

Keep CALM and Improve Visual Feature Attribution

Jae Myung Kim^{1*} Junsuk Choe^{2*} Zeynep Akata^{1,3,4} Seong Joon Oh^{5†}

¹University of Tübingen ²Department of Computer Science and Engineering, Sogang University

³Max Planck Institute for Intelligent Systems ⁴Max Planck Institute for Informatics ⁵NAVER AI Lab

Abstract

The class activation mapping, or CAM, has been the cornerstone of feature attribution methods for multiple vision tasks. Its simplicity and effectiveness have led to wide applications in the explanation of visual predictions and weakly-supervised localization tasks. However, CAM has its own shortcomings. The computation of attribution maps relies on ad-hoc calibration steps that are not part of the training computational graph, making it difficult for us to understand the real meaning of the attribution values. In this paper, we improve CAM by explicitly incorporating a latent variable encoding the location of the cue for recognition in the formulation, thereby subsuming the attribution map into the training computational graph. The resulting model, **class activation latent mapping**, or **CALM**, is trained with the expectation-maximization algorithm. Our experiments show that CALM identifies discriminative attributes for image classifiers more accurately than CAM and other visual attribution baselines. CALM also shows performance improvements over prior arts on the weakly-supervised object localization benchmarks. Our code is available at <https://github.com/naver-ai/calm>.

1. Introduction

Interpretable AI [25, 40, 24, 52] is becoming an absolute necessity in safety-critical and high-stakes applications of machine learning. Along with good recognition and prediction accuracies, we require models to be able to transparently communicate the inner mechanisms with human users. In visual recognition tasks, researchers have developed various feature attribution methods to inspect contributions of individual pixels or visual features towards the final model prediction. Input gradients [55, 56, 58, 5, 63, 39, 28] and input perturbation methods [59, 70, 20, 43, 23, 47] have been actively researched.

In this paper, we focus on the **class activation mapping (CAM)** [68] method, which has been the cornerstone of

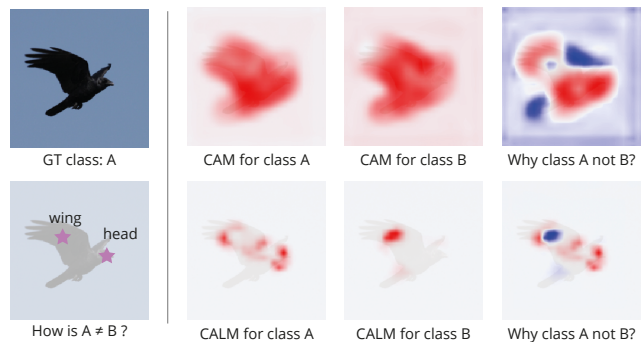


Figure 1. **CAM vs CALM**. CALM is better at locating the actual cues used for the recognition than CAM. Two bird classes A and B only differ in their head and wing attributes. Attributions for class A, B, and their difference are shown. While CAM fails to detect the head and wing, CALM captures them accurately.

the feature attribution research. CAM starts from the observation that many CNN classifiers make predictions by aggregating location-wise signals. For example, $p(y|x) = \text{softmax} \left(\frac{1}{HW} \sum_{hw} f_{yhw} \right)$ where $f = f(x)$ is the extracted feature map in $\mathbb{R}^{C \times H \times W}$ where C, H, W are the number of classes, height, and width of the feature map, respectively. CAM considers the pre-GAP feature map f_{yhw} as the attribution, after scaling it to the $[0, 1]$ range by dropping the negative values and dividing through by the maximum value: $s := (\max_{hw} f_{hw})^{-1} f^+ \in [0, 1]^{H \times W}$. Thanks to the algorithmic simplicity and reasonable effectiveness, CAM has been a popular choice as an attribution method with many follow-up variants [53, 9, 6, 69, 61, 45, 22].

Despite its popularity and contributions to the interpretability community, CAM still has its own limitations. What does the attribution map s really mean? We fail to find a reasonable linguistic description because s hardly encodes anything essential in the recognition process. s also violates key minimal requirements, or “axioms” [40, 59, 22], for an attribution method. For example, its dependence on the pre-softmax values f make it ill-defined: translating $f \mapsto f + c$ yields an identical model because of the translation invariance of softmax, but it changes the attribution map s .

We thus introduce a novel attribution method, **class ac-**

*Equal contribution. Majority of work done at NAVER AI Lab.

†Corresponding author.

tivation latent mapping (CALM). It builds a probabilistic graphical model on the last layers of CNNs, involving three variables: input image X , class label $Y \in \{1, \dots, C\}$, and the location of the cue for recognition $Z \in \{1, \dots, HW\}$. Since there is no observation for Z , we consider latent-variable training algorithms like marginal likelihood (ML) and expectation-maximization (EM). After learning the dependencies, we define the attribution map for image \hat{x} of class \hat{y} as $p(\hat{y}, z|\hat{x}) \in [0, 1]^{H \times W}$, the joint likelihood of the recognition cue being at z and the class being \hat{y} . CALM has many advantages over CAM. (1) It has a human-understandable probabilistic definition; (2) it satisfies the axiomatic requirements for attribution methods; (3) it is empirically more accurate and useful than CAM.

In our experimental analysis, we study how well CALM localizes the “correct cues” for recognizing the given class of interest. The “correct cues” for recognition are ill-defined in general, making the evaluation of attribution methods difficult. We build a novel evaluation benchmark on pairs of bird classes in CUB-200-2011 [62] where the true cue locations are given by the parts where the attributes for the class pair differ (Figure 1). Under this benchmark and a widely-used *remove-and-classify* type of benchmark, CALM shows better attribution performances than CAM and other baselines. We also show that CALM advances the state of the art in the weakly-supervised object localization (WSOL) task, where CAM has previously been one of the best [13, 12].

In summary, our paper contributes (1) analysis on the lack of interpretability for CAM, (2) a new attribution method CALM that is more interpretable and communicable than CAM, and (3) experimental results on real-world datasets where CALM outperforms CAM in multiple tasks. Our code is available at <https://github.com/naver-ai/calm>.

2. Related Work

Interpretable AI is a big field. The general aim is to enhance the transparency and trustworthiness of AI systems, but different sub-fields are concerned with different parts of the system and application domains. In this paper, we develop a visual feature attribution method for image classifiers based on deep neural networks. It is the task of answering the question: “how much does each pixel or visual feature contribute towards the model prediction?”

Gradient-based attribution. Feature attribution with gradients dates back to the pioneering works by Sung [60] and Baehrend *et al.* [7]. The first explicit application to CNNs is the work by Simonyan *et al.* [55]. Input gradients consider local linearization of the model, but it is often not suitable for CNNs because the local behavior hardly encodes the complex mechanisms in CNNs for *e.g.* more global perturbations on the input. Follow-up works have customized the backpropagation algorithm to improve the attribution performances: Guided Backprop [57], LRP [5], Deep Taylor

Decomposition [39], SmoothGrad [56], Full-Gradient [58], and others [65, 30, 54, 63, 3, 28]. We make an empirical comparison against key prior methods in this domain.

Perturbation-based attribution. Researchers have developed methods for measuring the model response to non-local perturbations. Integrated Gradients [59] measure the path integral of model responses to global input shifts. Another set of methods consider model responses to redacted input parts: sliding windows of an occlusion mask [70] and random-pixel occlusion masks [43]. Since the occluding patterns introduce artefacts that may mislead attributions, different options for redaction have been considered: “meaningful perturbations” like image blurring [20, 19], inpainting [70], and cutting-and-pasting a crop from another image of a different class [23]. Some of the key methods above are included as baselines for our experiments.

CAM-based attribution. Gradients and perturbations analyze the model by establishing the input-output relationships. Class activation mapping (CAM) [68] takes a different approach. Many CNNs have a global average pooling (GAP) layer towards the end. CAM argues that the pre-GAP features represent the discriminativeness in the image. Related works have considered variants of the last-layer modifications like max pooling [41] and various thresholding strategies [17, 16, 6]. GradCAM [53] and GradCAM++ [9] have later expanded the usability of CAM to networks of any last-layer modules by combining the widely-applicable gradient method with CAM. In this work, we identify issues with CAM and suggest an improvement.

Self-explainable models. Above attribution methods provide interpretations of a complex, black-box model in a post-hoc manner. Another paradigm is to design models that are interpretable by design in the first place [18]. There is a trade-off between interpretability and performance [34]; researchers have sought ways to push the boundary on both fronts. One line of work *distills* the complex, performant model into an interpretable surrogate model such as decision trees [21], sparse linear models [44, 46, 4]. Other works pursue a *hybrid* approach, where a small interpretable module of a neural network is exposed to humans, while the complex, less interpretable modelling is performed in the rest of the network. ProtoPNet [10] trains an interpretable linear map over prototype neural activations. Concept or semantic bottleneck models [36, 31] enforces an intermediate layer to explicitly encode semantic concepts. Our work is a *hybrid* self-explainable model based on the interpretable probabilistic treatment of the last layers of CNNs.

Evaluating attribution is challenging because of the lack of ground truths. Early works have resorted to qualitative [23] or human-in-the-loop evaluations [46, 53, 48, 33] with limited reproducibility. Wojciech *et al.* [51] and subsequent works [30, 43, 46, 26] have proposed a quantita-

tive measure based on the *remove-and-classify* framework. Along a different axis, researchers have focused on the *necessary conditions* for attribution methods. They can be either theoretical properties, referred to as “axioms” [59, 22] or empirical properties, referred to as “sanity checks” [2]. In our work, we analyze CALM in terms of the axioms (§4.2) and evaluate it on a remove-and-classify benchmark (§5.3). Additionally, we contribute a new type of evaluation; we compare the attribution map against the *known* ground-truth attributions on a real-world dataset [62] (§5.2).

Weakly-supervised object localization (WSOL) is similar yet different from the feature attribution task. While the latter is focused on detecting the small, class-discriminative cues in the input, the former necessitates the detection of the full object extents. Despite the discrepancy in the objective, CAM has been widely used for both tasks without modification [13, 12]. We show that CALM, despite being proposed for the attribution task, outperforms CAM on the WSOL task after some additional aggregation operations (§5.4).

3. Class Activation Mapping (CAM)

We cover the background for the class activation mapping (CAM) [68] and analyze its problems. CAM is a feature attribution method for CNN image classifiers. It is applicable to CNNs with the following last layers:

$$p(y|x) = \text{softmax} \left(\frac{1}{HW} \sum_{hw} f_{yhw}(x) \right) \quad (1)$$

where $f(x)$ is feature map from a fully-convolutional network [35] with dimensionality $C \times H \times W$; each channel corresponds to a class-wise feature map. The network is trained with the negative log-likelihood (NLL), also known as cross-entropy, loss.

At test time, the attribution map is computed by first fetching the pre-GAP feature map $f^{y=\hat{y}}(x) \in \mathbb{R}^{H \times W}$ for the ground-truth class \hat{y} . CAM then normalizes the feature map $f^{\hat{y}}$ to the interval $[0, 1]$ in either ways:

$$s = \begin{cases} (f_{\max}^{\hat{y}})^{-1} \max(0, f^{\hat{y}}) & \text{max [68]} \\ (f_{\max}^{\hat{y}} - f_{\min}^{\hat{y}})^{-1} (f^{\hat{y}} - f_{\min}^{\hat{y}}) & \text{min-max [53]} \end{cases} \quad (2)$$

where $f_{\{\min, \max\}} := \{\min_{hw}, \max_{hw}\} f_{hw}$.

Note that the original CAM paper [68] considers CNNs with an additional linear layer $W \in \mathbb{R}^{C \times L}$ after the pooling (e.g. ResNet). It is known that such networks are equivalent to Equation 1 when we swap the linear and the GAP layers (which are commutative) and treat the linear layer as a convolutional layer with 1×1 kernels [53].

3.1. Limitations of CAM

CAM lacks interpretability. How can we succinctly communicate the attribution value s_{hw} at pixel index (h, w) to others? The best we can come up with is:

“The pixel-wise pre-GAP, pre-softmax feature value at (h, w) , measured in relative scale within the range of values $[0, A]$ where A is the maximum of the feature values in the entire image.”

This description is hardly communicable even to experts in image recognition systems, not to mention general users. The difficulty of communication stems from the fact that the attribution scores s_{hw} are not the quantities used by the recognition system; the computational graph for CAM (Equation 2) is not part of the training graph (Equation 1).

We present the issues with CAM according to the set of axiomatic criteria for attribution methods [40, 59, 22].

Implementation-invariance axiom [59] states that two mathematically identical functions, $\phi_1 \equiv \phi_2$, shall possess the same attribution maps, regardless of their implementations. CAM violates this axiom. Assume $\phi_1(f) := \text{softmax}(\frac{1}{HW} \sum_{hw} f_{yhw})$ and $\phi_2(f) := \text{softmax}(\frac{1}{HW} \sum_{hw} f_{yhw} + C)$ for some constant C . Since the softmax function is translation invariant, $\phi_1 \equiv \phi_2$ for any C . However, the CAM attribution map for ϕ_2 varies arbitrarily with C : $s = (\max_{hw} f_{hw} + C)^{-1} (f + C)^+$. Min-max normalization is a solution to the problem, but it alone does not let CAM meet other axioms. This observation reveals the inherent limitation of utilizing feature values before softmax (often called “logits”) as attribution.

Sensitivity axiom [59, 22] states that if the function response $\phi(x)$ changes as the result of altering an input value x_{hw} at (h, w) , then the corresponding attribution value s_{hw} shall be non-zero. Conversely, if the response is not affected, then s_{hw} shall be zero. CAM fails to satisfy the sensitivity axiom. Depending on the normalization type, CAM assigns zero attributions to (h, w) where f_{hw} is either negative (for max normalization) or smallest (for min-max normalization). However, being assigned a negative or smallest feature value f_{hw} has little connection to the insensitivity of the model to the input value x_{hw} .

Completeness (or conservation) axiom [59, 22, 54, 5] states that the sum of attributions $\sum_{hw} s_{hw}$ shall add up to the function output $\phi(x) = p(y|x)$. The completeness criterion is violated by CAM in general because the summation $\sum_{hw} s_{hw}$ for s in Equation 2 do not match $\phi(x) = p(y|x)$ in Equation 1. In conclusion, CAM fails to satisfy key minimal requirements for an attribution method.

4. Class Activation Latent Mapping (CALM)

We fix the above issues by introducing a probabilistic learning framework involving **input image** X , **class label** Y , and the **cue location** Z . We set up a probabilistic graphical model and discuss how each component is parametrized with a CNN. We then introduce learning algorithms to account for the unobserved latent variable Z . An overview of our method is provided in Figure 2.

Probabilistic framework for CALM: inference and learning.

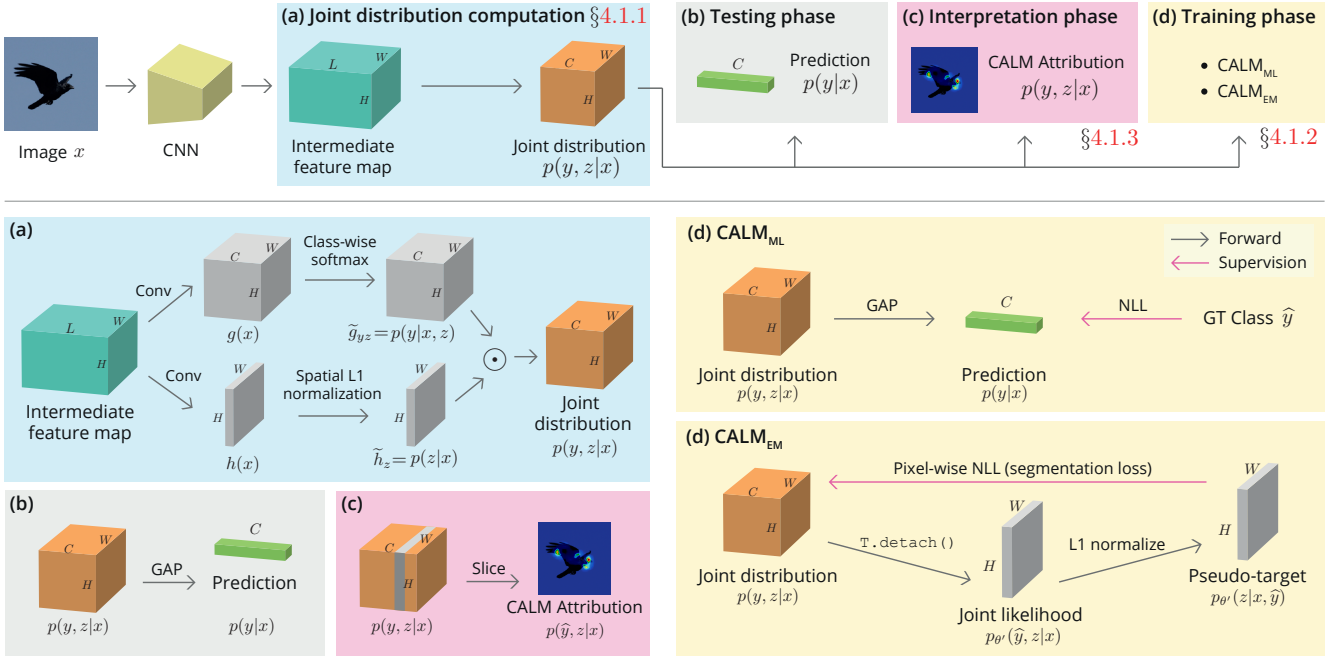
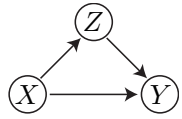


Figure 2. **Main components of CALM.** We show the computational pipeline for CALM during testing, interpretation, and training phases. We zoom into different components. See the relevant sections for more details.

4.1. Probabilistic inference with latent Z

We define Z as the location index (h, w) of the cue for recognizing the image X as the class Y . Our aim is to let the model explicitly depend its prediction on the features corresponding to the location Z and later on use the distribution of possible cue locations Z as the attribution provided by the model. Z is a random variable over indices (h, w) ; for simplicity, we use the integer indices $Z \in \{1, \dots, HW\}$.

Without loss of generality, we factorize $p(x, y, z)$ as $p(y, z|x)p(x) = p(y|x, z)p(z|x)p(x)$ (graph on the left). The recognition task is then performed via $p(y|x) = \sum_z p(y, z|x)$.



4.1.1 Representing joint distribution with CNNs.

We factorize the joint distribution $p(y, z|x)$ into $p(y, z|x) = p(y|x, z)p(z|x)$ and parametrize $p(y|x, z)$ and $p(z|x)$ as two convolutional branches of a CNN trunk (Figure 2a). Since $Y \in \{1, \dots, C\}$ and $Z \in \{1, \dots, HW\}$, we represent $p(y|x, z)$ as a CNN branch $g(x)$ with output dimensionality $C \times HW$. Likewise, we represent $p(z|x)$ with a CNN branch $h(x)$ with output dimensionality HW . To make sure that the outputs of the two branches are proper distributions, we normalize the outputs with softmax for g and ℓ_1 normalization followed by the softplus for h . We broadcast h to all class indices Y and multiply it element-wise with g to get $p(y, z|x)$ (Figure 2a).

4.1.2 Training algorithms

Training a latent variable model is challenging because of the unobserved variable Z . We consider two methods for training such a model: (1) marginal likelihood (ML) [37] and (2) expectation-maximization (EM) [15].

CALM_{ML} directly minimizes the marginal likelihood

$$-\log p_\theta(y|x) = -\log \sum_z p_\theta(y|x, z)p_\theta(z|x) \quad (3)$$

$$= -\log \sum_z g_{yz} \cdot h_z. \quad (4)$$

which is tractable for the discrete Z . See Figure 2d.

CALM_{EM} is based on the EM algorithm that generates pseudo-targets for Z to supervise the joint likelihood $p(y, z|x)$. The EM algorithm introduces two running copies of the parameter set: θ and θ' . The first signifies the model of interest, while the latter often refers to a slowly updated parameter used for generating the pseudo-targets for Z . The learning objective is

$$-\log p_\theta(y|x) \leq -\sum_z p_{\theta'}(z|x, y) \log p_\theta(y, z|x) \quad (5)$$

$$= -\sum_z \frac{g'_{yz} \cdot h'_z}{\sum_l g'_{yl} \cdot h'_l} \log (g_{yz} \cdot h_z) \quad (6)$$

where g' and h' denote the parametrization with θ' . Note that Equation 6 is the pixel-wise negative log likelihood, the loss function for semantic segmentation networks [11]. One may interpret the objective as self-supervising the pixel(z)-wise predictions $p(y, z|x)$ with its own estimation of the cue location z for the true class y : $p_{\theta'}(z|x, y)$. In practice, we use the current-iteration model parameter $\theta = \theta'$ to generate the pseudo-target for Z . See Figure 2d for an overview of the process. Even with $\theta = \theta'$, we need to apply `T.detach()` to block the gradient flow through the pseudo-target $p_{\theta'}(z|x, y)$, as required by Equation 6.

A similar framework appears in the weakly-supervised semantic segmentation task. Papandreou *et al.* [42] have generated pseudo-target label maps to train a segmentation network. CALM_{EM} is different because our location-encoding latent Z takes integer values, while their Z takes values in the space of all binary masks; our formulation admits an exact computation of Equation 6, while theirs require an additional approximation step.

4.1.3 Inferring feature attributions

Unlike CAM, our probabilistic formulation enables principled computation of the attribution map as part of the probabilistic inference on $p(y, z|x)$. Z is explicitly defined as the location of the cue for recognition. For CALM, the **attribution score** s_z for location z is naturally defined as the joint likelihood given the ground-truth class \hat{y}

$$s_z := p(\hat{y}, z|x), \tag{7}$$

or in human language,

“The probability that the cue for recognition was at z and the ground truth class \hat{y} was corrected predicted for the image x .”

Note that the definition is far more communicable than the one for CAM in §3.1. See Figure 2c for visualization.

Apart from the attribution map, one may compute additional interesting quantities. We show examples in Figure 3. Treating z as a free variable, the **conditional attribution** $p(y|x, z)$ is explained as the likelihood of the cue being at position z , given the prediction for image x as y . The **saliency** $p(z|x)$ encodes the likely location of any cue for recognizing classes $y \in \{1, \dots, C\}$ in image x . It is the sum over all attribution maps for classes y : $p(z|x) = \sum_y p(y, z|x)$. One may also compute the partial sum for classes $y \in \mathcal{Y}$ to obtain the **subset attribution** to highlight specific image regions of interest $p(z, \mathcal{Y}|x) := \sum_{y \in \mathcal{Y}} p(z, y|x)$. Above quantities are later utilized for the weakly-supervised object localization (WSOL) task in §5.4. It is also possible to reason why the class label for input x is \hat{y} instead of y' by computing the **counterfactual attribution** $p(\hat{y}, z|x) - p(y', z|x)$. Such counterfactual reasoning will be used in our analysis in §5.2.

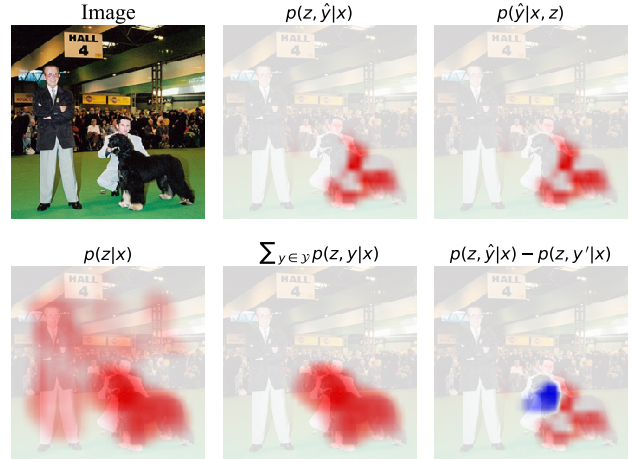


Figure 3. Various attribution maps by CALM on ImageNet. GT class is “Afghan hound”. For the subset attribution, the classes \mathcal{Y} correspond to all species of dogs in ImageNet. For the counterfactual attribution, the alternative class y' is “Gazelle hound”.

4.2. Theoretical properties of CALM

Now we revisit the axioms for attribution methods that CAM fails to fulfill (§3.1). **Implementation-invariance axiom** is satisfied by CALM because the attribution map $s := p(\hat{y}, z|x)$ is a mathematical object in the probabilistic graphical model. CALM attribution also does not depend on the fragile logit values. The **completeness axiom** trivially follows from CALM because the final prediction $p(y|x)$ is the sum of attribution values $p(y, z|x)$ over z . Likewise, the **sensitivity axiom** follows trivially from the fact that $p(y, z|x) > 0$ if and only if it contributes towards the sum $p(y|x) = \sum_z p(y, z|x)$.

The superior interpretability of CALM comes with a cost to pay. It alters the formulation of the usual structure for CNN classifiers where the loss function has the structure “NLL \circ SoftMax \circ Pool” on the feature map f into the one with the structure “Pool \circ NLL \circ SoftMax” on f . Compared to the former, CALM gain additional interpretability by making the last layer of the network as simple as a sum over the pixel-wise experts $p(y, z|x)$. The reduced complexity in turn increases the representational burden for the earlier layers $f(x)$ and induces a drop in the classification accuracies (§5.3).

The interpretability-performance trade-off is unavoidable [34]. Therefore, it benefits users to provide a diverse array of models with different degrees of interpretability and performance [49]. Our work contributes to this diversity of the ecosystem of models.

5. Experiments

We present experimental results for CALM. We present two experimental analyses on attribution qualities: evalua-

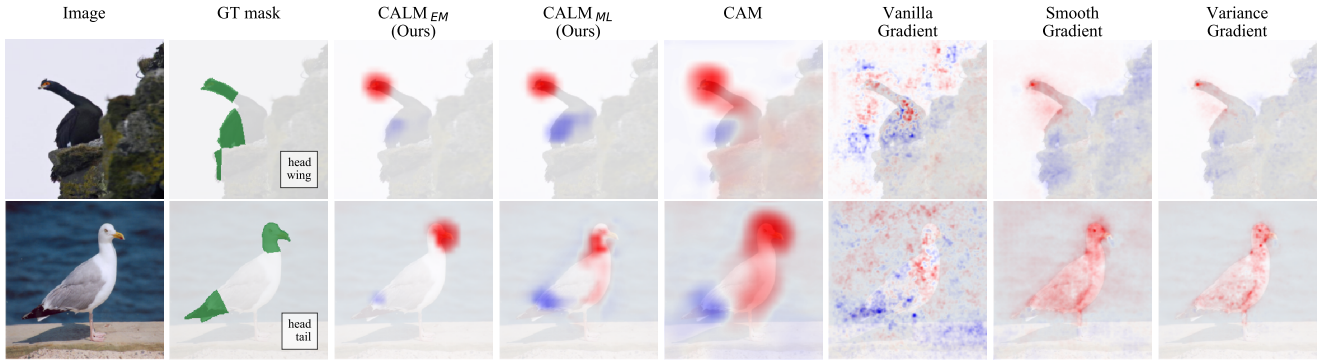


Figure 4. **Qualitative results on CUB.** We compare the counterfactual attributions from CALM and baseline methods against the GT attribution mask. The GT mask indicates the bird parts where the attributes for the class pair (A,B) differ. The counterfactual attributions denote the difference between the maps for classes A and B: $s^A - s^B$. Red: positive values. Blue: negative values.

tion with respect to estimated ground truth attributions on CUB (§5.2) and the remove-and-classify results on three image classification datasets (§5.3). We then show results on the weakly-supervised object localization task in §5.4. Naver Smart Machine Learning (NSML) platform [29] has been used in the experiments.

5.1. Implementation details and baselines

Backbone. CALM is backbone-agnostic, as long as it is fully convolutional. We use the ResNet50 as the feature extractor f unless stated otherwise. As discussed in §3, we move the final linear layer before the global average pooling layer as a convolutional layer with 1×1 kernels.

Datasets. Our experiments are built on three real-world image classification datasets: CUB-200-2011 [62], a subset of OpenImages [8], and ImageNet1K [50]. CUB is a fine-grained bird classification dataset with 200 bird classes. We use a subset of OpenImages curated by [12] that consists of 100 coarse-grained everyday objects. ImageNet1K has 1000 classes with mixed granularity, ranging from 116 fine-grained dog species to coarse-grained objects and concepts.

Pretraining. We use the ImageNet pre-trained weights for f . The two convolutional layers for computing $p(y, z|x)$ in Figure 2a are trained from scratch.

Attribution maps. The attribution maps for CAM and CALM are scaled up to the original image size via bilinear interpolation. For gradient-based baseline attribution methods, we apply Gaussian blurring and min-max normalization, following [12].

Other training details are in the Supplementary Materials.

5.2. Cue localization results

The difficulty of attribution evaluation comes from the fact that it is difficult to obtain the ground truth cue locations \hat{z} . We propose a way to estimate the true cue location using the rich attribute and part annotations on the bird images in CUB-200-2011 [62].

#part differences	1	2	3	
#class pairs	31	64	96	mean
Vanilla Gradient [55]	10.0	13.7	15.3	13.9
Integrated Gradient [59]	12.0	15.1	17.3	15.7
Smooth Gradient [56]	11.8	15.5	18.6	16.5
Variance Gradient [3]	16.7	21.1	23.1	21.4
CAM [68]	24.1	28.3	32.2	29.6
CALM _{ML} (Ours)	23.6	26.7	28.8	27.3
CALM _{EM} (Ours)	30.4	33.3	36.3	34.3

Table 1. **Attribution evaluation on CUB.** We use the estimated GT attribution masks (§5.2) to measure the performances of attribution methods. Mean pixel-wise average precision (mPxAP) values are reported. See Figure 4 for the setup and examples.

Estimating GT cue locations. We generate the ground-truth cue locations using following intuition: for two classes A and B differing only in one attribute a , the location z for the cue for predicting A instead of B will correspond to the object part containing the attribute a . We explain algorithmically how we build the ground-truth attribution mask for an image x with respect to two bird classes A and B in CUB. We first use the attribute annotations for 312 attributes in CUB to compute the set of attributes for each class: \mathcal{S}^A and \mathcal{S}^B . For example, $\mathcal{S}^{\text{Fish crow}} = \{\text{black crown, black wing, all-purpose bill-shape, } \dots\}$. We then compute the symmetric difference of the attributes for the two classes $\mathcal{S}^A \Delta \mathcal{S}^B = (\mathcal{S}^A \cup \mathcal{S}^B) \setminus (\mathcal{S}^A \cap \mathcal{S}^B)$. Now, we map each attribute in $a \in \mathcal{S}^A \Delta \mathcal{S}^B$ to the corresponding bird part $p \in \mathcal{P}$ among 7 bird parts annotated in CUB. For example, the attribute-mismatching bird parts for classes “Fish crow” and “Brandt cormorant” are $\mathcal{P}^{A,B} = \{\text{head, wing}\}$. We locate the parts $\mathcal{P}^{A,B}$ in samples x of classes A and B using the keypoint annotations in CUB: $\mathcal{K}^{A,B}(x)$. We expand the keypoint annotations to a binary mask $\mathcal{M}^{A,B}(x) \in \{0, 1\}^{H \times W}$ using the nearest-neighbor assignment of pixels to bird parts. The final mask $\mathcal{M}^{A,B}(x)$ for the input x is used as the ground-truth attribution map. See the “GT mask” column in Figure 4 for example binary

masks. For evaluation we use all class pairs in CUB with the number of attribute-differing parts $|\mathcal{P}^{A,B}| \leq 3$, resulting in $31 + 64 + 96 = 191$ class pairs.

Counterfactual attributions. To predict the difference in needed cues for recognizing classes A and B, we obtain the absolute values of counterfactual attributions from each method by computing the difference $|s^A - s^B| \in [0, 1]^{H \times W}$. The underlying assumption is that s^A and s^B point to cues corresponding to the attributes for A and B, respectively. Hence, by taking the difference, one removes the attributions on regions that are important for both A and B.

Evaluation metric: mean pixel-wise AP. To measure how well attribution maps retrieve the ground-truth part pixels $\mathcal{M}^{A,B}(x)$, we measure the average precision for the pixel retrieval task [1, 12]. Given a threshold $\tau \in [0, 1]$, we define the positive predictions as the set of pixels in over multiple images: $\{(n, h, w) \mid |s_{hw}^A(x_n) - s_{hw}^B(x_n)| \geq \tau\}$ for images x_n from classes A and B. With the pixel-wise binary labels $\mathcal{M}_{hw}^{A,B}(x_n)$, we compute the pixel-wise average precision (PxAP) for the class pair (A, B) by computing the area under the precision-recall curve. We then take the mean of PxAP over all the class pairs of interest (e.g. those with $|\mathcal{P}^{A,B}| = 1$) to compute the mPxAP.

Qualitative results. See Figure 4 for the qualitative examples of CALM and baselines including CAM. We observe that the counterfactual attribution maps $s^A - s^B$ generated by CALM_{EM} and CALM_{ML} are more accurate than CAM and gradient-based attribution methods; CALM_{EM} attributions are qualitatively more precise than CALM_{ML}. CALM tends to assign close-to-zero attributions on irrelevant regions, while the baseline methods tend to produce noisy attributions. The sparsity of CALM makes it qualitatively more interpretable than the baselines.

Quantitative results. Table 1 shows the mPxAP scores for CALM and baseline methods for retrieving relevant pixels as attribution regions. We examine CUB class pairs with the number of parts with attribute differences $|\mathcal{P}^{A,B}| \in \{1, 2, 3\}$. We observe that CALM_{EM} outperforms the baselines in all three sets of class pairs, confirming the qualitative superiority of CALM_{EM} in Figure 4. CALM_{EM} attains 4.7%p better mPxAP than CAM on average over the three sets. CALM_{ML} tends to be sub-optimal, compared to CALM_{EM} (27.3% vs 34.3% mPxAP). The variants of gradients perform below a mere 20% mPxAP on average. In conclusion, the counterfactual attribution by CALM generates precise localization of the important bird parts that matter for the recognition task.



Figure 5. **Remove-and-classify results.** Classification accuracies of CNNs when $k\%$ of pixels are erased according to the attribution values s_{hw} . We show the relative accuracies \mathcal{R}_k against the random-erasing baseline. Lower is better.

5.3. Remove-and-classify results

One of the most widely used frameworks for evaluating the attribution task is the remove-and-classify evaluation [51, 30, 43, 46, 26]. Image pixels x_{hw} are erased in descending order of importance dictated by the attribution values s_{hw} . We write x_{-k} for the image where the top- $k\%$ important pixels are removed. We use the meaningful perturbation of the “blur” type [20] for erasing the pixels. A good attribution method shall assign high attribution values on important pixels; erasing them will quickly drop the classification accuracy \mathcal{A}_k with increasing k . We set the base reference accuracy \mathcal{A}_k^r as the classifier’s accuracy with $k\%$ of the pixels erased at random. For each method, we report the relative accuracy $\mathcal{R}_k = \mathcal{A}_k / \mathcal{A}_k^r$ for different k .

Results. We show the remove-and-classify results in Figure 5 for three image classification datasets. We observe that CALM variants show the lowest relative accuracies (lower is better) \mathcal{R}_k on cue-removed images in CUB and OpenImages, compared to CAM and other baselines. For the two datasets, CALM_{EM} attains values even close to zero at $k \in [10, 50]$. On ImageNet, CALM_{EM} outperforms the baselines with a smaller margin. Overall, CALM_{EM} selects the important pixels for recognition best.

Classification performances. We study the trade-off between interpretability and performance. The improved attribution performances come at the cost of decreased classification accuracies. Our models will be useful in applications that require great attribution performances at a small cost in model accuracies.

5.4. Weakly-supervised object localization (WSOL)

WSOL is related to but different from the attribution task. For WSOL, one learns to detect *object foreground re-*

Methods	CUB	Open	ImNet
Baseline	70.6	72.1	74.5
CALM _{EM}	71.8	70.1	70.4
CALM _{ML}	59.6	70.9	70.6

Table 2. **Classification accuracy.**

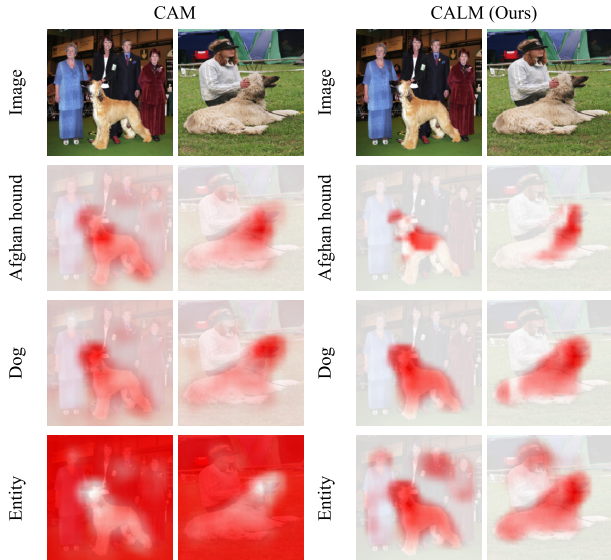


Figure 6. **Aggregating with class hierarchy.** “Afghan hound” is a descendant of “Dog”, which is a descendant of “Entity” in the WordNet hierarchy. Selecting a sensible superset \mathcal{Y} for aggregation lets CALM produce a high-quality foreground mask.

Methods	ImageNet	CUB	OpenImages
HaS [32]	62.6	60.6	57.4
ACoL [66]	61.1	60.0	56.3
SPG [67]	62.2	57.5	59.1
ADL [14]	61.6	61.1	56.9
CutMix [64]	62.2	60.8	59.5
InCA [27]	63.1	63.4	-
CAM [68]	62.4	61.1	60.0
CAM [68] + \mathcal{Y}	60.6	63.4	60.0
CALM _{EM}	62.5	52.5	62.7
CALM _{EM} + \mathcal{Y}	62.8	65.4	62.7
CALM _{ML}	62.6	61.3	62.3
CALM _{ML} + \mathcal{Y}	62.7	68.0	62.3

Table 3. **WSOL results on CUB, OpenImages, and ImageNet.** Average for ResNet, Inception, and VGG are reported for each dataset. CALM_{EM} and CALM_{EM} + \mathcal{Y} are compared against the baseline methods. “+ \mathcal{Y} ” denotes the aggregation.

gions with only image-label pairs. While the ingredients are identical ((X, Y) observed), the desired latent Z is different: the important cues for recognition may not necessarily agree with the object foreground regions. Nonetheless, the WSOL field benefits from the developments in attribution methods like CAM, which has remained the state of the art method for WSOL for the past few years [12].

We apply CALM to WSOL. Since attribution maps $p(\hat{y}, z|x)$ only point to sub-parts relevant for recognition, we aggregate the attributions them over multiple classes $p(\mathcal{Y}, z|x) = \sum_{y \in \mathcal{Y}} p(y, z|x)$ (subset attribution in §4.1.3) to fully cover the foreground regions.

Setting the superset \mathcal{Y} . For the ground-truth class \hat{y} , we set the superset \mathcal{Y} as the set of classes sharing the same part composition as \hat{y} . The intuition is that the attributions are mostly on object parts and that classes of such \mathcal{Y} have attributions spread across different object parts. For example, all bird classes in CUB [62] shall share the same superset $\mathcal{Y} = \{\text{all 200 birds}\}$, as they share the same body part composition. On the other extreme, 100 classes in OpenImages [8, 12] do not share part structures across classes. Thus, we always set $\mathcal{Y} = \{\hat{y}\}$. ImageNet1K [50] is mixed. Its 1000 classes include 116 dog species, but also many other objects and concepts that do not share the same part structure. For ImageNet, we have manually annotated the supersets \mathcal{Y} for every class \hat{y} , using the WordNet hierarchy [38]. Details in the Supplementary Materials.

Results. We evaluate WSOL performances based on the benchmarks and evaluation metrics in [12]. The benchmark considers 3 architectures (VGG, Inception, ResNet) and 3 datasets (CUB, OpenImages, ImageNet). Implementation details are in Supplementary Materials. We show results in Table 3. We observe that the aggregation significantly enhances the WSOL performances for CALM_{EM}: 52.5% to 65.4% on CUB. CALM_{EM} + \mathcal{Y} attains the best performances on CUB and OpenImages and second-best on ImageNet.

Analysis. We study the 116 fine-grained dog species in ImageNet more closely. We show the aggregation of attribution maps in Figure 6. CALM for \hat{y} fails to cover the full extent of the object. As the maps are aggregated over all dog species \mathcal{Y} , the map precisely covers the full extents of the dogs. However, if \mathcal{Y} covers all 1000 classes, the resulting saliency map $p(z|x)$ starts to include non-dog pixels.

6. Conclusion

Despite its great contributions to the field, the class activation mapping (CAM) is not as interpretable as it could be. It lacks communicability in practice and fails to meet key theoretical requirements for feature attribution methods. This paper has introduced a novel visual feature attribution method, class activation latent mapping (CALM). Based on the probabilistic treatment of the last layers of CNNs, CALM is interpretable by design. CALM satisfies the theoretical requirements as an attribution method and outperforms CAM and other baselines on attribution tasks.

Acknowledgement

This work has been partially funded by the ERC (853489 - DEXIM) and by the DFG (2064/1 – Project number 390727645). We thank members of NAVER AI Lab, especially Sanghyuk Chun for their helpful feedback and discussion. Kay Choi designed Figure 2.

References

- [1] Radhakrishna Achanta, Sheila Hemami, Francisco Estrada, and Sabine Süssstrunk. Frequency-tuned salient region detection. In *CVPR*, 2009. 7
- [2] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *NeurIPS*, 2018. 3
- [3] Chirag Agarwal and Sara Hooker. Estimating example difficulty using variance of gradients. *arXiv preprint arXiv:2008.11600*, 2020. 2, 6
- [4] David Alvarez-Melis and Tommi S Jaakkola. Towards robust interpretability with self-explaining neural networks. In *NeurIPS*, 2018. 2
- [5] 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*, 2015. 1, 2, 3
- [6] Wonho Bae, Junhyug Noh, and Gunhee Kim. Rethinking class activation mapping for weakly supervised object localization. In *ECCV*, 2020. 1, 2
- [7] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *JMLR*, 2010. 2
- [8] Rodrigo Benenson, Stefan Popov, and Vittorio Ferrari. Large-scale interactive object segmentation with human annotators. In *CVPR*, 2019. 6, 8
- [9] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *WACV*, 2018. 1, 2
- [10] Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K Su. This looks like that: Deep learning for interpretable image recognition. In *NeurIPS*, 2019. 2
- [11] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *PAMI*, 2017. 5
- [12] Junsuk Choe, Seong Joon Oh, Sanghyuk Chun, Zeynep Akata, and Hyunjung Shim. Evaluation for weakly supervised object localization: Protocol, metrics, and datasets. *arXiv preprint arXiv:2007.04178*, 2020. 2, 3, 6, 7, 8
- [13] Junsuk Choe, Seong Joon Oh, Seungho Lee, Sanghyuk Chun, Zeynep Akata, and Hyunjung Shim. Evaluating weakly supervised object localization methods right. In *CVPR*, 2020. 2, 3
- [14] Junsuk Choe and Hyunjung Shim. Attention-based dropout layer for weakly supervised object localization. In *CVPR*, 2019. 8
- [15] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *JRSS: Series B*, 1977. 4
- [16] Thibaut Durand, Taylor Mordan, Nicolas Thome, and Matthieu Cord. Wildcat: Weakly supervised learning of deep convnets for image classification, pointwise localization and segmentation. In *CVPR*, 2017. 2
- [17] Thibaut Durand, Nicolas Thome, and Matthieu Cord. Weldon: Weakly supervised learning of deep convolutional neural networks. In *CVPR*, 2016. 2
- [18] Daniel C Elton. Self-explaining ai as an alternative to interpretable ai. In *AGI*, 2020. 2
- [19] Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Understanding deep networks via extremal perturbations and smooth masks. In *ICCV*, 2019. 2
- [20] Ruth C Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *ICCV*, 2017. 1, 2, 7
- [21] Nicholas Frosst and Geoffrey Hinton. Distilling a neural network into a soft decision tree. *arXiv preprint arXiv:1711.09784*, 2017. 2
- [22] Ruigang Fu, Q. Hu, X. Dong, Yulan Guo, Y. Gao, and Biao Li. Axiom-based grad-cam: Towards accurate visualization and explanation of cnns. *ArXiv*, abs/2008.02312, 2020. 1, 3
- [23] Yash Goyal, Ziyang Wu, Jan Ernst, Dhruv Batra, Devi Parikh, and Stefan Lee. Counterfactual visual explanations. In *ICML*, 2019. 1, 2
- [24] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. A survey of methods for explaining black box models. *CSUR*, 2019. 1
- [25] David Gunning. Explainable artificial intelligence (xai). *DARPA*, 2017. 1
- [26] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability methods in deep neural networks. In *NeurIPS*, 2019. 2, 7
- [27] Minsong Ki, Youngjung Uh, Wonyoung Lee, and Hyeran Byun. In-sample contrastive learning and consistent attention for weakly supervised object localization. In *ACCV*, 2020. 8
- [28] 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 *ICML*, 2018. 1, 2
- [29] Hanjoo Kim, Minkyu Kim, Dongjoo Seo, Jinwoong Kim, Heungseok Park, Soeun Park, Hyunwoo Jo, KyungHyun Kim, Youngil Yang, Youngkwan Kim, et al. Nsm1: Meet the mlaas platform with a real-world case study. *arXiv preprint arXiv:1810.09957*, 2018. 6
- [30] Pieter-Jan Kindermans, Kristof T. Schütt, Maximilian Alber, Klaus-Robert Müller, Dumitru Erhan, Been Kim, and Sven Dähne. Learning how to explain neural networks: Patternnet and patternattribution. In *ICLR*, 2018. 2, 7
- [31] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. Concept bottleneck models. In *ICML*, 2020. 2
- [32] Krishna Kumar Singh and Yong Jae Lee. Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization. In *ICCV*, 2017. 8
- [33] Isaac Lage, Andrew Ross, Samuel J Gershman, Been Kim, and Finale Doshi-Velez. Human-in-the-loop interpretability prior. In *NeurIPS*, 2018. 2

- [34] Pantelis Linardatos, Vasilis Papastefanopoulos, and Sotiris Kotsiantis. Explainable ai: A review of machine learning interpretability methods. *Entropy*, 2021. 2, 5
- [35] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 3
- [36] Max Losch, Mario Fritz, and Bernt Schiele. Interpretability beyond classification output: Semantic bottleneck networks. *arXiv preprint arXiv:1907.10882*, 2019. 2
- [37] David JC MacKay and David JC Mac Kay. *Information theory, inference and learning algorithms*. Cambridge university press, 2003. 4
- [38] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 1995. 8
- [39] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 2017. 1, 2
- [40] Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 2018. 1, 3
- [41] Maxime Oquab, Léon Bottou, Ivan Laptev, and Josef Sivic. Is object localization for free?-weakly-supervised learning with convolutional neural networks. In *CVPR*, 2015. 2
- [42] George Papandreou, Liang-Chieh Chen, Kevin P. Murphy, and Alan L. Yuille. Weakly- and semi-supervised learning of a deep convolutional network for semantic image segmentation. In *ICCV*, 2015. 5
- [43] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. In *BMVC*, 2018. 1, 2, 7
- [44] Brett Poulin, Roman Eisner, Duane Szafron, Paul Lu, Russell Greiner, David S. Wishart, Alona Fyshe, Brandon Pearcy, Cam Macdonell, and John Anvik. Visual explanation of evidence with additive classifiers. In *AAAI*, 2006. 2
- [45] Harish Guruprasad Ramaswamy et al. Ablation-cam: Visual explanations for deep convolutional network via gradient-free localization. In *WACV*, 2020. 1
- [46] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *SIGKDD*, 2016. 2, 7
- [47] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. In *AAAI*, 2018. 1
- [48] Andrew Ross and Finale Doshi-Velez. Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients. In *AAAI*, 2018. 2
- [49] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *NMI*, 2019. 5
- [50] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 6, 8
- [51] Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE TNNLS*, 2016. 2, 7
- [52] Wojciech Samek, Grégoire Montavon, Sebastian Lapuschkin, Christopher J Anders, and Klaus-Robert Müller. Explaining deep neural networks and beyond: A review of methods and applications. *Proceedings of the IEEE*, 2021. 1
- [53] 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 *ICCV*, 2017. 1, 2, 3
- [54] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *ICML*, 2017. 2, 3
- [55] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *ICLRW*, 2014. 1, 2, 6
- [56] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017. 1, 2, 6
- [57] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. In *ICLRW*, 2015. 2
- [58] Suraj Srinivas and François Fleuret. Full-gradient representation for neural network visualization. In *NeurIPS*, 2019. 1, 2
- [59] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *ICML*, 2017. 1, 2, 3, 6
- [60] AH Sung. Ranking importance of input parameters of neural networks. *Expert Systems with Applications*, 1998. 2
- [61] Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and Xia Hu. Score-cam: Score-weighted visual explanations for convolutional neural networks. In *CVPRW*, 2020. 1
- [62] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010. 2, 3, 6, 8
- [63] Yiding Yang, Jiayan Qiu, Mingli Song, Dacheng Tao, and Xinchao Wang. Learning propagation rules for attribution map generation. In *ECCV*. Springer, 2020. 1, 2
- [64] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, 2019. 8
- [65] Jianming Zhang, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. Top-down neural attention by excitation backprop. In *ECCV*, 2016. 2
- [66] Xiaolin Zhang, Yunchao Wei, Jiashi Feng, Yi Yang, and Thomas S Huang. Adversarial complementary learning for weakly supervised object localization. In *CVPR*, 2018. 8
- [67] Xiaolin Zhang, Yunchao Wei, Guoliang Kang, Yi Yang, and Thomas Huang. Self-produced guidance for weakly-supervised object localization. In *ECCV*, 2018. 8
- [68] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016. 1, 2, 3, 6, 8

- [69] Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Interpretable basis decomposition for visual explanation. In *ECCV*, 2018. [1](#)
- [70] Luisa M Zintgraf, Taco S Cohen, Tameem Adel, and Max Welling. Visualizing deep neural network decisions: Prediction difference analysis. In *ICLR*, 2017. [1](#), [2](#)