects with intra-class discriminator for weakly-supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4283–4292, 2020.

[29] Jiwoon Ahn and Suha Kwak. Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4981–4990, 2018.

[30] Junsong Fan, Zhaoxiang Zhang, Tieniu Tan, Chunfeng Song, and Jun Xiao. Cian: Cross-image affinity net for weakly supervised semantic segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10762–10769, 2020.

[31] Yu-Ting Chang, Qiaosong Wang, Wei-Chih Hung, Robinson Piramuthu, Yi-Hsuan Tsai, and Ming-Hsuan Yang. Mixup-cam: Weakly-supervised semantic segmentation via uncertainty regularization. *arXiv preprint arXiv:2008.01201*, 2020.

[32] Yu-Ting Chang, Qiaosong Wang, Wei-Chih Hung, Robinson Piramuthu, Yi-Hsuan Tsai, and Ming-Hsuan Yang. Weakly-supervised semantic segmentation via sub-category exploration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8991–9000, 2020.

[33] Kunpeng Li, Ziyan Wu, Kuan-Chuan Peng, Jan Ernst, and Yun Fu. Tell me where to look: Guided attention inference network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9215–9223, 2018.

[34] Jungbeom Lee, Eunji Kim, Sungmin Lee, Jangho Lee, and Sungroh Yoon. Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5267–5276, 2019.

[35] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[36] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.

[37] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.

[38] Sanhyun Jo and In-Jae Yu. Puzzle-cam: Improved localization via matching partial and full features. *arXiv preprint arXiv:2101.11253*, 2021.

[39] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.

[40] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *2011 International Conference on Computer Vision*, pages 991–998. IEEE, 2011.

[41] Rachel Lea Draelos, David Dov, Maciej A. Mazurowski, Joseph Y. Lo, Ricardo Henao, Geoffrey D. Rubin, and Lawrence Carin. Machine-learning-based multiple abnormality prediction with large-scale chest computed tomography volumes. *Medical Image Analysis*, 67, 2021.

[42] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[43] Rachel Lea Draelos and Lawrence Carin. Explainable multiple abnormality classification of chest ct volumes with axialnet and hirescam. *arXiv preprint*, 2021.

# 7 Appendix

## 7.1 ResNet-34v and DenseNet-121v Models

In the main paper we report results on PASCAL VOC 2012 for two models, ResNet-34v and DenseNet-121v. The ResNet-34 variant (ResNet-34v) begins with all the standard ResNet-34 convolutional layers pretrained on ImageNet. The DenseNet-121 variant (DenseNet-121v) begins with all the standard DenseNet-121 convolutional layers pretrained on ImageNet. Both models end with a randomly-initialized custom convolutional layer (512 output feature maps, 2×2 kernel, 1×1 stride, no padding) followed by a randomly-initialized fully connected layer that produces the final predictions. All explanations were calculated at the last convolutional layer. All layers were refined during training on the PASCAL VOC 2012 classification task. The models do not include the global average pooling of the CAM architecture because otherwise HiResCAM and Grad-CAM both reduce to CAM and no meaningful comparison between the methods can be made.

## 7.2 Details of the step-by-step figure

Figure 7 shows a step-by-step example of how Grad-CAM and HiResCAM explanations were calculated for a PASCAL VOC 2012 validation set image and a ResNet-34v model. This section provides additional details explaining the figure.

The figure illustrates only the top features for Grad-CAM and HiResCAM. For Grad-CAM, top features are defined as those for which the average value of ($\alpha_m \times$ activation map) is highest. For HiResCAM, top features are those for which the average value of (gradient $\times$ activation map) is highest.

In the figure, a feature column is titled with a "g" if it was one of Grad-CAM's top features, an "h" if it was one of HiResCAM's top features, and a "gh" if it was a top feature for both explanation methods. For example, the leftmost feature in Figure 7 is "gh 340" meaning this is feature 340/512 and was a top contributor to both Grad-CAM and HiResCAM for this image and class. For Grad-CAM, the "(-)" part of "gh 340 (-)" indicates that the $\alpha_m$ value was negative for this feature.

The CAM row at the bottom shows the raw CAM first (the sum over the feature dimension of $\alpha_m \times$ activation map or gradient $\times$ activation map; all 512 features, even those not shown in the figure, are included in the CAM), then the CAM rescaled to [0,1] which has the same relative colors when visualized, then the CAM upsampled to the dimensions of the input image, and last the final explanation which is the upsampled CAM overlaid on the input image itself.

## 7.3 The bias term of the final fully connected layer has no effect on CAM visualizations

This section includes a miscellaneous observation about CAM explanations that was noted during the course of writing this paper. The authors of the CAM paper state, "we ignore the bias term: we explicitly set the input bias of the softmax to 0 as it has little to no impact on the classification performance." We found that there is no need to explicitly ignore the bias term because it drops out of the visualization on its own. Section 3.6.2 demonstrated that HiResCAM is a generalization of CAM. Calculating the CAM map from the HiResCAM perspective illustrates why the bias term is irrelevant to the explanation in a CAM architecture. Consider the score $s_m$ written in terms of the feature maps, for a final fully connected layer that does have a nonzero bias term:

$$
\begin{aligned}
s_m = {} & w_m^1 \left( \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \mathbf{A}_{d_1 d_2}^1 \right) \\
& + w_m^2 \left( \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \mathbf{A}_{d_1 d_2}^2 \right) \\
& + \cdots + w_m^F \left( \frac{1}{D_1 D_2} \sum_{d_1=1}^{D_1} \sum_{d_2=1}^{D_2} \mathbf{A}_{d_1 d_2}^F \right) + b_m.
\end{aligned}
\tag{15}
$$

When calculating the gradient $\frac{\partial s_m}{\partial \mathbf{A}}$, the bias term disappears, because it has nothing to do with $\mathbf{A}$. Thus, the formula for CAM (or HiResCAM) remains the same, with or without a bias term.

There is some intuition about why the bias term of the final fully connected layer should not affect the visualizations: the point of CAM is to show an explanation for a particular class and a particular input image. However, the bias term will be the same for every image in an arbitrary collection of images, because the bias term is part of the model and is not input-dependent. Therefore, the bias term must not contribute anything input-specific to the explanation.