

Grad-CAM: Why did you say that?

Visual Explanations from Deep Networks via Gradient-based Localization

Ramprasaath R. Selvaraju

Abhishek Das

Ramakrishna Vedantam

Michael Cogswell

Devi Parikh

Dhruv Batra

Virginia Tech

{ram21, abhshkdz, vrama91, cogswell, parikh, dbatra}@vt.edu

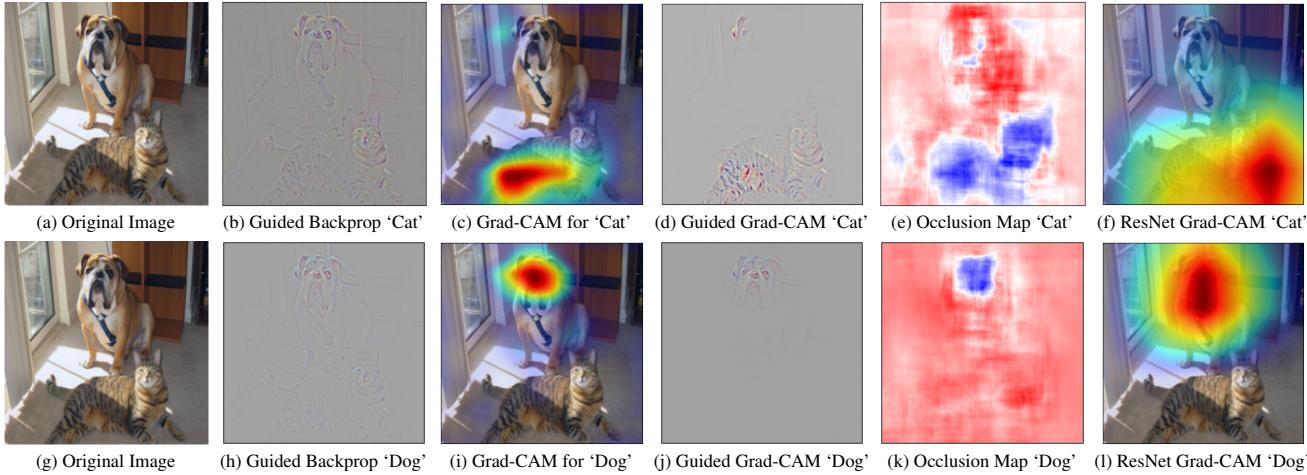


Figure 1. (a) Original image with a cat and a dog. (b-f) Support for the cat category according to various visualizations. (b) Guided Backpropagation [42]: provides high-resolution visualization of contributing features, (c, f) Grad-CAM (Ours): localizes class-discriminative regions, (d) Combining (b) and (c) gives Guided Grad-CAM, which gives high-resolution visualizations that are class-discriminative. Interestingly, the localizations achieved by our Grad-CAM technique (c) are very similar to results from occlusion sensitivity (e), while being much cheaper to compute. (f, l) are Grad-CAM visualizations for ResNet-18 layer. Note that in (e, k), blue corresponds to evidence for the class while in (d, f, i, l) blue indicates regions with low score for the class. Figure best viewed in color.

Abstract

We propose a technique for making Convolutional Neural Network (CNN)-based models more transparent by visualizing the input regions that are ‘important’ for predictions from these models – producing visual explanations.

Our approach, called Gradient-weighted Class Activation Mapping (Grad-CAM), uses the class-specific gradient information flowing into the final convolutional layer of a CNN to produce a coarse localization map of the important regions in the image for each class. Grad-CAM is a strict generalization of Class Activation Mapping (CAM) [47]. While CAM is limited to a narrow class of CNN models, Grad-CAM is broadly applicable to any CNN-based architectures and needs no re-training. We also show how Grad-CAM may be combined with existing pixel-space visualizations (such as Guided Backpropagation [42]) to create a high-resolution class-discriminative visualization (Guided Grad-CAM).

We generate Grad-CAM and Guided Grad-CAM visual explanations to better understand image classification, image captioning, and visual question answering (VQA) models, including Res-Net based architectures. In the context of image classification models, our visualizations (a) lend

insight into failure modes of these models (showing that seemingly unreasonable predictions have reasonable explanations), and (b) outperform pixel-space gradient visualizations (Guided Backpropagation [42] and Deconvolution [45]) on the ILSVRC-15 weakly supervised localization task. For image captioning and VQA, our visualizations expose the somewhat surprising insight that common CNN + Long Short Term Memory (LSTM) models can often be good at localizing discriminative input image regions despite not being trained on grounded image-text pairs.

Finally, we design and conduct human studies to measure if Guided Grad-CAM explanations help users establish trust in the predictions made by deep networks. Interestingly, we show that Guided Grad-CAM helps untrained users successfully discern a ‘stronger’ deep network from a ‘weaker’ one even when both networks make identical predictions, simply on the basis of their different explanations.

Our code is available at <https://github.com/rampsr/grad-cam/> and a demo is available on CloudCV [2]¹. Video of the demo can be found at youtu.be/COjUB9Izk6E.

¹<http://gradcam.cloudcv.org>

1. Introduction

Convolutional Neural Networks (CNNs) and other deep networks have enabled unprecedented breakthroughs in a variety of computer vision tasks, from image classification [25, 16] to object detection [15], semantic segmentation [29], image captioning [43, 6, 12, 21], and more recently, visual question answering [3, 14, 33, 37]. While these deep neural networks enable superior performance, their lack of decomposability into *intuitive and understandable* components makes them hard to interpret [28]. Consequently, when today’s intelligent systems fail, they fail spectacularly disgracefully, without warning or explanation, leaving a user staring at incoherent output, wondering why.

Interpretability Matters. In order to build trust in intelligent systems and move towards their meaningful integration into our everyday lives, it is clear that we must build ‘transparent’ models that explain *why they predict what they do*. Broadly speaking, this transparency is useful at three different stages of Artificial Intelligence (AI) evolution. First, when AI is significantly weaker than humans and not yet reliably ‘deployable’ (e.g. visual question answering [3]), the goal of transparency and explanations is to identify the failure modes [1, 17], thereby helping researchers focus their efforts on the most fruitful research directions. Second, when AI is on par with humans and reliably ‘deployable’ (e.g., image classification [22] on a set of categories trained on sufficient data), the goal is to establish trust with users. Third, when AI is significantly stronger than humans (e.g. chess or Go playing bots [39]), the goal of transparency and explanations is machine teaching [20] – *i.e.*, a machine teaching a human on how to make accurate predictions.

There typically exists a trade-off between accuracy and simplicity or interpretability. Classical rule-based or expert systems [18] are highly interpretable but not very accurate (or robust). Decomposable pipelines where each stage is hand-designed are thought to be more interpretable as each individual component assumes a natural intuitive explanation. By using deep models, we sacrifice interpretable modules for uninterpretable ones that achieve greater performance through greater abstraction (more layers) and tighter integration (end-to-end training). Recently introduced deep residual networks (ResNets) [16] are over 200-layers deep and have shown state-of-the-art accuracy in several challenging tasks. Due to their complexity they have been highly uninterpretable. As such deep models are beginning to explore the spectrum between interpretability and accuracy.

Zhou *et al.* [47] recently proposed a technique called Class Activation Mapping (CAM) for identifying discriminative regions used by a particular class of modified image classification CNNs (not containing any fully-connected layers). In essence, this work trades off model complexity for more transparency into the working of the model. In contrast, we make existing state-of-the-art deep models interpretable

without altering the architecture, thus avoiding a tradeoff between interpretability and accuracy. Our approach is a generalization of CAM [47] to any CNN-based architecture (CNNs with fully-connected layers, CNNs stacked with Recurrent Neural Networks (RNNs), ResNets *etc.*).

What makes a good visual explanation? Consider image classification [9] – a ‘good’ visual explanation from the model justifying a predicted class should be (a) class-discriminative (*i.e.* localize the category in the image) and (b) high-resolution (*i.e.* capture fine-grained detail).

Fig. 1 shows outputs from a number of visualizations for the ‘tiger cat’ class (top) and ‘boxer’ (dog) class (bottom). Pixel-space gradient visualizations such as Guided Backpropagation [42] and Deconvolution [45] are high-resolution and highlight fine-grained details in the image, but are not class-discriminative (for example, the visualization for both ‘cat’ and ‘dog’ in Figures 1b and 1h are very similar).

In contrast, our approach (Grad-CAM) shown in Figures 1c, 1f and 1i, 1l, is highly class-discriminative (*i.e.* the ‘cat’ explanation exclusively highlights the ‘cat’ regions, and not the ‘dog’ regions and *vice versa*). Note that these very closely match with the occlusion maps generated through multiple forward passes (Figures 1e and 1k). The spatial resolution of the most class-discriminative Grad-CAM maps is the size of the last convolution layer in the CNN, which is typically small (e.g. 14×14 in VGGNet [41]) and hence does not show fine-grained details.

In order to combine the best of both worlds, we show that it is possible to fuse existing pixel-space gradient visualizations with Grad-CAM to create Guided Grad-CAM visualizations that are both high-resolution and class-discriminative. As a result, important regions of the image which correspond to a class of interest are visualized in high-resolution detail even if the image contains multiple classes, as shown in Figures 1d and 1j. When visualized for ‘tiger cat’, Guided Grad-CAM not only highlights the cat regions, but also highlights the stripes on the cat which is important for predicting that particular variety of cat.

To summarize, our contributions are as follows:

(1) We propose a class-discriminative localization technique called Gradient-weighted Class Activation Mapping (Grad-CAM) that can be used to generate visual explanations from *any* CNN-based network without requiring architectural changes. We evaluate Grad-CAM for weakly-supervised image localization on ImageNet (Section 4.1) where it outperforms pixel-space gradients and contrastive-Marginal Winning Probability (c-MWP) [46].

(2) To illustrate the broad applicability of our technique across tasks, we apply Grad-CAM to state-of-the-art classification, image captioning (Section 7.1), and visual question answering (Section 7.2), effectively visualizing the image support for predictions from such networks. For image classification, our visualizations lend insight into failure modes

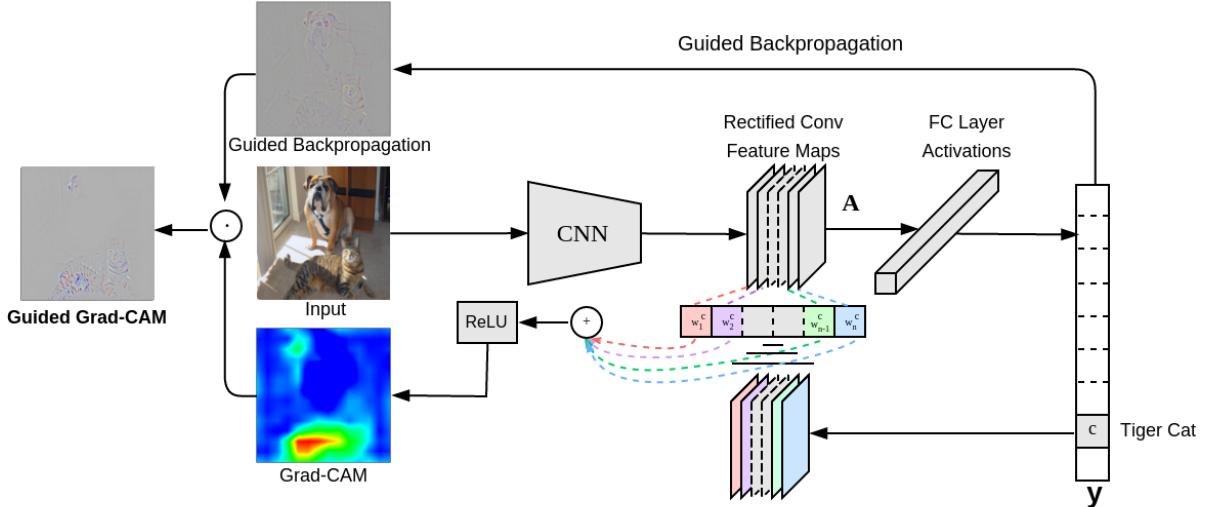


Figure 2: Grad-CAM overview: Given an image, and a category (‘tiger cat’) as input, we forward propagate the image through the model to obtain the raw class scores before softmax. The gradients are set to zero for all classes except the desired class (tiger cat), which is set to 1. This signal is then backpropagated to the rectified convolutional feature map of interest, where we can compute the coarse Grad-CAM localization (blue heatmap). Finally, we pointwise multiply the heatmap with guided backpropagation to get Guided Grad-CAM visualizations which are both high-resolution and class-discriminative.

of current generation CNNs (Section 6), showing that seemingly unreasonable predictions have reasonable explanations. For image captioning and VQA, our visualizations expose the somewhat surprising insight that common CNN + Long Short Term Memory (LSTM) models can often be good at localizing discriminative input image regions despite not being trained on grounded image-text pairs.

(3) We visualize the recently introduced ResNets [16] for the task of image classification and VQA (Section 7.2). We observe that while visualizing different layers from top to bottom, the discriminative ability of Grad-CAM significantly reduces as we encounter skip connections connecting layers with different output dimensionality. We analyze this further in the supplementary Section G.

(4) We design and conduct human studies to show that Guided Grad-CAM explanations are class-discriminative and help humans not only establish trust, but also helps untrained users successfully discern a ‘stronger’ network from a ‘weaker’ one *even when both make identical predictions, simply on the basis of their different explanations*.

2. Related Work

Our work draws on recent work in CNN visualizations, model trust assessment, and weakly-supervised localization.

Visualizing CNNs. A number of previous works [40, 42, 45, 13] have visualized CNN predictions by highlighting ‘important’ pixels (*i.e.* change in intensities of these pixels have the most impact on the prediction’s score). Specifically, Simonyan *et al.* [40] visualize partial derivatives of predicted class scores w.r.t. pixel intensities, while Guided Backpropagation [42] and Deconvolution [45] make modifications to ‘raw’ gradients that result in qualitative improvements. Despite producing fine-grained visualizations, these methods are not class-discriminative. Visualizations with respect to different classes are nearly identical (see Figures 1b and 1h).

Other visualization methods synthesize images to maximally activate a network unit [40, 11] or invert a latent representation [32, 10]. Although these can be high-resolution and class-discriminative, they visualize a model overall and not predictions for specific input images.

Assessing Model Trust. Motivated by notions of interpretability [28] and assessing trust in models [38], we evaluate Grad-CAM visualizations in a manner similar to [38] via human studies to show that they can be important tools for users to evaluate and place trust in automated systems.

Weakly supervised localization. Another relevant line of work is weakly supervised localization in the context of CNNs, where the task is to localize objects in images using only whole image class labels [7, 34, 35, 47].

Most relevant to our approach is the Class Activation Mapping (CAM) approach to localization [47]. This approach modifies image classification CNN architectures replacing fully-connected layers with convolutional layers and global average pooling [26], thus achieving class-specific feature maps. Others have investigated similar methods using global max pooling [35] and log-sum-exp pooling [36].

A drawback of CAM is that it requires feature maps to directly precede softmax layers, so it is only applicable to a particular kind of CNN architectures performing global average pooling over convolutional maps immediately prior to prediction (*i.e.* conv feature maps → global average pooling → softmax layer). Such architectures may achieve inferior accuracies compared to general networks on some tasks (*e.g.* image classification) or may simply be inapplicable to other tasks (*e.g.* image captioning or VQA). We introduce a new way of combining feature maps using the gradient signal that does not require *any* modification in the network architecture. This allows our approach to be applied to any CNN-based architecture, including those for image caption-

ing and visual question answering. For a fully-convolutional architecture, Grad-CAM reduces to CAM (ignoring rectification/normalization for visualization purposes). Thus, Grad-CAM is a strict generalization to CAM.

Other methods approach localization by classifying perturbations of the input image. Zeiler and Fergus [45] perturb inputs by occluding patches and classifying the occluded image, typically resulting in lower classification scores for relevant objects when those objects are occluded. This principle is applied for localization in [4]. Oquab *et al.* [34] classify many patches containing a pixel then average these patch class-wise scores to provide the pixel’s class-wise score. Unlike these, our approach achieves localization in one shot; it only requires a single forward and a partial backward pass per image and thus is typically an order of magnitude more efficient. In recent work Zhang *et al.* [46] introduce contrastive Marginal Winning Probability (c-MWP), a probabilistic Winner-Take-All formulation for modelling the top-down attention for neural classification models which can highlight discriminative regions. This is slower than Grad-CAM and only works for Image Classification CNNs. Moreover, qualitative and quantitative results are worse than for Grad-CAM (see Sec. 4.1 and supplementary Section F).

3. Approach

We briefly recap the localization approach presented in CAM [47], and then describe our generalization, Grad-CAM. Then we describe how this class-discriminative but coarse localization technique can be combined with high-resolution visualizations obtained using Deconvolution and Guided Backpropagation to obtain both desirable properties (high-resolution and class-discrimination).

Class Activation Mapping (CAM). Recall that CAM [47] produces a localization map for an image classification CNN with a specific kind of architecture where global average pooled convolutional feature maps are fed directly into a softmax. Specifically, let the penultimate layer produce K feature maps $A^k \in \mathbb{R}^{u \times v}$ of width u and height v . These feature maps are then spatially pooled using Global Average Pooling (GAP) and linearly transformed to produce a score y^c for each class c

$$y^c = \sum_k w_k^c \underbrace{\frac{1}{Z} \sum_i \sum_j}_{\text{global average pooling}} A_{ij}^k \quad (1)$$

class feature weights feature map

To produce a localization map $L_{\text{CAM}}^c \in \mathbb{R}^{u \times v}$ for class c , CAM computes the linear combination of the final feature maps using the learned weights of the final layer:

$$L_{\text{CAM}}^c = \underbrace{\sum_k w_k^c A^k}_{\text{linear combination}}. \quad (2)$$

This is normalized to lie between 0 and 1 for visualization purposes. To apply CAM to a network which uses multi-

ple fully-connected layers before the final layer, the fully connected layers are replaced with GAP and the network is re-trained.

Gradient-weighted Class Activation Mapping. In order to obtain the class-discriminative localization map Grad-CAM $L_{\text{Grad-CAM}}^c \in \mathbb{R}^{u \times v}$ in general architectures, we first compute the gradient of y^c with respect to feature maps A of a convolutional layer, *i.e.* $\frac{\partial y^c}{\partial A_{ij}^k}$. These gradients flowing back are global-average-pooled to obtain weights α_k^c :

$$\alpha_k^c = \underbrace{\frac{1}{Z} \sum_i \sum_j}_{\text{global average pooling}} \underbrace{\frac{\partial y^c}{\partial A_{ij}^k}}_{\text{gradients via backprop}} \quad (3)$$

This weight α_k^c represents a *partial linearization* of the deep network downstream from A , and captures the ‘importance’ of feature map k for a target class c ². Empirically, using the averaged gradient through Global Average Pooling (GAP) is more robust to noise in the gradients and thus leads to better localizations than other choices like taking the Global Max Pooling, as shown in the supplementary Section C. In general, y^c need not be the class score produced by an image classification CNN, and could be any differentiable activation. To be concrete we introduce Grad-CAM using the notion ‘class’ from image classification (e.g., cat or dog), but visual explanations can be considered for any differentiable node in a computational graph, including words from a caption or the answer to a question.

As in CAM, our Grad-CAM heat-map is a weighted combination of feature maps, but we follow this by a ReLU:

$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left(\underbrace{\sum_k \alpha_k^c A^k}_{\text{linear combination}} \right) \quad (4)$$

Notice that this results in a coarse heat-map of the same size as the convolutional feature maps (14×14 in the case of last convolutional layers of VGG [41] and AlexNet [25] networks). For visualization as in Figure 1c, 1i, and others, $L_{\text{Grad-CAM}}^c$ is normalized so the values lie between 0 and 1. For architectures where CAM is applicable – *i.e.*, fully-convolutional CNNs with A being the output of the final conv layer, followed by global average pooling and a linear classification layer with softmax – the weights used in CAM w_k^c are precisely α_k^c . Other than the ReLU in (4), this makes Grad-CAM a generalization of CAM. See supplementary Section A for derivation. The motivation for the ReLU is the following – we are only interested in the features that have a *positive* influence on the class of interest, *i.e.* pixels whose intensity should be *increased* in order to increase

² We performed some initial experiments using variants of gradients used in the computation of α_k^c . The “gradients” computed using Deconvolution led to non-discriminative localizations while Guided Backpropagation led to somewhat discriminative localizations, though much less so than normal unmodified gradients.

y^c . Since the feature maps A are already non-negative, we rectify the heat-map obtained in (4) to highlight such pixels. And as expected, without this ReLU, localization maps sometimes highlight more than just the desired class and achieve lower localization performance (see supplementary Sec. C). Intuitively, negative values indicate pixels that are likely to belong to other categories in the image, and the application of ReLU excludes them. Figures 1c, 1f and 1i, 11 show Grad-CAM visualizations for ‘tiger cat’ and ‘boxer (dog)’ respectively. More Grad-CAM visualizations can be found in supplementary Sec. B.

Note that the above generalization also allows us to generate visual explanations from CNN-based models that cascade convolutional layers with more complex interactions. Indeed, we apply Grad-CAM to “beyond classification” tasks and models that utilize CNNs for image captioning and Visual Question Answering (VQA) (Sec. 7.2).

Guided Grad-CAM. While Grad-CAM visualizations are class-discriminative and localize relevant image regions well, they lack the ability to show fine-grained importance like pixel-space gradient visualization methods (Guided Backpropagation and Deconvolution). For example in Figure 1c, Grad-CAM can easily localize the cat region; however, it is unclear from the low-resolutions of the heat-map why the network predicts this particular instance is ‘tiger cat’. In order to combine the best aspects of both, we fuse Guided Backpropagation and Grad-CAM visualizations via pointwise multiplication (L_{CAM}^c is first up-sampled to the input image resolution using bi-linear interpolation). Fig. 2 bottom-left illustrates this fusion. This visualization is both high-resolution (when the class of interest is ‘tiger cat’, it identifies important ‘tiger cat’ features like stripes, pointy ears and eyes) and class-discriminative (it shows the ‘tiger cat’ but not the ‘boxer (dog)’). Replacing Guided Backpropagation with Deconvolution in the above gives similar results, but we found Deconvolution to have artifacts (and Guided Backpropagation visualizations were generally noise-free), so we chose Guided Backpropagation over Deconvolution. A number of works have asserted that as the depth of a CNN increases, higher-level visual constructs are captured [5, 32]. Furthermore, convolutional layers naturally retain spatial information which is lost in fully-connected layers, so we expect the last convolutional layers to have the best compromise between high-level semantics and detailed spatial information. We provide Grad-CAM visualizations computed at various convolutional layers in the supplementary Sec. C to demonstrate this effect.

4. Evaluating Localization

4.1. Weakly-supervised Localization

In this section, we evaluate the localization capability of Grad-CAM in the context of image classification. The ImageNet localization challenge [9] requires competing ap-

proaches to provide bounding boxes in addition to classification labels. Similar to classification, evaluation is performed for both the top-1 and top-5 predicted categories. Similar to Zhou *et al.* [47], given an image, we first obtain class predictions from our network. Next, we generate Grad-CAM localization maps for each of the predicted classes and binarize with threshold of 15% of the max intensity. This results in connected segments of pixels and we draw our bounding box around the single largest segment.

We evaluate the pretrained off-the-shelf VGG-16 [41] model from the Caffe [19] Model Zoo. Following ILSVRC evaluation, we report both top-1 and top-5 localization error on ILSVRC-15 validation set in Table. 1. Grad-CAM localization errors are significantly lower than those achieved by c-MWP [46] and Simonyan *et al.* [40] for the VGG-16 model, which uses grabcut to post-process image space gradients into heat maps. We also see that CAM achieves a slightly better localization, but requires a change in the VGG architecture, necessitates re-training, and achieves worse top-1 val classification error (2.76% increase) [47], whereas our model makes no compromise on classification accuracy.

Method	Top-1 loc error	Top-5 loc error	Top-1 cls error	Top-5 cls error
Backprop on VGG-16 [40]	61.12	51.46	30.64	11.03
c-MWP on VGG-16 [46]	70.92	63.04	30.64	11.03
Grad-CAM on VGG-16	57.80	48.03	30.64	11.03
VGG-16-GAP [47]	57.20	45.14	33.40	12.20

Table 1: Classification and Localization results on ILSVRC-15 val. Grad-CAM outperforms [40] and c-MWP [46]. CAM [47] achieves slightly better localization via a modified architecture which results in higher classification error. Note that these networks are off-the-shelf classification CNNs, trained only with class labels and not bounding box annotations.

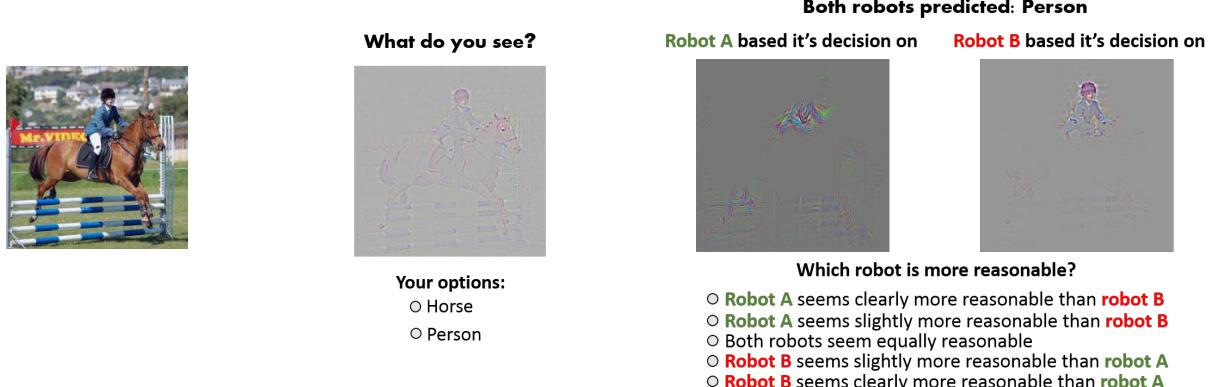
4.2. Weakly-supervised Segmentation

We also perform weakly-supervised Segmentation using the architecture from SEC [24]. We provide more details along with qualitative examples in the supplementary Sec. D.

4.3. Pointing Game

Zhang *et al.* [46] introduced the Pointing Game experiment to evaluate the discriminativeness of different attention maps for localizing target objects in scenes. Given a pretrained CNN classifier, their evaluation protocol cues each classifier with the ground-truth object label. It extracts the maximum point on the generated heatmap and evaluates if it lies in one of the annotated instances of the cued object category, thereby a hit or a miss is counted. The localization accuracy is then calculated as $\text{Acc} = \frac{\# \text{Hits}}{\# \text{Hits} + \# \text{Misses}}$. However this evaluation only measures the precision aspect of the CNN classifier. Hence we modify the pointing experiment to also measure the recall, as follows. We compute the visualization for the top-5 class predictions from the CNN classifiers³ and evaluate them using the pointing game setup with an additional option that a visualization may reject any of the top-5 predictions from the model if the max value in the visualization is below a threshold (in which case, it gets that as a hit, not a miss). More details on the experimental

³We use the GoogLeNet CNN finetuned on COCO provided in [46].



(a) Raw input image. Note that this is not a part of the tasks (b) and (c)

(b) AMT interface for evaluating the class-discriminative property

(c) AMT interface for evaluating if our visualizations instill trust in an end user

Figure 3: We evaluate visualizations extracted from image (a) for class discrimination (b) and trust worthiness (c). The class discrimination study measures if people can tell from the visualization which class is being visualized. Trust evaluations check if a visualization helps humans place trust in a more accurate classifier. Our results show that our Grad-CAM approach outperforms baseline approaches (Guided Backpropagation and Deconvolution) on both tasks.

setup can be found in the supplementary Sec. E. We find that our approach Grad-CAM outperforms c-MWP [46] by a significant margin (70.58% vs 60.30%).

We observe that c-MWP highlights arbitrary regions for predicted but non-existent categories, unlike Grad-CAM maps which seem more reasonable. Qualitative examples comparing c-MWP [46] and Grad-CAM on COCO, imageNet, and PASCAL categories can be found in supplementary Sec. F.

5. Evaluating Visualizations

Our first human study evaluates the main premise of our approach: are Grad-CAM visualizations more class-discriminative than previous techniques? Moreover, we want to understand if our class-discriminative interpretations can lead an end user to trust the visualized models.

For these experiments, we use VGG and AlexNet CNNs fine-tuned on PASCAL VOC 2007 train set, and the validation set is used to generate visualizations.

5.1. Evaluating Class Discrimination

We select images from PASCAL VOC 2007 val set that contain exactly two annotated categories, and create visualizations for one of the classes. For both VGG-16 and AlexNet CNNs, we obtain visualizations using four techniques: Deconvolution, Guided Backpropagation, and Grad-CAM versions of each these methods (Deconvolution Grad-CAM and Guided Grad-CAM). We show visualizations to workers on Amazon Mechanical Turk (AMT) and ask them “Which of the two object categories is depicted in the image?”, as shown in Fig. 3a. This task measures if people can tell from the visualization which class is being visualized.

Intuitively, a good prediction explanation is one that produces discriminative visualizations for the class of interest. The experiment was conducted using all 4 visualizations for 90 image-category pairs (*i.e.* 360 visualizations); 9 ratings were collected for each image, evaluated against the ground truth and averaged to obtain the accuracy. When viewing Guided Grad-CAM, human subjects can correctly

identify the category being visualized in 61.23% of cases (compared to 44.44% for Guided Backpropagation; thus, Grad-CAM improves human performance by 16.79%). Similarly, we also find that Grad-CAM helps make Deconvolution more class-discriminative (from 53.33% to 61.23%). Guided Grad-CAM performs the best among all the methods. Interestingly, our results seem to indicate that Deconvolution is more class discriminative than Guided Backpropagation, although Guided Backpropagation is more aesthetically pleasing than Deconvolution. To the best of our knowledge, our evaluations are the first to quantify these subtle differences.

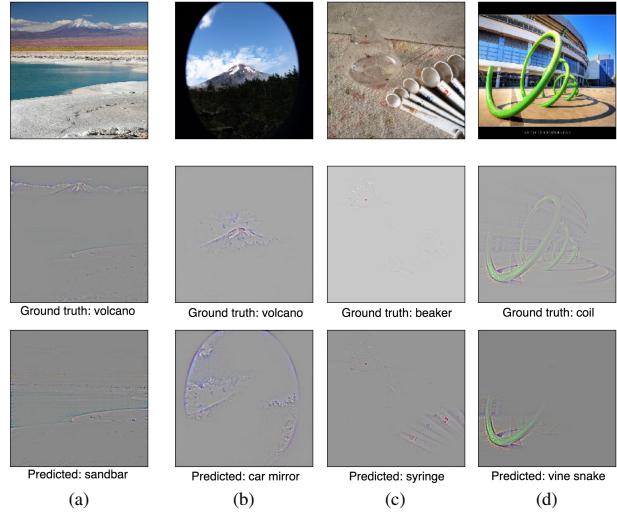


Figure 4: In these cases the model (VGG-16) failed to predict the correct class as its top 1 prediction, but it even failed to predict the correct class in its top 5 for figure b. All of these errors are due in part to class ambiguity. (a) For example, the network predicts ‘sandbar’ based on the foreground of (a), but it also knows where the correct label ‘volcano’ is located. (c-d) In other cases, these errors are still reasonable but not immediately apparent. For example, humans would find it hard to explain the predicted ‘syringes’ in (c) without looking at the visualization for the predicted class.

5.2. Evaluating Trust

Given two prediction explanations, we want to evaluate which seems more trustworthy. We use AlexNet and VGG-16 to compare Guided Backpropagation and Guided Grad-CAM visualizations, noting that VGG-16 is known to be

more reliable than AlexNet with an accuracy of 79.09 mAP *vs.* 69.20 mAP for AlexNet. In order to tease apart the efficacy of the visualization from the accuracy of the model being visualized, we consider only those instances where *both* models made the same prediction as ground truth. Given a visualization from AlexNet and one from VGG-16, and the predicted object category, workers are instructed to rate the reliability of the models relative to each other on a scale of clearly more/less reliable (+/-2), slightly more/less reliable (+/-1), and equally reliable (0). This interface is shown in Fig. 3c. To eliminate any biases, VGG and AlexNet were assigned to be *model1* with approximately equal probability. Remarkably, we find that human subjects are able to identify the more accurate classifier (VGG over AlexNet) *despite viewing identical predictions from the two, simply from the different explanations* generated from the two. With Guided Backpropagation, humans assign VGG an average score of 1.00 which means that it is slightly more reliable than AlexNet, while Guided Grad-CAM achieves a higher score of 1.27 which is closer to the option saying that VGG is clearly more reliable. Thus our Guided Grad-CAM visualization can help users place trust in a model that can generalize better, based on individual prediction explanations.

5.3. Faithfulness *vs.* Interpretability

Faithfulness of a visualization to a model is its ability to accurately explain the function learned by the model. Naturally, there exists a tradeoff between the interpretability and faithfulness of a visualization: a more faithful visualization is typically less interpretable and *vice versa*. In fact, one could argue that a fully faithful explanation is the entire description of the model, which in the case of deep models is not interpretable/easy to visualize. We have verified in previous sections that our visualizations are reasonably interpretable. We now evaluate how faithful they are to the underlying model. One expectation is that our explanations should be locally accurate, *i.e.* in the vicinity of the input data point, our explanation should be faithful to the model [38].

For comparison, we need a reference explanation with high local-faithfulness. One obvious choice for such a visualization is image occlusion [45], where we measure the difference in CNN scores when patches of the input image are masked out. Interestingly, patches which change the CNN score are also patches to which Guided Grad-CAM assigns high intensity, achieving rank correlation 0.261 (*vs.* 0.168 achieved by Guided Backpropagation) averaged over 2510 images in PASCAL 2007 val set. This shows that Guided Grad-CAM is more faithful to the original model than Guided Backpropagation.

6. Analyzing Failure Modes for VGG-16

We use Guided Grad-CAM to analyze failure modes of the VGG-16 CNN on ImageNet classification [9]. In order to see what mistakes a network is making we first get a list of ex-

amples that the network (VGG-16) fails to classify correctly. For the misclassified examples, we use Guided Grad-CAM to visualize both the correct and the predicted class. A major advantage of Guided Grad-CAM visualization over other methods that allows for this analysis is its high-resolution and its ability to be highly class-discriminative. As seen in Fig. 4, some failures are due to ambiguities inherent in ImageNet classification. We can also see that *seemingly unreasonable predictions have reasonable explanations*, which is a similar observation to HOGgles [44].

7. Image Captioning and VQA

In this subsection we apply our Grad-CAM technique to the image captioning [6, 21, 43] and Visual Question Answering (VQA) [3, 14, 33, 37] tasks. We find that Grad-CAM leads to interpretable visual explanations for these tasks as compared to baseline visualizations which do not change noticeably across different predictions. Note that existing visualization techniques such as Guided Backpropagation, Deconvolution, or c-MWP are either not class-discriminative, or simply cannot be used for these tasks or architectures, or both.

7.1. Image Captioning

In this section, we visualize spatial support for a simple image captioning model (without attention) using Grad-CAM visualizations. More specifically, we build on top of the publicly available ‘neuraltalk2’⁴ implementation [23] that makes use of a finetuned VGG-16 CNN for images and an LSTM-based language model. Given a caption, we compute the gradient of its log probability w.r.t. units in the last convolutional layer of the CNN (*conv5_3* for VGG-16) and generate Grad-CAM visualizations as described in Section 3. See Fig. 5a. For first example, the Grad-CAM maps for the generated caption localize every occurrence of both the kites and people inspite of their relatively small size. In the next example, see how Grad-CAM correctly highlights the pizza and the man, but ignores the woman nearby, since ‘woman’ is not mentioned in the caption. More qualitative examples can be found in the supplementary Sec. B.

Comparison to dense captioning. Johnson *et al.* [21] recently introduced the Dense Captioning (DenseCap) task that requires a system to jointly localize and caption salient regions in a given image. The model proposed in [21] consists of a Fully Convolutional Localization Network (FCLN) and an LSTM-based language model that produces both bounding boxes for regions of interest and associated captions in a single forward pass. Using the DenseCap model, we generate region-specific captions. Next, we visualize Grad-CAM localizations for these region-specific captions using the simple captioning model described earlier (neuraltalk2). Interestingly, we observe that Grad-CAM localizations correspond to regions in the image that the DenseCap model described, even though the holistic captioning model was not

⁴<https://github.com/karpathy/neuraltalk2>



(a) Image captioning explanations



(b) Comparison to DenseCap

Figure 5: Interpreting image captioning models: We use our class-discriminative localization technique, Grad-CAM to find spatial support regions for captions in images. Fig. 5a Visual explanations from image captioning model [23] highlighting image regions considered to be important for producing the captions. Fig. 5b Grad-CAM localizations of a *global* or *holistic* captioning model for captions generated by a dense captioning model [21] for the three bounding box proposals marked on the left. We can see that we get back Grad-CAM localizations (right) that agree with those bounding boxes – even though the captioning model and Grad-CAM techniques do not use any bounding box annotations.

trained with any region or bounding-box level annotations.

Results are shown in Fig. 5b.

7.2. Visual Question Answering

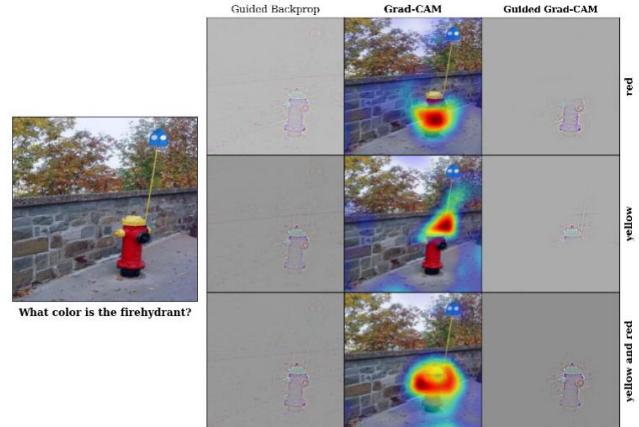
Typical VQA pipelines [3, 14, 33, 37] consist of a CNN to model images and an RNN language model for questions. The image and the question representations are fused to predict the answer, typically with a 1000-way classification. Since this is a classification problem, we pick an answer (the score y_c in (1)) and use its score to compute Grad-CAM to show image evidence that supports the answer. Despite the complexity of the task, involving both visual and language components, the explanations (of the VQA model from [30]) described in Fig. 6 are surprisingly intuitive and informative.

Comparison to Human Attention. Das *et al.* [8] collected human attention maps for a subset of the VQA dataset [3]. These maps have high intensity where humans looked in the image in order to answer a visual question. Human attention maps are compared to Grad-CAM visualizations for the simple VQA model introduced above [30] on 1374 val question-image (QI) pairs from [3] using the rank correlation evaluation protocol developed in [8]. Grad-CAM and human attention maps have a correlation of 0.136, which is statistically higher than chance or random attention maps (zero correlation). This shows that despite not being trained on grounded image-text pairs, CNN+LSTM based VQA models are surprisingly good at localizing discriminative regions required to output a particular answer.

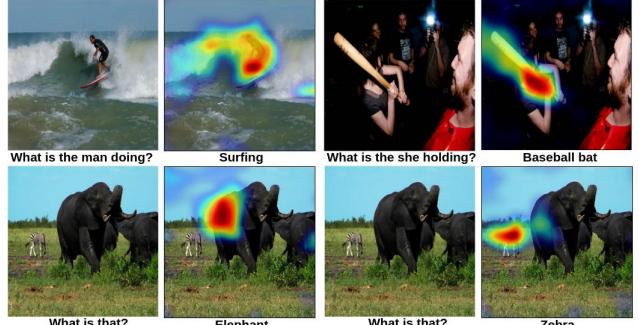
Visualizing ResNet-based VQA model with attention. Lu *et al.* [31] use a 200 layer ResNet [16] to encode the image, and jointly learn a hierarchical attention mechanism based on parses of question and image. Fig. 6b shows Grad-CAM visualization for this network. To the best of our knowledge, we are the first to visualize decisions made by ResNet-based architectures.

8. Conclusion

In this work, we proposed a novel class-discriminative localization technique – Gradient-weighted Class Activation Mapping (Grad-CAM) – for making CNN-based models more transparent by producing visual explanations. Further, we combined our Grad-CAM localizations with existing high-resolution visualizations having poor localization ability to obtain high-resolution class-discriminative visualizations, Guided Grad-CAM. Extensive human studies with



(a) Visualizing baseline VQA model from [30]



(b) Visualizing ResNet-based Hierarchical co-attention VQA model from [31]

Figure 6: Qualitative Results for our VQA experiments: (a) Given the image on the left and the question “What color is the firehydrant?”, we visualize Grad-CAMs and Guided Grad-CAMs for the answers “red”, “yellow” and “yellow and red”. Grad-CAM localizations are highly interpretable and help explain the model’s predictions – for “red”, the model focuses on the bottom red part of the firehydrant; when forced to answer “yellow”, it looks at the whole firehydrant! (b) Grad-CAM visualizations for the provided question and answer.

visualizations reveal that our localization-augmented visualizations can discriminate between classes more accurately and better reveal the trustworthiness of a classifier. Finally, we provide some quantitative and qualitative results on interpreting predictions from image classification, visual question answering and image captioning models including visualizations from very deep architectures such as ResNets and their variants. We believe that a true AI system should not only be intelligent, but also be able to reason about its beliefs and actions for humans to trust it. Future work includes explaining the decisions made by deep networks in domains such as reinforcement learning and natural language processing.

References

- [1] A. Agrawal, D. Batra, and D. Parikh. Analyzing the Behavior of Visual Question Answering Models. In *EMNLP*, 2016. 2
- [2] H. Agrawal, C. S. Mathialagan, Y. Goyal, N. Chavali, P. Banik, A. Mohapatra, A. Osman, and D. Batra. CloudCV: Large Scale Distributed Computer Vision as a Cloud Service. In *Mobile Cloud Visual Media Computing*, pages 265–290. Springer, 2015. 1
- [3] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh. VQA: Visual Question Answering. In *ICCV*, 2015. 2, 7, 8
- [4] L. Bazzani, A. Bergamo, D. Anguelov, and L. Torresani. Self-taught object localization with deep networks. In *WACV*, 2016. 4
- [5] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013. 5
- [6] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick. Microsoft COCO captions: Data Collection and Evaluation Server. *arXiv preprint arXiv:1504.00325*, 2015. 2, 7
- [7] R. G. Cinbis, J. Verbeek, and C. Schmid. Weakly supervised object localization with multi-fold multiple instance learning. *IEEE transactions on pattern analysis and machine intelligence*, 2016. 3
- [8] A. Das, H. Agrawal, C. L. Zitnick, D. Parikh, and D. Batra. Human Attention in Visual Question Answering: Do Humans and Deep Networks Look at the Same Regions? In *EMNLP*, 2016. 8
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009. 2, 5, 7, 15
- [10] A. Dosovitskiy and T. Brox. Inverting Convolutional Networks with Convolutional Networks. In *CVPR*, 2015. 3
- [11] D. Erhan, Y. Bengio, A. Courville, and P. Vincent. Visualizing Higher-layer Features of a Deep Network. *University of Montreal*, 1341, 2009. 3
- [12] H. Fang, S. Gupta, F. Iandola, R. K. Srivastava, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. C. Platt, et al. From Captions to Visual Concepts and Back. In *CVPR*, 2015. 2
- [13] C. Gan, N. Wang, Y. Yang, D.-Y. Yeung, and A. G. Hauptmann. Devnet: A deep event network for multimedia event detection and evidence recounting. In *CVPR*, 2015. 3
- [14] H. Gao, J. Mao, J. Zhou, Z. Huang, L. Wang, and W. Xu. Are You Talking to a Machine? Dataset and Methods for Multilingual Image Question Answering. In *NIPS*, 2015. 2, 7, 8
- [15] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In *CVPR*, 2014. 2
- [16] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 2, 3, 8, 19
- [17] D. Hoiem, Y. Chodpathumwan, and Q. Dai. Diagnosing Error in Object Detectors. In *ECCV*, 2012. 2
- [18] P. Jackson. *Introduction to Expert Systems*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 3rd edition, 1998. 2
- [19] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional Architecture for Fast Feature Embedding. In *ACM MM*, 2014. 5
- [20] E. Johns, O. Mac Aodha, and G. J. Brostow. Becoming the Expert - Interactive Multi-Class Machine Teaching. In *CVPR*, 2015. 2
- [21] J. Johnson, A. Karpathy, and L. Fei-Fei. DenseCap: Fully Convolutional Localization Networks for Dense Captioning. In *CVPR*, 2016. 2, 7, 8
- [22] A. Karpathy. What I learned from competing against a ConvNet on ImageNet. <http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/>, 2014. 2
- [23] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015. 7, 8
- [24] A. Kolesnikov and C. H. Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *ECCV*, 2016. 5, 10, 17, 18
- [25] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 2, 4, 12
- [26] M. Lin, Q. Chen, and S. Yan. Network in network. In *ICLR*, 2014. 3
- [27] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 11
- [28] Z. C. Lipton. The Mythos of Model Interpretability. *ArXiv e-prints*, June 2016. 2, 3
- [29] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 2
- [30] J. Lu, X. Lin, D. Batra, and D. Parikh. Deeper LSTM and normalized CNN Visual Question Answering model. https://github.com/VT-vision-lab/VQA_LSTM_CNN, 2015. 8, 11
- [31] J. Lu, J. Yang, D. Batra, and D. Parikh. Hierarchical question-image co-attention for visual question answering. In *NIPS*, 2016. 8
- [32] A. Mahendran and A. Vedaldi. Visualizing deep convolutional neural networks using natural pre-images. *International Journal of Computer Vision*, pages 1–23, 2016. 3, 5
- [33] M. Malinowski, M. Rohrbach, and M. Fritz. Ask your neurons: A neural-based approach to answering questions about images. In *ICCV*, 2015. 2, 7, 8
- [34] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *CVPR*, 2014. 3, 4
- [35] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Is object localization for free? – weakly-supervised learning with convolutional neural networks. In *CVPR*, 2015. 3
- [36] P. O. Pinheiro and R. Collobert. From image-level to pixel-level labeling with convolutional networks. In *CVPR*, 2015. 3
- [37] M. Ren, R. Kiros, and R. Zemel. Exploring models and data for image question answering. In *NIPS*, 2015. 2, 7, 8
- [38] M. T. Ribeiro, S. Singh, and C. Guestrin. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *SIGKDD*, 2016. 3, 7
- [39] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016. 2
- [40] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *CoRR*, abs/1312.6034, 2013. 3, 5
- [41] K. Simonyan and A. Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *ICLR*, 2015. 2, 4, 5, 11, 12, 16
- [42] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. A. Riedmiller. Striving for Simplicity: The All Convolutional Net. *CoRR*, abs/1412.6806, 2014. 1, 2, 3, 17
- [43] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015. 2, 7
- [44] C. Vondrick, A. Khosla, T. Malisiewicz, and A. Torralba. HOGgles: Visualizing Object Detection Features. *ICCV*, 2013. 7
- [45] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014. 1, 2, 3, 4, 7, 17
- [46] J. Zhang, Z. Lin, J. Brandt, X. Shen, and S. Sclaroff. Top-down Neural Attention by Excitation Backprop. In *ECCV*, 2016. 2, 4, 5, 6, 19
- [47] B. Zhou, A. Khosla, L. A., A. Oliva, and A. Torralba. Learning Deep Features for Discriminative Localization. In *CVPR*, 2016. 1, 2, 3, 4, 5, 19

Appendix

In this supplementary document, we provide

- Section **A**: Derivation to show that Grad-CAM is a generalization to CAM for *any* CNN-based architecture and hence doesn't require any architectural change or retraining.
- Section **B**: Qualitative results showing Grad-CAM and Guided Grad-CAM visualizations for image classification, image captioning, and visual question answering (VQA). For image captioning and VQA, our visualizations (Grad-CAM, and Guided Grad-CAM) expose the somewhat surprising insight that common CNN+LSTM models can often be good at localizing discriminative input image regions despite not being trained on grounded image-text pairs.
- Section **C**: Ablation studies to explore and validate our design choices for computing Grad-CAM visualizations.
- Section **D**: Weakly-supervised segmentation results on PASCAL VOC 2012 by using weak-localization cues from Grad-CAM as a seed for SEC [24].
- Section **E**: More details on the pointing game setup.
- Section **F**: Comparison to existing visualization techniques, CAM and c-MWP on PASCAL and COCO, where we find that our visualizations are superior, while being faster to compute and at the same time being possible to visualize a *wide variety of CNN-based models, including but not limited to, CNNs with fully-connected layers, CNNs stacked with Recurrent Neural Networks (RNNs), ResNets etc..*
- Section **G**: Analysis of Grad-CAM visualizations for 200-layer Residual Network.

A. Grad-CAM as generalization of CAM

In this section we formally prove that Grad-CAM is a generalization of CAM, as mentioned in Section 3 in the main paper. Recall that the CAM architecture consists of fully-covolutional CNNs, followed by global average pooling, and linear classification layer with softmax.

Let the final convolutional layer produce K feature maps A^k , with each element indexed by i, j . So A_{ij}^k refers to the activation at location (i, j) of the feature map A^k .

CAM computes a global average pooling (GAP) on A_{ij}^k . Let us define F^k to be the global average pooled output, So,

$$F^k = \frac{1}{Z} \sum_i \sum_j A_{ij}^k \quad (5)$$

CAM computes the final scores by,

$$Y^c = \sum_k w_k^c \cdot F^k \quad (6)$$

where w_k^c is the weight connecting the k^{th} feature map with the c^{th} class.

Taking the gradient of the score for class c (Y^c) with respect to the feature map F^k we get,

$$\text{(From Chain Rule)} \frac{\partial Y^c}{\partial F^k} = \frac{\frac{\partial Y^c}{\partial A_{ij}^k}}{\frac{\partial A_{ij}^k}{\partial F^k}} \quad (7)$$

Taking partial derivative of (5) w.r.t. A_{ij}^k , we can see that $\frac{\partial F^k}{\partial A_{ij}^k} = \frac{1}{Z}$. Substituting this in (7), we get,

$$\frac{\partial Y^c}{\partial F^k} = \frac{\partial Y^c}{\partial A_{ij}^k} \cdot Z \quad (8)$$

From (6) we get that, $\frac{\partial Y^c}{\partial F^k} = w_k^c$. Hence,

$$w_k^c = Z \cdot \frac{\partial Y^c}{\partial A_{ij}^k} \quad (9)$$

Now, we can sum both sides of this expression in (5) over all pixels (i, j) to get:

$$\sum_i \sum_j w_k^c = \sum_i \sum_j Z \cdot \frac{\partial Y^c}{\partial A_{ij}^k}, \quad \text{which can be rewritten as} \quad (10)$$

$$Zw_k^c = Z \sum_i \sum_j \frac{\partial Y^c}{\partial A_{ij}^k} \quad (\text{Since } Z \text{ and } w_k^c \text{ do not depend on } (i, j)) \quad (11)$$

Note that Z is the number of pixels in the feature map (or $Z = \sum_i \sum_j 1$). Thus, we can re-order terms and see that:

$$w_k^c = \sum_i \sum_j \frac{\partial Y^c}{\partial A_{ij}^k} \quad (12)$$

(13)

We can see that up to a proportionality constant ($1/Z$) that is normalized out during visualization, the expression for w_k^c is identical to α_k^c used by Grad-CAM (as described in the main paper).

Thus Grad-CAM is a generalization of CAM to arbitrary CNN-based architectures, while maintaining the computational efficiency of CAM.

B. Experimental Results

In this section we provide more qualitative results for Grad-CAM and Guided Grad-CAM applied to the task of image classification, image captioning and VQA.

B.1. Image Classification

We use Grad-CAM and Guided Grad-CAM to visualize the regions of the image that provide support for a particular prediction. The results reported in Fig. A1 correspond to the VGG-16 [41] network trained on ImageNet.

Fig. A1 shows randomly sampled examples from COCO [27] validation set. COCO images typically have multiple objects per image and Grad-CAM visualizations show precise localization to support the model’s prediction.

Guided Grad-CAM can even localize tiny objects. For example our approach correctly localizes the predicted class “torch” (Fig. A1.a) inspite of its size and odd location in the image. Our method is also class-discriminative – it places attention *only* on the “toilet seat” even when a popular ImageNet category “dog” exists in the image (Fig. A1.e).

We also visualized Grad-CAM, Guided Backpropagation (GB), Deconvolution (DC), GB + Grad-CAM (Guided Grad-CAM), DC + Grad-CAM (Deconvolution Grad-CAM) for images from the ILSVRC13 detection val set that have at least 2 unique object categories each. The visualizations for the mentioned class can be found in the following links.

“computer keyboard, keypad” class: <http://i.imgur.com/QMhsRzf.jpg>

“sunglasses, dark glasses, shades” class: <http://i.imgur.com/a1C7DGh.jpg>

B.2. Image Captioning

We use the publicly available Neuraltalk2 code and model⁵ for our image captioning experiments. The model uses VGG-16 to encode the image. The image representation is passed as input at the first time step to an LSTM that generates a caption for the image. The model is trained end-to-end along with CNN finetuning using the COCO [27] Captioning dataset. We feedforward the image to the image captioning model to obtain a caption. We use Grad-CAM to get a coarse localization and combine it with Guided Backpropagation to get a high-resolution visualization that highlights regions in the image that provide support for the generated caption.

B.3. Visual Question Answering (VQA)

We use Grad-CAM and Guided Grad-CAM to explain why a publicly available VQA model [30] answered what it answered. The VQA model by Lu *et al.* uses a standard CNN followed by a fully connected layer to transform the image to 1024-dim to match the LSTM embeddings of the question. Then the transformed image and LSTM embeddings are pointwise multiplied

⁵<https://github.com/karpathy/neuraltalk2>

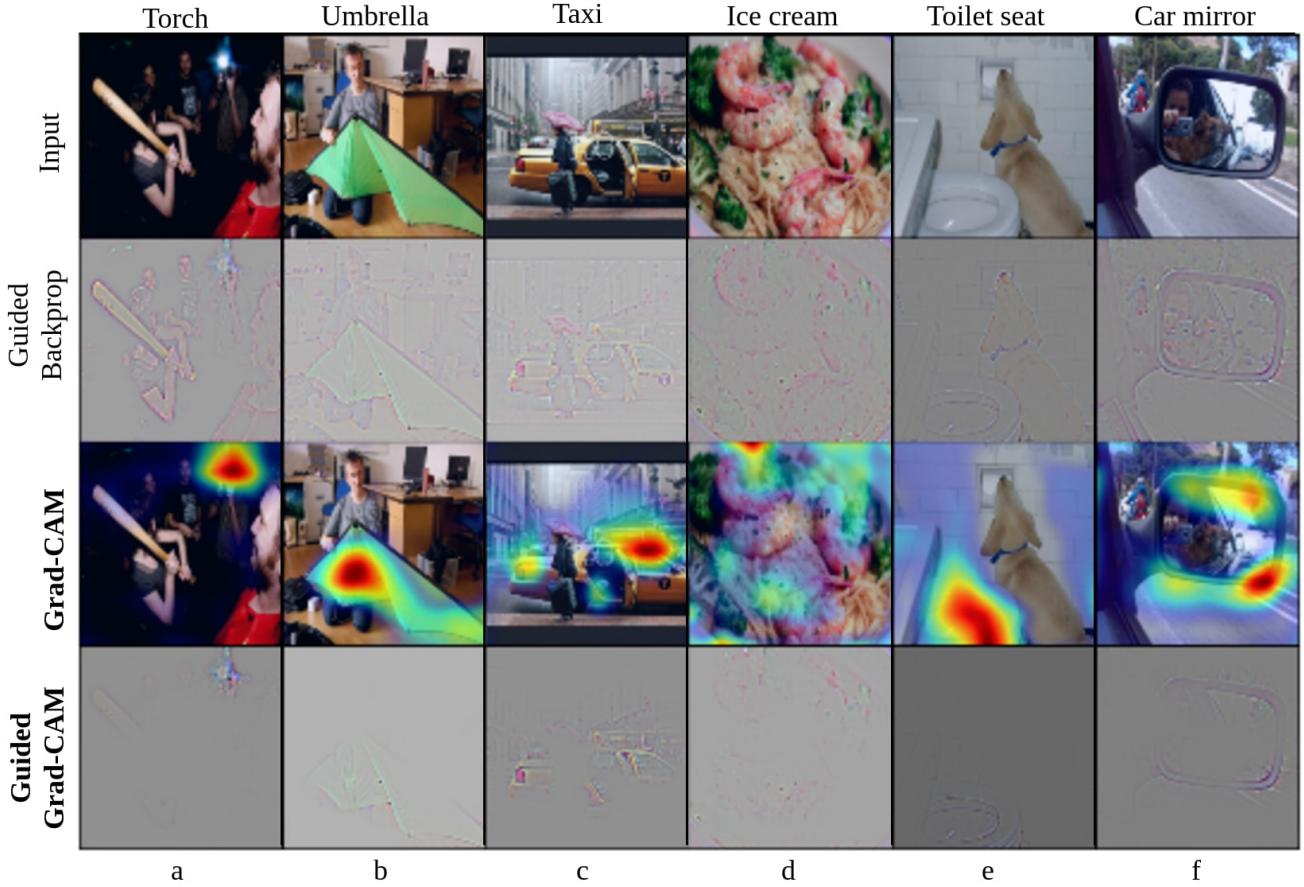


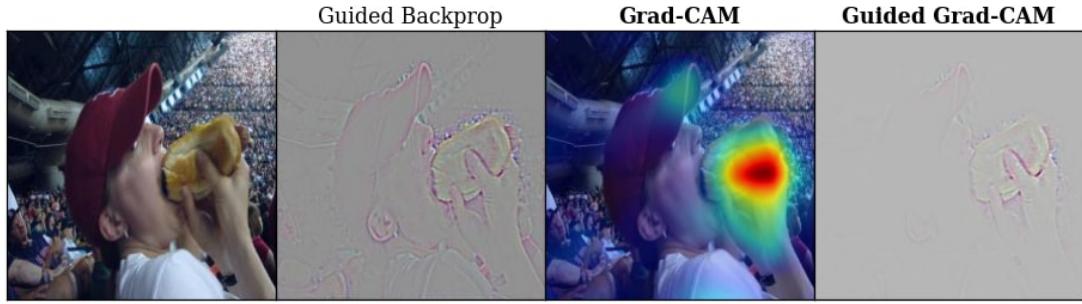
Figure A1: Visualizations for randomly sampled images from the COCO validation dataset. Predicted classes are mentioned at the top of each column.

to get a combined representation of the image and question and a multi-layer perceptron is trained on top to predict one among 1000 answers. We show visualizations for the VQA model trained with 3 different CNNs - AlexNet [25], VGG-16 and VGG-19 [41]. Even though the CNNs were not finetuned for the task of VQA, it is interesting to see how our approach can serve as a tool to understand these networks better by providing a localized high-resolution visualization of the regions the model is looking at. Note that these networks were trained with no explicit attention mechanism enforced.

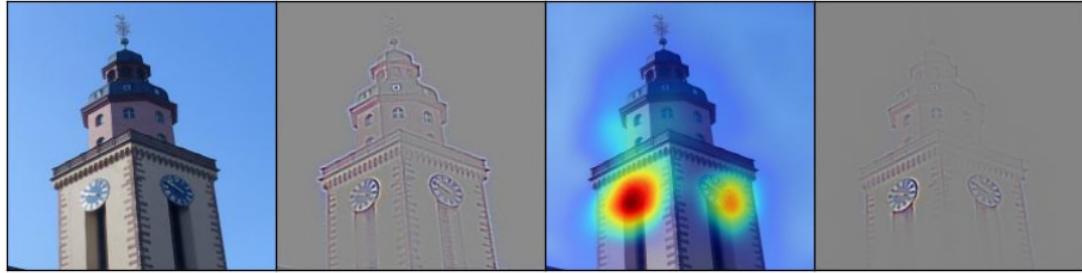
Notice in the first row of Fig. A3, for the question, “*Is the person riding the waves?*”, the VQA model with AlexNet and VGG-16 answered “No”, as they concentrated on the person mainly, and not the waves. On the other hand, VGG-19 correctly answered “Yes”, and it looked at the regions around the man in order to answer the question. In the second row, for the question, “*What is the person hitting?*”, the VQA model trained with AlexNet answered “Tennis ball” just based on context without looking at the ball. Such a model might be risky when employed in real-life scenarios. It is difficult to determine the trustworthiness of a model just based on the predicted answer. Our visualizations provide an accurate way to explain the model’s predictions and help in determining which model to trust, without making any architectural changes or sacrificing accuracy. Notice in the last row of Fig. A3, for the question, “*Is this a whole orange?*”, the model looks for regions around the orange to answer “No”.

C. Ablation studies

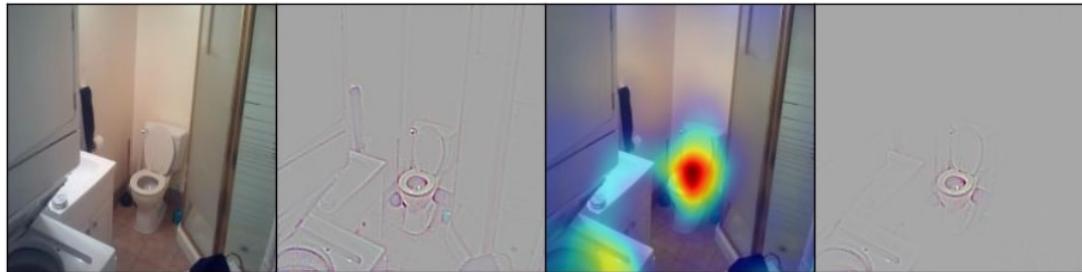
In this section we provide details of the ablation studies we performed.



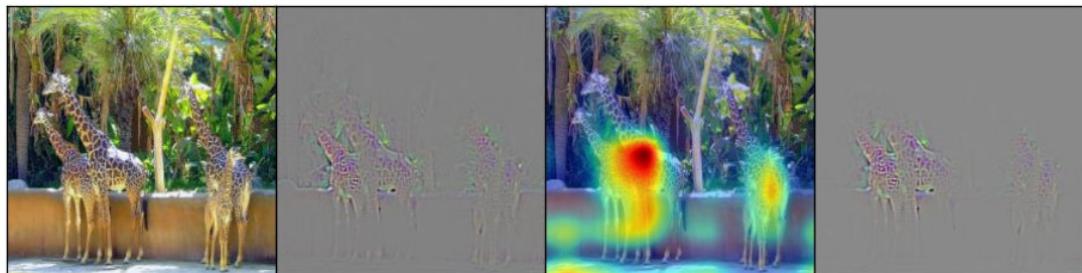
A man is holding a hot dog in his hand



A large clock tower with a clock on the top of it



A bathroom with a toilet and a sink



Two giraffes standing in a zoo enclosure with a fence



A stop sign on a street corner with a sign on it

Figure A2: Guided Backpropagation, Grad-CAM and Guided Grad-CAM visualizations for the captions produced by the NeuralTalk2 image captioning model.

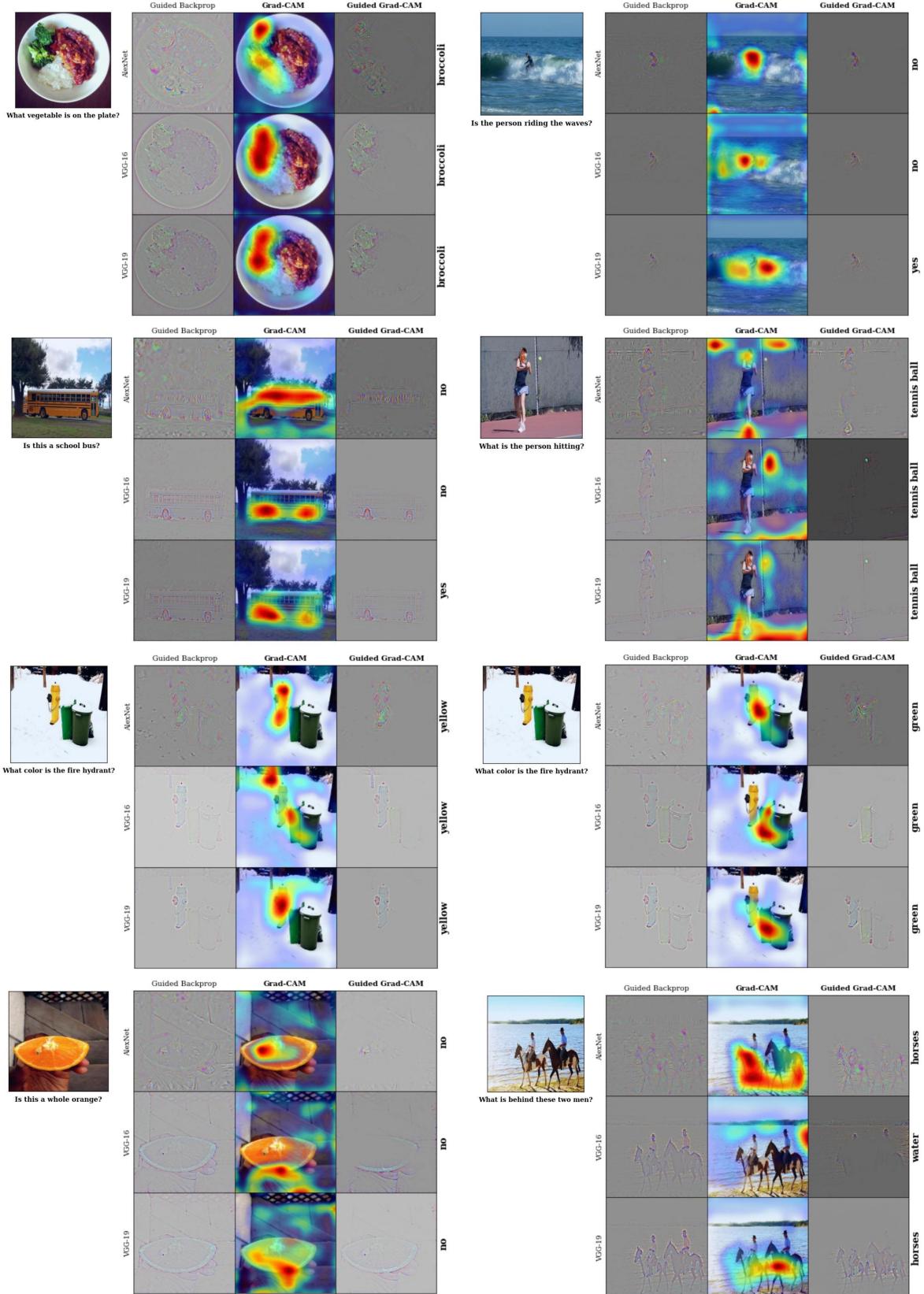


Figure A3: Guided Backpropagation, Grad-CAM and Guided Grad-CAM visualizations for the answers from a VQA model. For each image-question pair, we show visualizations for AlexNet, VGG-16 and VGG-19. Notice how the attention changes in row 3, as we change the answer from *Yellow* to *Green*.

C.1. Varying mask size for occlusion

Fig. 1 (e,k) of main paper show the results of occlusion sensitivity for the “cat” and “dog” class. We compute this occlusion map by repeatedly masking regions of the image and forward propagate each masked image. At each location of the occlusion map we store the difference in the original score for the particular class and the score obtained after forward propagating the masked image. Our choices for mask sizes include (10×10 , 15×15 , 25×25 , 35×35 , 45×45 , and 90×90). We zero-pad the images so that the resultant occlusion map is of the same size as the original image. The resultant occlusion maps can be found in Fig. A4. Note that blue regions correspond to a decrease in score for a particular class (“tiger cat” in the case of Fig. A4) when the region around that pixel is occluded. Hence it serves as an evidence for the class. Whereas the red regions correspond to an increase in score as the region around that pixel is occluded. Hence these regions might indicate existence of other confusing classes. We observe that 35×35 is a good trade-off between sharp results and a smooth appearance.

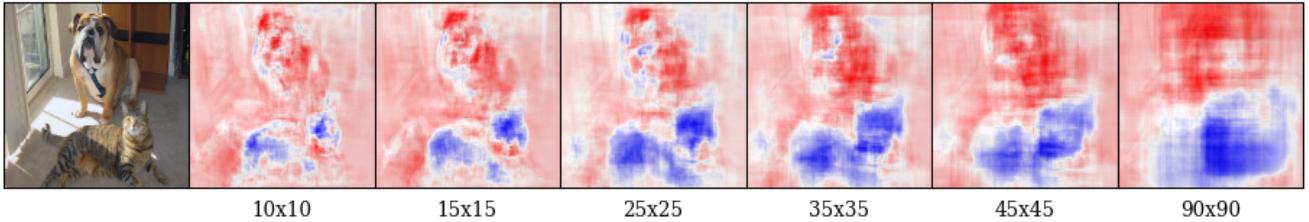


Figure A4: Occlusion maps with different mask sizes for the “tiger cat” category.

C.2. Guided Grad-CAM on different layers

We show results of applying Grad-CAM for the “Tiger-cat” category on different convolutional layers in AlexNet and VGG-16 CNN. As expected, the results from Fig. A5 show that localization becomes progressively worse as we move to shallower convolutional layers. This is because the later convolutional layers capture high-level semantic information and at the same time retain spatial information, while the shallower layers have smaller receptive fields and only concentrate on local features that are important for the next layers.

C.3. Design choices

Method	Top-1 error
Grad-CAM	59.65
Grad-CAM without ReLU in Eq.1	74.98
Grad-CAM with Absolute gradients	58.19
Grad-CAM with GMP gradients	59.96
Grad-CAM with Deconv ReLU	83.95
Grad-CAM with Guided ReLU	59.14

Table A1: Localization results on ILSVRC-15 val for the ablation studies. Note that the visualizations were created for single-crop, compared to the 10-crop evaluation reported in the main paper.

We evaluate design choices via top-1 localization error on the ILSVRC15 val set [9].

C.3.1 Importance of ReLU in Eq. 1 in main paper

Removing ReLU (Eq. 1 in main paper) increases error by 15.3%. See Table. A1. Negative values in Grad-CAM indicate confusion between multiple occurring classes. Thus, localization improves when we suppress them (see Fig. A7).

C.3.2 Absolute value of each derivative in Eq. 2 in main paper

Taking the absolute value of each derivative in Eq. 2 in main paper decreases the error by 1.5% (see Table. A1). But qualitatively maps look a bit worse (see Fig. A7), and this evaluation does not fully capture class discriminability (most ImageNet images have only 1 class).

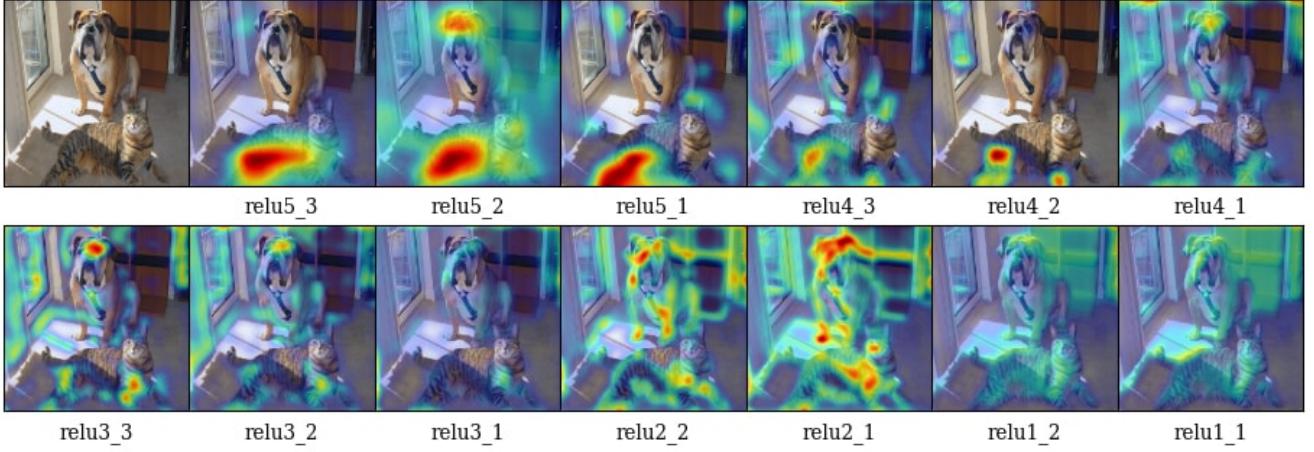


Figure A5: Grad-CAM at different convolutional layers for the ‘tiger cat’ class. This figure analyzes how localizations change qualitatively as we perform Grad-CAM with respect to different feature maps in a CNN (VGG16 [41]). We find that the best looking visualizations are often obtained after the deepest convolutional layer in the network, and localizations get progressively worse at shallower layers. This is consistent with our intuition described in Section 3 of main paper.

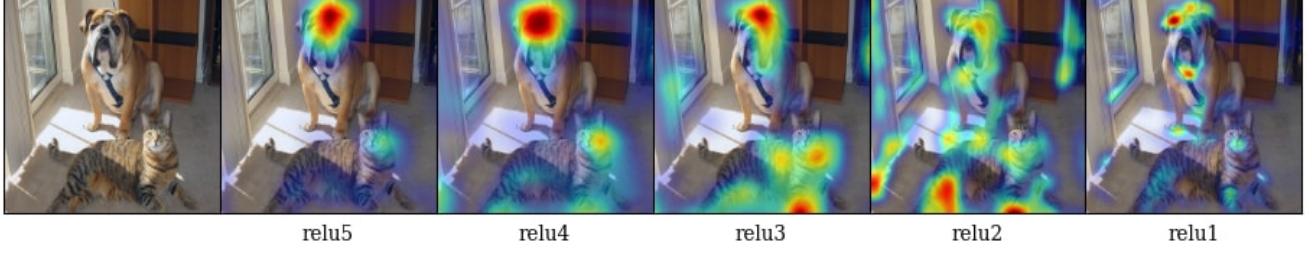


Figure A6: Grad-CAM localizations for “tiger cat” category for different rectified convolutional layer feature maps for AlexNet.

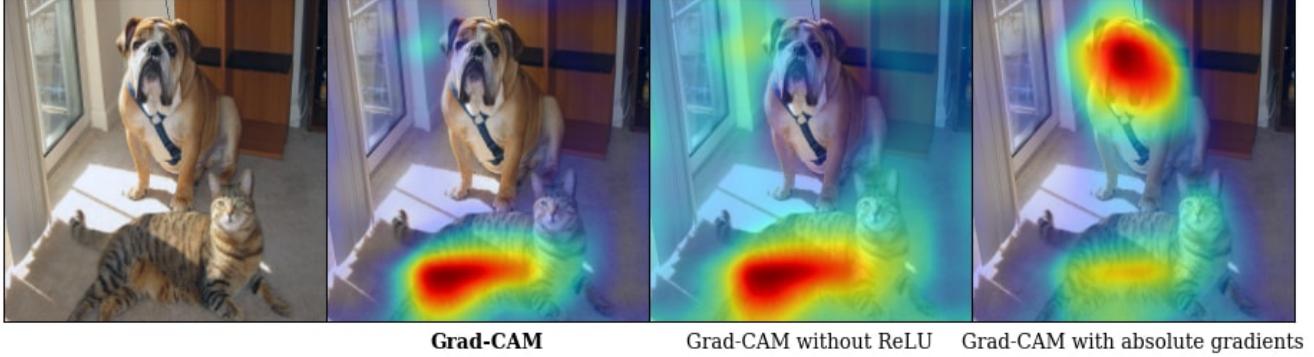


Figure A7: Grad-CAM visualizations for “tiger cat” category stating the importance of ReLU and effect of using absolute gradients in Eq. 1 of main paper.

C.3.3 Global Average Pooling vs. Global Max Pooling

Instead of Global Average Pooling (GAP) the incoming gradients to the convolutional layer, we tried Global Max Pooling (GMP) them. We observe that using GMP lowers the localization ability of our Grad-CAM technique. An example can be found in Fig. A8 below. This observation is also summarized in Table. A1. This may be due to the fact that *max* is statistically less robust to noise compared to the *averaged* gradient.

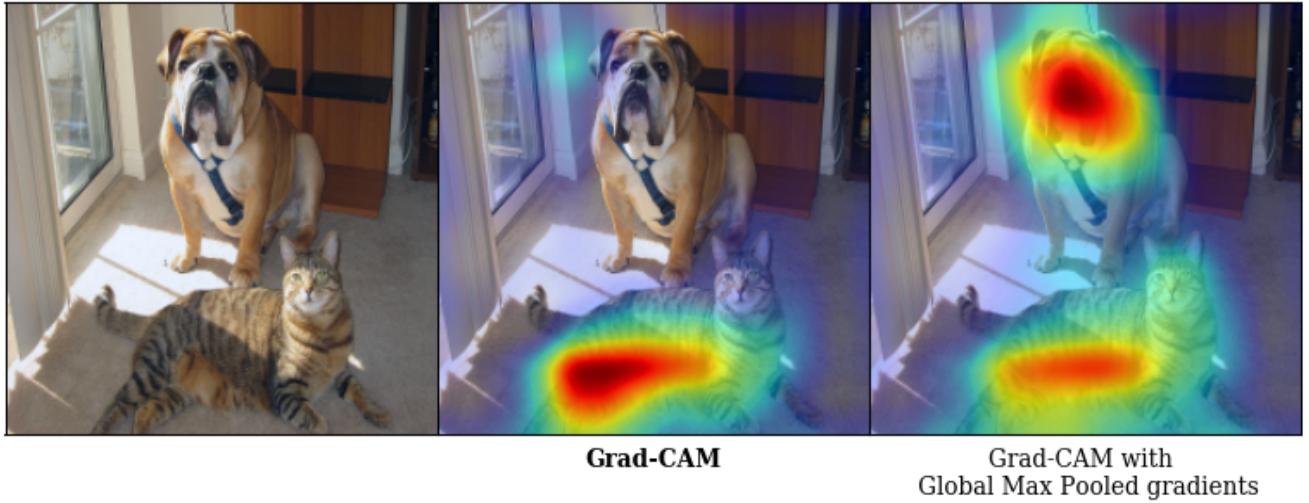


Figure A8: Grad-CAM visualizations for “tiger cat” category with Global Average Pooling and Global Max Pooling.

C.3.4 Effect of different ReLU on Grad-CAM

We experiment with different modifications to the backward pass of ReLU, namely, using Guided-ReLU [42] and Deconv-ReLU [45].

Effect of Guided-ReLU:

Springenberg *et al.* [42] introduced Guided Backprop, where they modified the backward pass of ReLU to pass only positive gradients to regions with positive activations. Applying this change to the computation of our Grad-CAM maps introduces a drop in the class-discriminative ability of Grad-CAM as can be seen in Fig. A9, but it gives a slight improvement in the localization ability on ILSVRC’14 localization challenge (see Table. A1).

Effect of Deconv-ReLU:

Zeiler and Fergus [45] in their Deconvolution work introduced a slight modification to the backward pass of ReLU, to pass only the positive gradients from higher layers. Applying this modification to the computation of our Grad-CAM gives worse results as shown in Fig. A9.

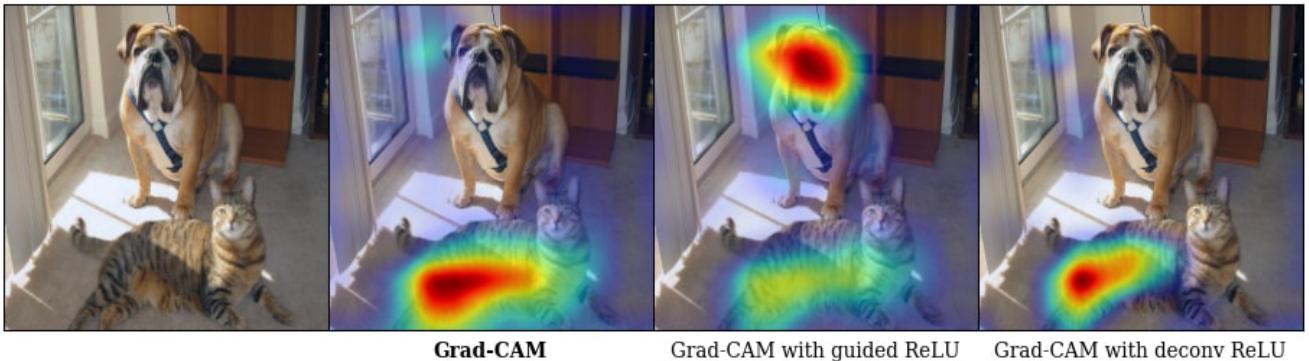


Figure A9: Grad-CAM visualizations for “tiger cat” category for different modifications to the ReLU backward pass. The best results are obtained when we use the actual gradients during the computation of Grad-CAM.

D. Weakly-supervised segmentation

In recent work Kolesnikov *et al.* [24] introduced a new loss function for training weakly-supervised image segmentation models. Their loss function is based on three principles: 1. to seed with weak localization cues, 2. to expand object seeds to

regions of reasonable size, 3. to constrain segmentations to object boundaries. They showed that their proposed loss function leads to better segmentation.

They showed that their algorithm is very sensitive to seed loss, without which the segmentation network fails to localize the objects correctly [24]. In their work, they used CAM for weakly localizing foreground classes. We replaced CAM with Grad-CAM and show results in Fig. A10. The last row shows 2 failure cases. In the bottom left image, the clothes of the 2 person weren't highlighted correctly. This could be because the most discriminative parts are their faces, and hence Grad-CAM maps only highlights those. This results in a segmentation that only highlights the faces of the 2 people. In the bottom right image, the bicycles, being extremely thin aren't highlighted. This could be because the resolution of the Grad-CAM maps are low (14×14) which makes it difficult to capture thin areas.

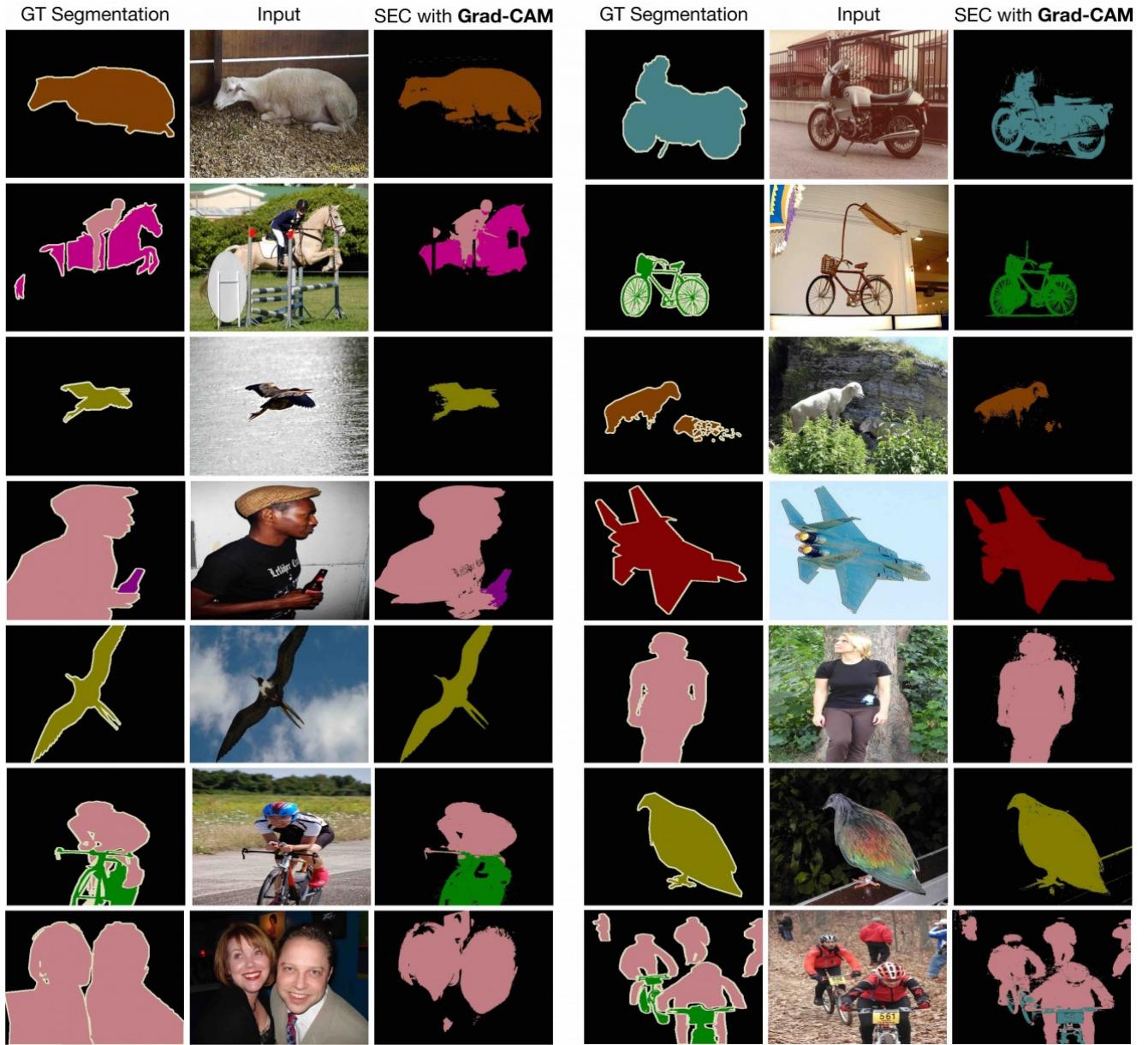


Figure A10: PASCAL VOC 2012 Segmentation results with Grad-CAM as seed for SEC [24].

E. More details of Pointing Game

In [46], the pointing game was setup to evaluate the discriminativeness of different attention maps for localizing ground-truth categories. In a sense, this evaluates the precision of a visualization, *i.e.* how often does the attention map intersect the segmentation map of the ground-truth category. This does not evaluate how often the visualization technique produces maps which do not correspond to the category of interest. For example this evaluation does not penalize the visualization in Fig. A12 top-left, for highlighting a zebra when visualizing the bird category.

Hence we propose a modification to the pointing game to evaluate visualizations of the top-5 predicted category. In this case the visualizations are given an additional option to reject any of the top-5 predictions from the CNN classifiers. For each of the two visualizations, Grad-CAM and c-MWP, we choose a threshold on the max value of the visualization, that can be used to determine if the category being visualized exists in the image.

We compute the maps for the top-5 categories, and based on the maximum value in the map, we try to classify if the map is of the GT label or a category that is absent in the image. As mentioned in Section 4.2 of the main paper, we find that our approach Grad-CAM outperforms c-MWP by a significant margin (70.58% vs 60.30%). Fig. A12 shows the maps computed for the top-5 categories using c-MWP and Grad-CAM.

F. Qualitative comparison to Excitation Backprop (c-MWP) and CAM

In this section we provide more qualitative results comparing Grad-CAM with CAM [47] and c-MWP [46].

F.1. PASCAL

We compare Grad-CAM, CAM and c-MWP visualizations from ImageNet trained VGG-16 models finetuned on PASCAL VOC 2012 dataset. While Grad-CAM and c-MWP visualizations can be directly obtained from existing models, CAM requires an architectural change, and requires re-training, which leads to loss in accuracy. Also, unlike Grad-CAM, c-MWP and CAM can only be applied for image classification networks. Visualizations for the ground-truth categories can be found in Fig. A11.

F.2. COCO

We compare Grad-CAM and c-MWP visualizations from ImageNet trained VGG-16 models finetuned on COCO dataset. Visualizations for the top-5 predicted categories can be found in Fig. A12. It can be seen that c-MWP highlights arbitrary regions for predicted but non-existent categories, unlike Grad-CAM which seem much more reasonable. We quantitatively evaluate this through the pointing experiment.

G. Analyzing Residual Networks

In this section, we perform Grad-CAM on Residual Networks (ResNets). In particular, we analyze the 200-layer architecture trained on ImageNet⁶.

Current ResNets [16] typically consist of residual blocks. One set of blocks use identity skip connections (shortcut connections between two layers having identical output dimensions). These sets of residual blocks are interspersed with downsampling modules that alter dimensions of propagating signal. As can be seen in Fig. A13 our visualizations applied on the last convolutional layer can correctly localize the cat and the dog. Grad-CAM can also visualize the cat and dog correctly in the residual blocks of the last set. However, as we go towards earlier sets of residual blocks with different spatial resolution, we see that Grad-CAM fails to localize the category of interest (see last row of Fig. A13). We observe similar trends for other ResNet architectures (18 and 50-layer ResNets).

⁶We use the 200-layer ResNet architecture from <https://github.com/facebook/fb.resnet.torch>.

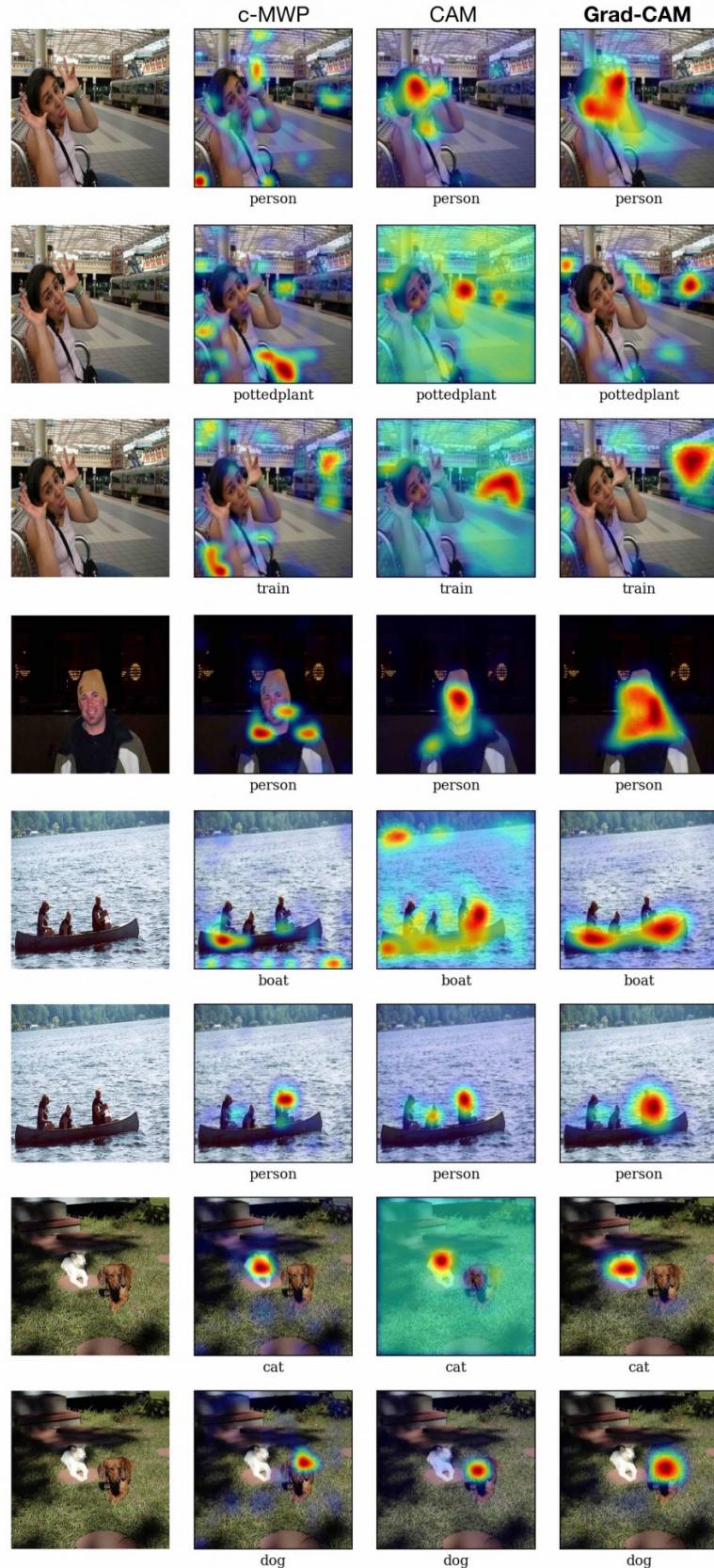


Figure A11: Visualizations for ground-truth categories (shown below each image) for images sampled from the PASCAL validation set.

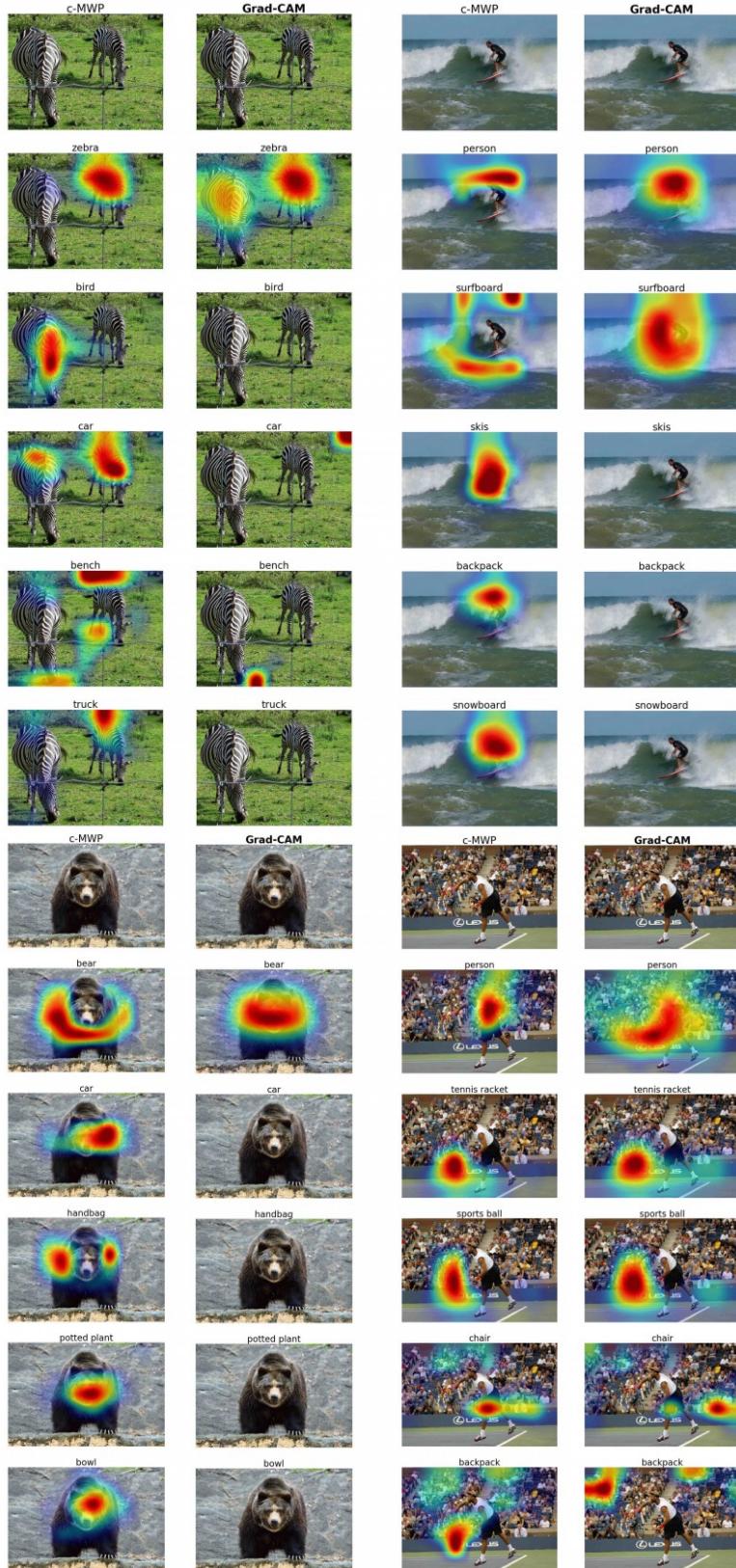
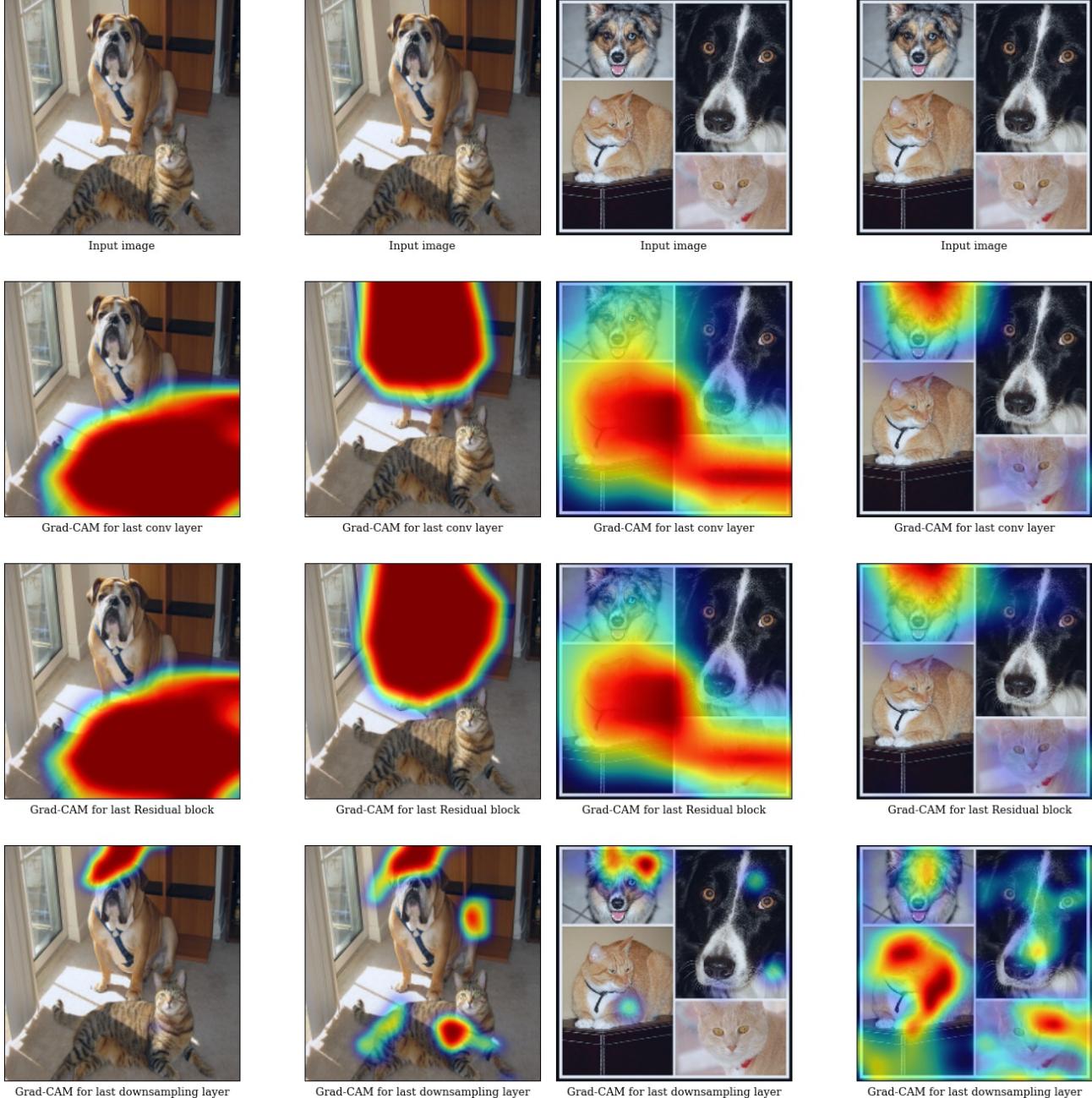


Figure A12: c-MWP and Grad-CAM visualizations for the top-5 predicted categories (shown above each image) for images sampled from the COCO validation set.



(a) Grad-CAM visualizations for the ResNet-200 layer architecture for 'tiger cat' (left) and 'boxer' (right) category.

(b) Grad-CAM visualizations for the ResNet-200 layer architecture for 'tabby cat' (left) and 'boxer' (right) category.

Figure A13: We observe that the discriminative ability of Grad-CAM significantly reduces as we encounter the downsampling layer.