

Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization

Ramprasaath R. Selvaraju · Michael Cogswell · Abhishek Das · Ramakrishna Vedantam · Devi Parikh · Dhruv Batra

Abstract We propose a technique for producing ‘visual explanations’ for decisions from a large class of Convolutional Neural Network (CNN)-based models, making them more transparent and explainable.

Our approach – Gradient-weighted Class Activation Mapping (Grad-CAM), uses the gradients of any target concept (say ‘dog’ in a classification network or a sequence of words in captioning network) flowing into the final convolutional layer to produce a coarse localization map highlighting the important regions in the image for predicting the concept.

Unlike previous approaches, Grad-CAM is applicable to a wide variety of CNN model-families: (1) CNNs with fully-connected layers (*e.g.* VGG), (2) CNNs used for structured outputs (*e.g.* captioning), (3) CNNs used in tasks with multi-modal inputs (*e.g.* visual question answering) or reinforcement learning, all *without architectural changes or re-training*. We combine Grad-CAM with existing fine-grained visualizations to create a high-resolution class-discriminative vi-

sualization, Guided Grad-CAM, and apply it to image classification, image captioning, and visual question answering (VQA) models, including ResNet-based architectures.

In the context of image classification models, our visualizations (a) lend insights into failure modes of these models (showing that seemingly unreasonable predictions have reasonable explanations), (b) outperform previous methods on the ILSVRC-15 weakly-supervised localization task, (c) are robust to adversarial perturbations, (d) are more faithful to the underlying model, and (e) help achieve model generalization by identifying dataset bias.

For image captioning and VQA, our visualizations show that even non-attention based models learn to localize discriminative regions of input image.

We devise a way to identify important neurons through Grad-CAM and combine it with neuron names [4] to provide textual explanations for model decisions. Finally, we design and conduct human studies to measure if Grad-CAM explanations help users establish appropriate trust in predictions from deep networks and show that Grad-CAM helps untrained users successfully discern a ‘stronger’ deep network from a ‘weaker’ one even when both make identical predictions. Our code is available at <https://github.com/ramprs/grad-cam/>, along with a demo on CloudCV [2]¹, and a video at youtu.be/COjUB9Izk6E.

Ramprasaath R. Selvaraju
Georgia Institute of Technology, Atlanta, GA, USA
E-mail: ramprs@gatech.edu

Michael Cogswell
Georgia Institute of Technology, Atlanta, GA, USA
E-mail: cogswell@gatech.edu

Abhishek Das
Georgia Institute of Technology, Atlanta, GA, USA
E-mail: abhshkdz@gatech.edu

Ramakrishna Vedantam
Georgia Institute of Technology, Atlanta, GA, USA
E-mail: vrama@gatech.edu

Devi Parikh
Georgia Institute of Technology, Atlanta, GA, USA
Facebook AI Research, Menlo Park, CA, USA
E-mail: parikh@gatech.edu

Dhruv Batra
Georgia Institute of Technology, Atlanta, GA, USA
Facebook AI Research, Menlo Park, CA, USA
E-mail: dbatra@gatech.edu

1 Introduction

Deep neural models based on Convolutional Neural Networks (CNNs) have enabled unprecedented breakthroughs in a variety of computer vision tasks, from image classification [33, 24], object detection [21], semantic segmentation [37] to image captioning [55, 7, 18, 29], visual question answering [3, 20, 42, 46] and more recently, visual dialog [11, 13, 12] and embodied question answering [10, 23]. While

¹ <http://gradcam.cloudcv.org>

these models enable superior performance, their lack of decomposability into *individually intuitive* components makes them hard to interpret [36]. Consequently, when today’s intelligent systems fail, they often fail spectacularly disgracefully without warning or explanation, leaving a user staring at an incoherent output, wondering why the system did what it did.

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 have the ability to explain *why they predict what they predict*. Broadly speaking, this transparency and ability to explain 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, 25], 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 [30] trained on sufficient data), the goal is to establish appropriate trust and confidence in users. Third, when AI is significantly stronger than humans (e.g. chess or Go [50]), the goal of explanations is in machine teaching [28] – i.e., a machine teaching a human about how to make better decisions.

There typically exists a trade-off between accuracy and simplicity or interpretability. Classical rule-based or expert systems [26] 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) [24] are over 200-layers deep and have shown state-of-the-art performance in several challenging tasks. Such complexity makes these models hard to interpret. As such, deep models are beginning to explore the spectrum between interpretability and accuracy.

Zhou *et al.* [59] recently proposed a technique called Class Activation Mapping (CAM) for identifying discriminative regions used by a restricted class of image classification CNNs which do not contain any fully-connected layers. In essence, this work trades off model complexity and performance for more transparency into the working of the model. In contrast, we make existing state-of-the-art deep models interpretable without altering their architecture, thus avoiding the interpretability vs. accuracy trade-off. Our approach is a generalization of CAM [59] and is applicable to a significantly broader range of CNN model families: (1) CNNs with fully-connected layers (e.g. VGG), (2) CNNs used for structured outputs (e.g. captioning), (3) CNNs used in tasks with

multi-modal inputs (e.g. VQA) or reinforcement learning, without requiring architectural changes or re-training.

What makes a good visual explanation? Consider image classification [14] – a ‘good’ visual explanation from the model for justifying any target category 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 Back-propagation [53] and Deconvolution [57] are high-resolution and highlight fine-grained details in the image, but are not class-discriminative (Fig. 1b and Fig. 1h are very similar).

In contrast, localization approaches like CAM or our proposed method Gradient-weighted Class Activation Mapping (Grad-CAM), are highly class-discriminative (the ‘cat’ explanation exclusively highlights the ‘cat’ regions but not ‘dog’ regions in Fig. 1c, and *vice versa* in Fig. 1i).

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 any decision of interest are visualized in high-resolution detail even if the image contains evidence for multiple possible concepts, 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 introduce Grad-CAM, a class-discriminative localization technique that generates visual explanations for *any* CNN-based network without requiring architectural changes or re-training. We evaluate Grad-CAM for localization (Sec. 4.1), and faithfulness to model (Sec. 5.3), where it outperforms baselines.
- (2) We apply Grad-CAM to existing top-performing classification, captioning (Sec. 8.1), and VQA (Sec. 8.2) models. For image classification, our visualizations lend insight into failures of current CNNs (Sec. 6.1), showing that seemingly unreasonable predictions have reasonable explanations. For captioning and VQA, our visualizations expose that common CNN + LSTM models are often surprisingly good at localizing discriminative image regions despite not being trained on grounded image-text pairs.
- (3) We show a proof-of-concept of how interpretable Grad-CAM visualizations help in diagnosing failure modes by uncovering biases in datasets. This is important not just for generalization, but also for fair and bias-free outcomes as more and more decisions are made by algorithms in society.
- (4) We present Grad-CAM visualizations for ResNets [24] applied to image classification and VQA (Sec. 8.2).

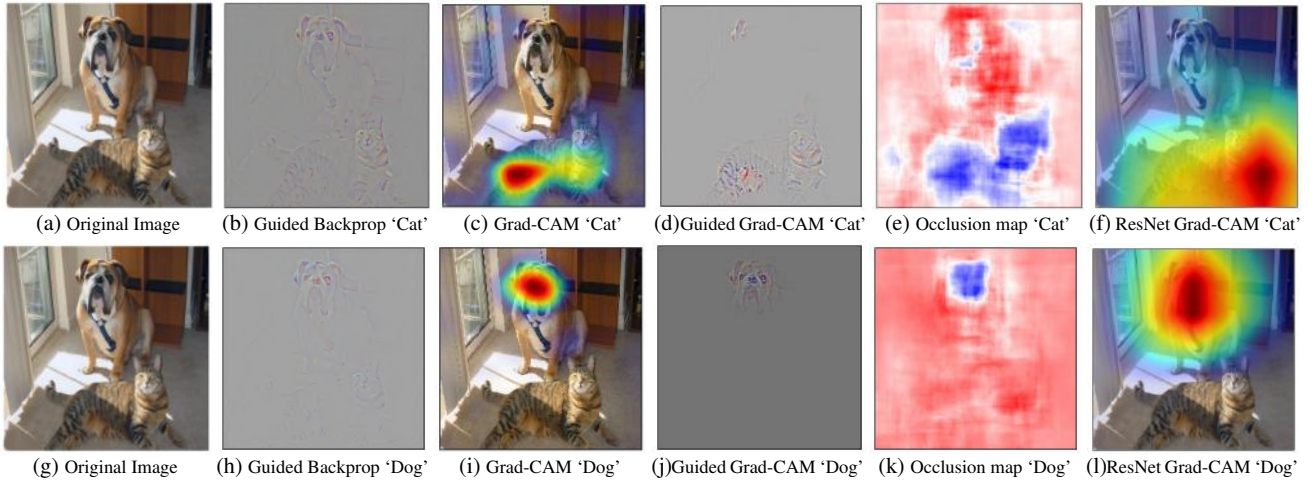


Fig. 1: (a) Original image with a cat and a dog. (b-f) Support for the cat category according to various visualizations for VGG-16 and ResNet. (b) Guided Backpropagation [53]: highlights all contributing features. (c, f) Grad-CAM (Ours): localizes class-discriminative regions, (d) Combining (b) and (c) gives Guided Grad-CAM, which gives high-resolution class-discriminative visualizations. Interestingly, the localizations achieved by our Grad-CAM technique, (c) are very similar to results from occlusion sensitivity (e), while being orders of magnitude cheaper to compute. (f, l) are Grad-CAM visualizations for ResNet-18 layer. Note that in (c, f, i, l), red regions corresponds to high score for class, while in (e, k), blue corresponds to evidence for the class. Figure best viewed in color.

(5) We use neuron importance from Grad-CAM and neuron names from [4] and obtain textual explanations for model decisions (Sec. 7).

(6) We conduct human studies (Sec. 5) that show Guided Grad-CAM explanations are class-discriminative and not only help humans establish trust, but also help untrained users successfully discern a ‘stronger’ network from a ‘weaker’ one, *even when both make identical predictions*.

Paper Organization: The rest of the paper is organized as follows. In section 3 we propose our approach Grad-CAM and Guided Grad-CAM. In sections 4 and 5 we evaluate the localization ability, class-discriminateness, trustworthiness and faithfulness of Grad-CAM. In section 6 we show certain use cases of Grad-CAM such as diagnosing image classification CNNs and identifying biases in datasets. In section 7 we provide a way to obtain textual explanations with Grad-CAM. In section 8 we show how Grad-CAM can be applied to vision and language models – image captioning and Visual Question Answering (VQA).

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 [51, 53, 57, 19] have visualized CNN predictions by highlighting ‘important’ pixels (*i.e.* change in intensities of these pixels have the most impact on the prediction score). Specifically, Simonyan *et al.* [51] visualize partial derivatives of predicted class scores w.r.t. pixel intensities, while Guided Backpropagation [53] and Deconvolution [57] make modifications to ‘raw’ gradients that result in qualitative improvements. These approaches are compared in [40]. Despite produc-

ing 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 [51, 16] or invert a latent representation [41, 15]. Although these can be high-resolution and class-discriminative, they are not specific to a single input image and visualize a model overall.

Assessing Model Trust. Motivated by notions of interpretability [36] and assessing trust in models [47], we evaluate Grad-CAM visualizations in a manner similar to [47] via human studies to show that they can be important tools for users to evaluate and place trust in automated systems.

Aligning Gradient-based Importances. Selvaraju *et al.* [48] proposed an approach that uses the gradient-based neuron importances introduced in our work, and maps it to class-specific domain knowledge from humans in order to learn classifiers for novel classes. In future work, Selvaraju *et al.* [49] proposed an approach to align gradient-based importances to human attention maps in order to ground vision and language models.

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 holistic image class labels only [8, 43, 44, 59].

Most relevant to our approach is the Class Activation Mapping (CAM) approach to localization [59]. This approach modifies image classification CNN architectures replacing fully-connected layers with convolutional layers and global average pooling [34], thus achieving class-specific feature maps. Others have investigated similar methods using global max pooling [44] and log-sum-exp pooling [45].

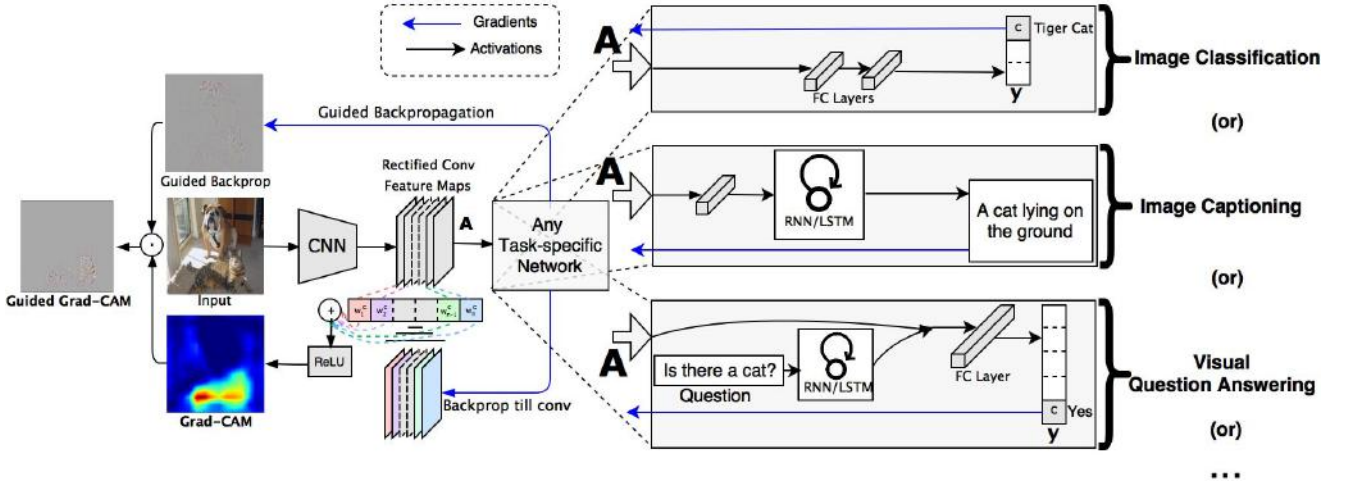


Fig. 2: Grad-CAM overview: Given an image and a class of interest (e.g., ‘tiger cat’ or any other type of differentiable output) as input, we forward propagate the image through the CNN part of the model and then through task-specific computations to obtain a raw score for the category. The gradients 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 maps of interest, which we combine to compute the coarse Grad-CAM localization (blue heatmap) which represents where the model has to look to make the particular decision. Finally, we pointwise multiply the heatmap with guided backpropagation to get Guided Grad-CAM visualizations which are both high-resolution and concept-specific.

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 \rightarrow global average pooling \rightarrow 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 any 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 off-the-shelf CNN-based architectures, including those for image captioning and visual question answering. For a fully-convolutional architecture, CAM is a special case of Grad-CAM.

Other methods approach localization by classifying perturbations of the input image. Zeiler and Fergus [57] 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 [5]. Oquab *et al.* [43] classify many patches containing a pixel then average these patch-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.* [58] 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 computationally more expensive than Grad-CAM and only works for image classification CNNs. Moreover, Grad-CAM outperforms c-MWP in quantitative and qualitative evaluations (see Sec. 4.1 and Sec. D).

3 Grad-CAM

A number of previous works have asserted that deeper representations in a CNN capture higher-level visual constructs [6, 41]. Furthermore, convolutional layers naturally retain spatial information which is lost in fully-connected layers, so we can expect the last convolutional layers to have the best compromise between high-level semantics and detailed spatial information. The neurons in these layers look for semantic class-specific information in the image (say object parts). Grad-CAM uses the gradient information flowing into the last convolutional layer of the CNN to assign importance values to each neuron for a particular decision of interest. Although our technique is fairly general in that it can be used to explain activations in any layer of a deep network, in this work, we focus on explaining output layer decisions only.

As shown in Fig. 2, in order to obtain the class-discriminative localization map Grad-CAM $L_{\text{Grad-CAM}}^c \in \mathbb{R}^{u \times v}$ of width u and height v for any class c , we first compute the gradient of the score for class c , y^c (before the softmax), with respect to feature map activations A^k of a convolutional layer, *i.e.* $\frac{\partial y^c}{\partial A^k}$. These gradients flowing back are global-average-pooled² over the width and height dimensions (indexed by i and j respectively) to obtain the neuron importance weights α_k^c :

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

During computation of α_k^c while backpropagating gradients with respect to activations, the exact computation amounts

² Empirically we found global-average-pooling to work better than global-max-pooling as can be found in the Appendix.

to successive matrix products of the weight matrices and the gradient with respect to activation functions till the final convolution layer that the gradients are being propagated to. Hence, 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 .

We perform a weighted combination of forward activation maps, and follow it by a ReLU to obtain,

$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left(\underbrace{\sum_k \alpha_k^c A^k}_{\text{linear combination}} \right) \quad (2)$$

Notice that this results in a coarse heatmap of the same size as the convolutional feature maps (14×14 in the case of last convolutional layers of VGG [52] and AlexNet [33] networks)³. We apply a ReLU to the linear combination of maps because 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 y^c . Negative pixels are likely to belong to other categories in the image. As expected, without this ReLU, localization maps sometimes highlight more than just the desired class and perform worse at localization. Figures 1c, 1f and 1i, 1l show Grad-CAM visualizations for ‘tiger cat’ and ‘boxer (dog)’ respectively. Ablation studies are available in Sec. B.

In general, y^c need not be the class score produced by an image classification CNN. It could be any differentiable activation including words from a caption or answer to a question.

3.1 Grad-CAM generalizes CAM

In this section, we discuss the connections between Grad-CAM and Class Activation Mapping (CAM) [59], and formally prove that Grad-CAM generalizes CAM for a wide variety of CNN-based architectures. Recall that CAM 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 softmax. Specifically, let the penultimate layer produce K feature maps, $A^k \in \mathbb{R}^{u \times v}$, 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 . 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 ,

$$Y^c = \sum_k \underbrace{w_k^c}_{\text{class feature weights}} \underbrace{\frac{1}{Z} \sum_i \sum_j A_{ij}^k}_{\substack{\text{global average pooling} \\ \text{feature map}}} \quad (3)$$

³ We find that Grad-CAM maps become progressively worse as we move to earlier convolutional layers as they have smaller receptive fields and only focus on less semantic local features.

Let us define F^k to be the global average pooled output,

$$F^k = \frac{1}{Z} \sum_i \sum_j A_{ij}^k \quad (4)$$

CAM computes the final scores by,

$$Y^c = \sum_k w_k^c \cdot F^k \quad (5)$$

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,

$$\frac{\partial Y^c}{\partial F^k} = \frac{\frac{\partial Y^c}{\partial A_{ij}^k}}{\frac{\partial F^k}{\partial A_{ij}^k}} \quad (6)$$

Taking partial derivative of (4) w.r.t. A_{ij}^k , we can see that $\frac{\partial F^k}{\partial A_{ij}^k} = \frac{1}{Z}$. Substituting this in (6), we get,

$$\frac{\partial Y^c}{\partial F^k} = \frac{\partial Y^c}{\partial A_{ij}^k} \cdot Z \quad (7)$$

From (5) we get that, $\frac{\partial Y^c}{\partial F^k} = w_k^c$. Hence,

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

Summing both sides of (8) over all pixels (i, j) ,

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

Since Z and w_k^c do not depend on (i, j) , rewriting this as

$$Z w_k^c = Z \sum_i \sum_j \frac{\partial Y^c}{\partial A_{ij}^k} \quad (10)$$

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

$$w_k^c = \sum_i \sum_j \frac{\partial Y^c}{\partial A_{ij}^k} \quad (11)$$

Up to a proportionality constant ($1/Z$) that gets normalized-out during visualization, the expression for w_k^c is identical to α_k^c used by Grad-CAM (1). Thus, Grad-CAM is a strict generalization of CAM. This generalization allows us to generate visual explanations from CNN-based models that cascade convolutional layers with much more complex interactions, such as those for image captioning and VQA (Sec. 8.2).

3.2 Guided Grad-CAM

While Grad-CAM is class-discriminative and localizes relevant image regions, it lacks the ability to highlight fine-grained details like pixel-space gradient visualization methods (Guided Backpropagation [53], Deconvolution [57]). Guided Backpropagation visualizes gradients with respect to the image where negative gradients are suppressed when backpropagating through ReLU layers. Intuitively, this aims to capture pixels detected by neurons, not the ones that suppress neurons. See Figure 1c, where Grad-CAM can easily localize the cat; however, it is unclear from the coarse heatmap why the network predicts this particular instance as ‘tiger cat’. In order to combine the best aspects of both, we fuse Guided Backpropagation and Grad-CAM visualizations via element-wise multiplication ($L_{\text{Grad-CAM}}^c$ is first upsampled to the input image resolution using bilinear 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 highlights the ‘tiger cat’ but not the ‘boxer (dog)’). Replacing Guided Backpropagation with Deconvolution gives similar results, but we found Deconvolution visualizations to have artifacts and Guided Backpropagation to be generally less noisy.

3.3 Counterfactual Explanations

Using a slight modification to Grad-CAM, we can obtain explanations that highlight support for regions that would make the network change its prediction. As a consequence, removing concepts occurring in those regions would make the model more confident about its prediction. We refer to this explanation modality as counterfactual explanations.

Specifically, we negate the gradient of y^c (score for class c) with respect to feature maps A of a convolutional layer. Thus the importance weights α_k^c now become

$$\alpha_k^c = \underbrace{\frac{1}{Z} \sum_i \sum_j}_{\text{global average pooling}} - \underbrace{\frac{\partial y^c}{\partial A_{ij}^k}}_{\text{Negative gradients}} \quad (12)$$

As in (2), we take a weighted sum of the forward activation maps, A , with weights α_k^c , and follow it by a ReLU to obtain counterfactual explanations as shown in Fig. 3.

4 Evaluating Localization Ability of Grad-CAM

4.1 Weakly-supervised Localization

In this section, we evaluate the localization capability of Grad-CAM in the context of image classification. The ImageNet

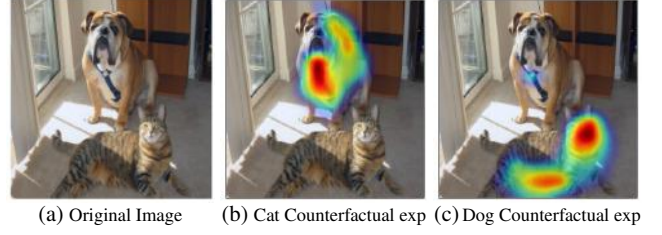


Fig. 3: Counterfactual Explanations with Grad-CAM

localization challenge [14] requires approaches 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.

Given an image, we first obtain class predictions from our network and then generate Grad-CAM maps for each of the predicted classes and binarize them with a threshold of 15% of the max intensity. This results in connected segments of pixels and we draw a bounding box around the single largest segment. Note that this is weakly-supervised localization – the models were never exposed to bounding box annotations during training.

We evaluate Grad-CAM localization with off-the-shelf pre-trained VGG-16 [52], AlexNet [33] and GoogleNet [54] (obtained from the Caffe [27] Zoo). Following ILSVRC-15 evaluation, we report both top-1 and top-5 localization errors on the val set in Table. 1. Grad-CAM localization errors are significantly better than those achieved by c-MWP [58] and Simonyan *et al.* [51], which use grab-cut to post-process image space gradients into heat maps. Grad-CAM for VGG-16 also achieves better top-1 localization error than CAM [59], which requires a change in the model architecture, necessitates re-training and thereby achieves worse classification errors (2.98% worse top-1), while Grad-CAM does not compromise on classification performance.

		Classification		Localization	
		Top-1	Top-5	Top-1	Top-5
VGG-16	Backprop [51]	30.38	10.89	61.12	51.46
	c-MWP [58]	30.38	10.89	70.92	63.04
	Grad-CAM (ours)	30.38	10.89	56.51	46.41
	CAM [59]	33.40	12.20	57.20	45.14
AlexNet	c-MWP [58]	44.2	20.8	92.6	89.2
	Grad-CAM (ours)	44.2	20.8	68.3	56.6
GoogleNet	Grad-CAM (ours)	31.9	11.3	60.09	49.34
	CAM [59]	31.9	11.3	60.09	49.34

Table 1: Classification and localization error % on ILSVRC-15 val (lower is better) for VGG-16, AlexNet and GoogleNet. We see that Grad-CAM achieves superior localization errors without compromising on classification performance.

4.2 Weakly-supervised Segmentation

Semantic segmentation involves the task of assigning each pixel in the image an object class (or background class). Be-

ing a challenging task, this requires expensive pixel-level annotation. The task of weakly-supervised segmentation involves segmenting objects with just image-level annotation, which can be obtained relatively cheaply from image classification datasets. In recent work, Kolesnikov *et al.* [32] 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, encouraging segmentation network to match these cues, 2) to expand object seeds to regions of reasonable size based on information about which classes can occur in an image, 3) to constrain segmentations to object boundaries that alleviates the problem of imprecise boundaries already at training time. They showed that their proposed loss function, consisting of the above three losses leads to better segmentation.

However, their algorithm is sensitive to the choice of weak localization seed, without which the network fails to localize objects correctly. In their work, they used CAM maps from a VGG-16 based network which are used as object seeds for weakly localizing foreground classes. We replaced the CAM maps with Grad-CAM obtained from a standard VGG-16 network and obtain a Intersection over Union (IoU) score of 49.6 (compared to 44.6 obtained with CAM) on the PASCAL VOC 2012 segmentation task. Fig. 4 shows some qualitative results.

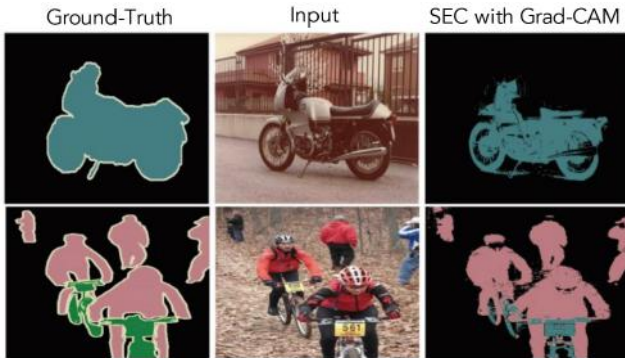


Fig. 4: PASCAL VOC 2012 Segmentation results with Grad-CAM as seed for SEC [32].

4.3 Pointing Game

Zhang *et al.* [58] introduced the Pointing Game experiment to evaluate the discriminativeness of different visualization methods for localizing target objects in scenes. Their evaluation protocol first cues each visualization technique with the ground-truth object label and extracts the maximally activated point on the generated heatmap. It then evaluates if the point lies within one of the annotated instances of the target object category, thereby counting it as a hit or a miss.

The localization accuracy is then calculated as

$Acc = \frac{\#Hits}{\#Hits + \#Misses}$. However, this evaluation only measures precision of the visualization technique. We modify the protocol to also measure recall – we compute localization

maps for top-5 class predictions from the CNN classifiers⁴ and evaluate them using the pointing game setup with an additional option to reject any of the top-5 predictions from the model if the maximally activated point in the map is below a threshold, *i.e.* if the visualization correctly rejects the predictions which are absent from the ground-truth categories, it gets that as a hit. We find that Grad-CAM outperforms c-MWP [58] by a significant margin (70.58% vs. 60.30%). Qualitative examples comparing c-MWP [58] and Grad-CAM on can be found in Sec. D⁵.

5 Evaluating Visualizations

In this section, we describe the human studies and experiments we conducted to understand the interpretability vs. faithfulness tradeoff of our approach to model predictions. Our first human study evaluates the main premise of our approach – are Grad-CAM visualizations more class discriminative than previous techniques? Having established that, we turn to understanding whether it can lead an end user to trust the visualized models appropriately. For these experiments, we compare VGG-16 and AlexNet finetuned on PASCAL VOC 2007 *train* and visualizations evaluated on *val*.

5.1 Evaluating Class Discrimination

In order to measure whether Grad-CAM helps distinguish between classes, we select images from the PASCAL VOC 2007 *val* set, which contain exactly 2 annotated categories and create visualizations for each one of them. For both VGG-16 and AlexNet CNNs, we obtain category-specific visualizations using four techniques: Deconvolution, Guided Backpropagation, and Grad-CAM versions of each of these methods (Deconvolution Grad-CAM and Guided Grad-CAM). We show these visualizations to 43 workers on Amazon Mechanical Turk (AMT) and ask them “Which of the two object categories is depicted in the image?” (shown in Fig. 5).

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 in Table. 2. 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% → 60.37%). Guided Grad-CAM performs the best among all methods.

⁴ We use GoogLeNet finetuned on COCO, as provided by [58].

⁵ c-MWP [58] highlights arbitrary regions for predicted but non-existent categories, unlike Grad-CAM maps which typically do not.

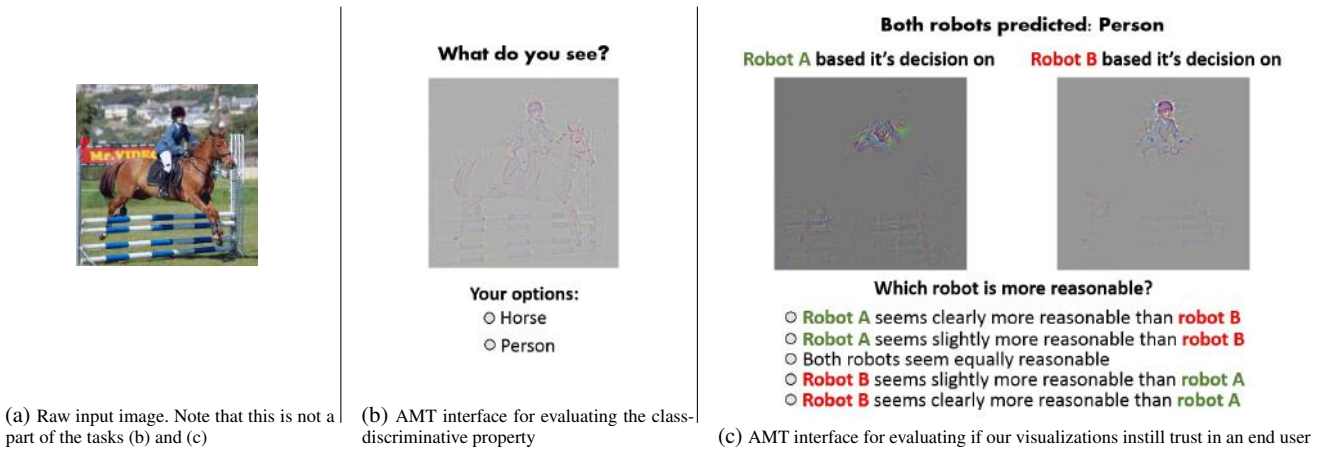


Fig. 5: AMT interfaces for evaluating different visualizations for class discrimination (b) and trustworthiness (c). Guided Grad-CAM outperforms baseline approaches (Guided-backprop and Deconvolution) showing that our visualizations are more class-discriminative and help humans place trust in a more accurate classifier.

Interestingly, our results indicate that Deconvolution is more class-discriminative than Guided Backpropagation (53.33% vs. 44.44%), although Guided Backpropagation is more aesthetically pleasing. To the best of our knowledge, our evaluations are the first to quantify this subtle difference.

Method	Human Accuracy	Classification	Relative Reliability	Rank Correlation w/ Occlusion
Guided Backpropagation	44.44		+1.00	0.168
Guided Grad-CAM	61.23		+1.27	0.261

Table 2: Quantitative Visualization Evaluation. Guided Grad-CAM enables humans to differentiate between visualizations of different classes (Human Classification Accuracy) and pick more reliable models (Relative Reliability). It also accurately reflects the behavior of the model (Rank Correlation w/ Occlusion).

5.2 Evaluating Trust

Given two prediction explanations, we 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) on PASCAL classification. 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, 54 AMT workers were 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. 5. To eliminate any biases, VGG-16 and AlexNet were assigned to be ‘model-1’ with approximately equal probability. Remarkably, as can be seen in Table. 2, we find that human subjects are able to identify the more accurate classifier (VGG-16 over AlexNet) *simply from the prediction explanations, despite both models making identical predictions*. With Guided Backpropagation, humans assign VGG-16 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 saying that VGG-16 is clearly more reliable. Thus, our visualizations can help users place trust in a model that generalizes better, just 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 trade-off 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 [47].

For comparison, we need a reference explanation with high local-faithfulness. One obvious choice for such a visualization is image occlusion [57], where we measure the difference in CNN scores when patches of the input image are masked. Interestingly, patches which change the CNN score are also patches to which Grad-CAM and Guided Grad-CAM assign high intensity, achieving rank correlation 0.254 and 0.261 (vs. 0.168, 0.220 and 0.208 achieved by Guided Backpropagation, c-MWP and CAM respectively) averaged over 2510 images in the PASCAL 2007 val set. This shows that Grad-CAM is more faithful to the original model compared to prior methods. Through localization experiments and human studies, we see that Grad-CAM visualizations are *more interpretable*, and through correlation with occlusion maps, we see that Grad-CAM is *more faithful* to the model.

6 Diagnosing image classification CNNs with Grad-CAM

In this section we further demonstrate the use of Grad-CAM in analyzing failure modes of image classification CNNs, understanding the effect of adversarial noise, and identifying and removing biases in datasets, in the context of VGG-16 pretrained on imagenet.

6.1 Analyzing failure modes for VGG-16

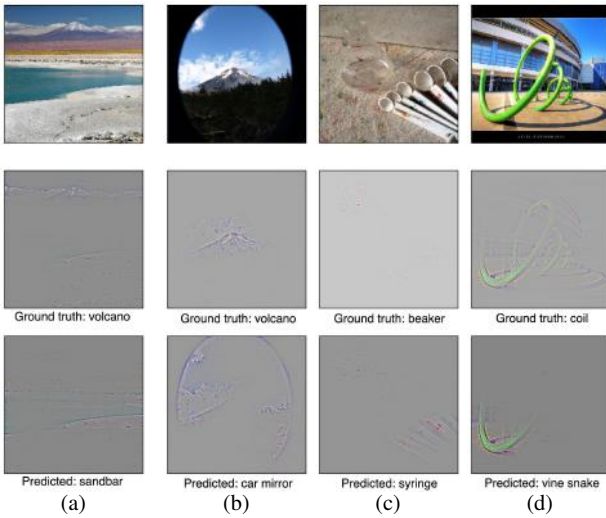


Fig. 6: In these cases the model (VGG-16) failed to predict the correct class in its top 1 (a and d) and top 5 (b and c) predictions. Humans would find it hard to explain some of these predictions without looking at the visualization for the predicted class. But with Grad-CAM, these mistakes seem justifiable.

In order to see what mistakes a network is making, we first get a list of examples that the network (VGG-16) fails to classify correctly. For these misclassified examples, we use Guided Grad-CAM to visualize both the correct and the predicted class. As seen in Fig. 6, some failures are due to ambiguities inherent in ImageNet classification. We can also see that *seemingly unreasonable predictions have reasonable explanations*, an observation also made in HOGgles [56]. A major advantage of Guided Grad-CAM visualizations over other methods is that due to its high-resolution and ability to be class-discriminative, it readily enables these analyses.

6.2 Effect of adversarial noise on VGG-16

Goodfellow *et al.* [22] demonstrated the vulnerability of current deep networks to adversarial examples, which are slight imperceptible perturbations of input images that fool the network into misclassifying them with high confidence. We generate adversarial images for an ImageNet-pretrained VGG-16 model such that it assigns high probability (> 0.9999) to a

category that is not present in the image and low probabilities to categories that are present. We then compute Grad-CAM visualizations for the categories that are present. As shown in Fig. 7, despite the network being certain about the absence of these categories ('tiger cat' and 'boxer'), Grad-CAM visualizations can correctly localize them. This shows that Grad-CAM is fairly robust to adversarial noise.

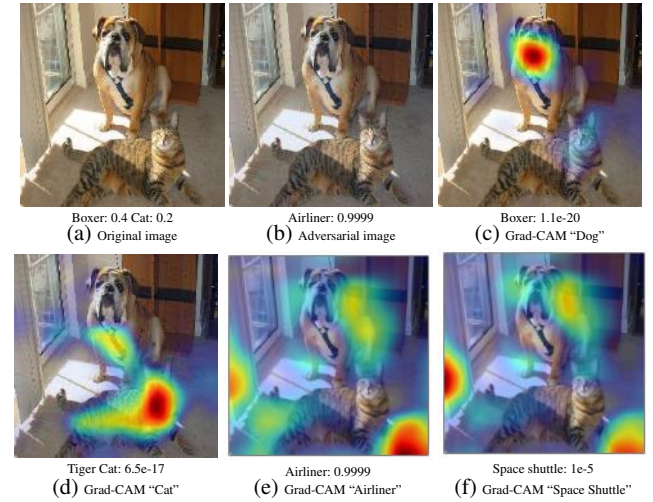


Fig. 7: (a-b) Original image and the generated adversarial image for category "airliner". (c-d) Grad-CAM visualizations for the original categories "tiger cat" and "boxer (dog)" along with their confidence. Despite the network being completely fooled into predicting the dominant category label of "airliner" with high confidence (>0.9999), Grad-CAM can localize the original categories accurately. (e-f) Grad-CAM for the top-2 predicted classes "airliner" and "space shuttle" seems to highlight the background.

6.3 Identifying bias in dataset

In this section, we demonstrate another use of Grad-CAM: identifying and reducing bias in training datasets. Models trained on biased datasets may not generalize to real-world scenarios, or worse, may perpetuate biases and stereotypes (w.r.t. gender, race, age, *etc.*). We finetune an ImageNet-pretrained VGG-16 model for a "doctor" vs. "nurse" binary classification task. We built our training and validation splits using the top 250 relevant images (for each class) from a popular image search engine. And the test set was controlled to be balanced in its distribution of genders across the two classes. Although the trained model achieves good validation accuracy, it does not generalize well (82% test accuracy).

Grad-CAM visualizations of the model predictions (see the red box⁶ regions in the middle column of Fig. 8) revealed that the model had learned to look at the person's face / hairstyle to distinguish nurses from doctors, thus learning

⁶ The green and red boxes are drawn manually to highlight correct and incorrect focus of the model.

a gender stereotype. Indeed, the model was misclassifying several female doctors to be a nurse and male nurses to be a doctor. Clearly, this is problematic. Turns out the image search results were gender-biased (78% of images for doctors were men, and 93% images for nurses were women).

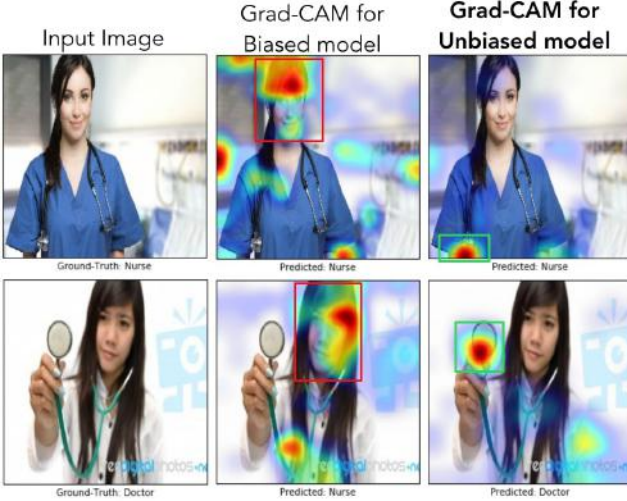


Fig. 8: In the first row, we can see that even though both models made the right decision, the biased model (model1) was looking at the face of the person to decide if the person was a nurse, whereas the unbiased model was looking at the short sleeves to make the decision. For the example image in the second row, the biased model made the wrong prediction (misclassifying a doctor as a nurse) by looking at the face and the hairstyle, whereas the unbiased model made the right prediction looking at the white coat, and the stethoscope.

Through these intuitions gained from Grad-CAM visualizations, we reduced bias in the training set by adding in images of male nurses and female doctors, while maintaining the same number of images per class as before. The re-trained model not only generalizes better (90% test accuracy), but also looks at the right regions (last column of Fig. 8). This experiment demonstrates a proof-of-concept that Grad-CAM can help detect and remove biases in datasets, which is important not just for better generalization, but also for fair and ethical outcomes as more algorithmic decisions are made in society.

7 Textual Explanations with Grad-CAM

Equation. (1) gives a way to obtain neuron-importance, α , for each neuron in a convolutional layer for a particular class. There have been hypotheses presented in the literature [60, 57] that neurons act as concept ‘detectors’. Higher positive values of the neuron importance indicate that the presence of that concept leads to an increase in the class score, whereas higher negative values indicate that its absence leads to an increase in the score for the class.

Given this intuition, let’s examine a way to generate textual explanations. In recent work, Bau *et al.* [4] proposed

an approach to automatically name neurons in any convolutional layer of a trained network. These names indicate concepts that the neuron looks for in an image. Using their approach, we first obtain neuron names for the last convolutional layer. Next, we sort and obtain the top-5 and bottom-5 neurons based on their class-specific importance scores, α_k . The names for these neurons can be used as text explanations.

Fig. 9 shows some examples of visual and textual explanations for the image classification model (VGG-16) trained on the Places365 dataset [61]. In (a), the positively important neurons computed by (1) look for intuitive concepts such as book and shelf that are indicative of the class ‘Book-store’. Also note that the negatively important neurons look for concepts such as sky, road, water and car which don’t occur in ‘Book-store’ images. In (b), for predicting ‘waterfall’, both visual and textual explanations highlight ‘water’ and ‘stratified’ which are descriptive of ‘waterfall’ images. (e) is a failure case due to misclassification as the network predicted ‘rope-bridge’ when there is no rope, but still the important concepts (water and bridge) are indicative of the predicted class. In (f), while Grad-CAM correctly looks at the door and the staircase on the paper to predict ‘Elevator door’, the neurons detecting doors did not pass the IoU threshold⁷ of 0.05 (chosen in order to suppress the noise in the neuron names), and hence are not part of the textual explanations. More qualitative examples can be found in the Sec. F.

8 Grad-CAM for Image Captioning and VQA

Finally, we apply Grad-CAM to vision & language tasks such as image captioning [7, 29, 55] and Visual Question Answering (VQA) [3, 20, 42, 46]. We find that Grad-CAM leads to interpretable visual explanations for these tasks as compared to baseline visualizations which do not change noticeably across changing predictions. Note that existing visualization techniques either are not class-discriminative (Guided Back-propagation, Deconvolution), or simply cannot be used for these tasks/architectures, or both (CAM, c-MWP).

8.1 Image Captioning

In this section, we visualize spatial support for an image captioning model using Grad-CAM. We build Grad-CAM on top of the publicly available neuraltalk2⁸ implementation [31] that uses a finetuned VGG-16 CNN for images and an LSTM-based language model. Note that this model does not have an explicit attention mechanism. Given a caption, we compute the gradient of its log probability w.r.t. units in

⁷ Area of overlap between ground truth concept annotation and neuron activation over area of their union. More details of this metric can be found in [4]

⁸ <https://github.com/karpathy/neuraltalk2>

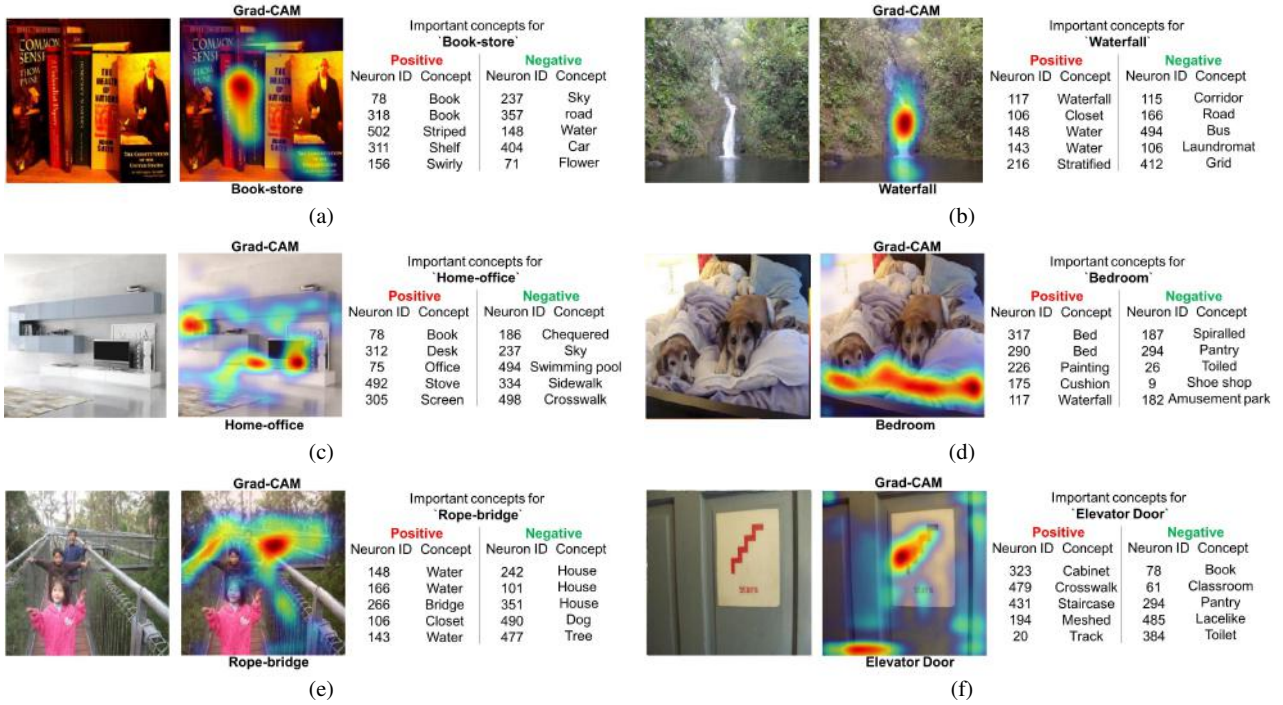


Fig. 9: Examples showing visual explanations and textual explanations for VGG-16 trained on Places365 dataset [61]. For textual explanations we provide the most important neurons for the predicted class along with their names. Important neurons can be either be persuasive (positive importance) or inhibitive (negative importance). The first 2 rows show success cases, and the last row shows 2 failure cases. We see that in (a), the important neurons computed by (1) look for concepts such as book and shelf which are indicative of class ‘Book-store’ which is fairly intuitive.

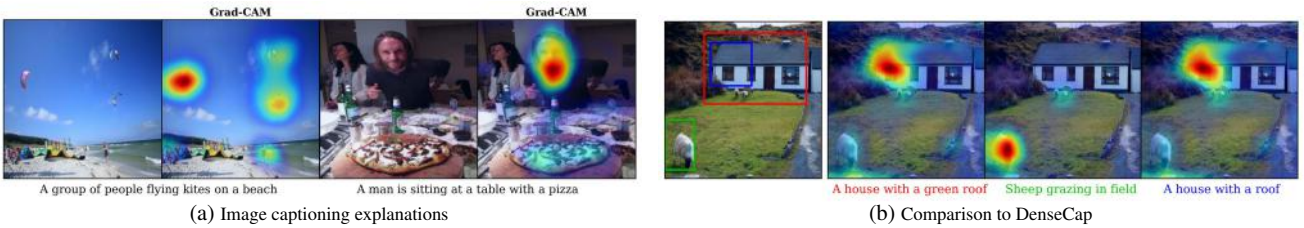


Fig. 10: Interpreting image captioning models: We use our class-discriminative localization technique, Grad-CAM to find spatial support regions for captions in images. Fig. 10a Visual explanations from image captioning model [31] highlighting image regions considered to be important for producing the captions. Fig. 10b Grad-CAM localizations of a *global* or *holistic* captioning model for captions generated by a dense captioning model [29] 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.

the last convolutional layer of the CNN (*conv5_3* for VGG-16) and generate Grad-CAM visualizations as described in Sec. 3. See Fig. 10a. In the first example, Grad-CAM maps for the generated caption localize every occurrence of both the kites and people despite their relatively small size. In the next example, Grad-CAM correctly highlights the pizza and the man, but ignores the woman nearby, since ‘woman’ is not mentioned in the caption. More examples are in Sec. C.

Comparison to dense captioning. Johnson *et al.* [29] recently introduced the Dense Captioning (DenseCap) task that requires a system to jointly localize and caption salient regions in a given image. Their model consists of a Fully Convolutional Localization Network (FCLN) that produces bounding boxes for regions of interest and an LSTM-based language model that generates associated captions, all in a single forward pass. Using DenseCap, we generate 5 region-

specific captions per image with associated ground truth bounding boxes. Grad-CAM for a whole-image captioning model (neuraltalk2) should localize the bounding box the region-caption was generated for, which is shown in Fig. 10b. We quantify this by computing the ratio of mean activation inside *vs.* outside the box. Higher ratios are better because they indicate stronger attention to the region the caption was generated for. Uniformly highlighting the whole image results in a baseline ratio of 1.0 whereas Grad-CAM achieves 3.27 ± 0.18 . Adding high-resolution detail gives an improved baseline of 2.32 ± 0.08 (Guided Backpropagation) and the best localization at 6.38 ± 0.99 (Guided Grad-CAM). Thus, Grad-CAM is able to localize regions in the image that the DenseCap model describes, even though the holistic captioning model was never trained with bounding-box annotations.

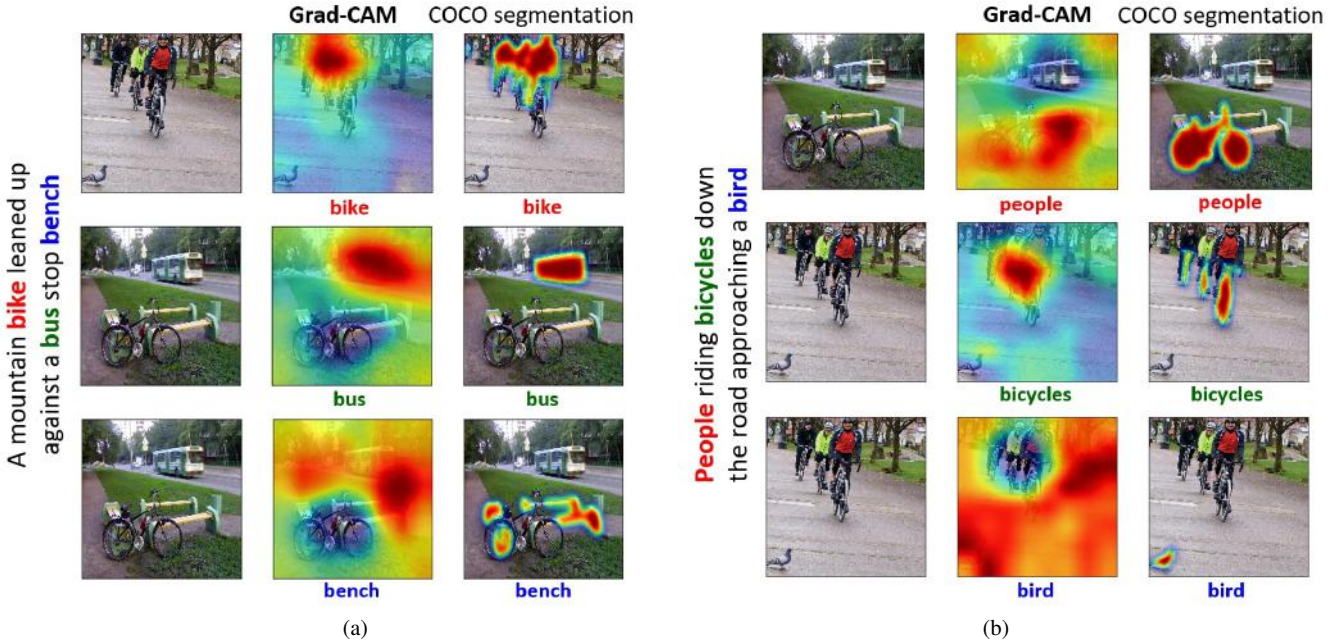


Fig. 11: Qualitative Results for our word-level captioning experiments: (a) Given the image on the left and the caption, we visualize Grad-CAM maps for the visual words “bike”, “bench” and “bus”. Note how well the Grad-CAM maps correlate with the COCO segmentation maps on the right column. (b) shows a similar example where we visualize Grad-CAM maps for the visual words “people”, “bicycle” and “bird”.

8.1.1 Grad-CAM for individual words of caption

In our experiment we use the Show and Tell model [55] pre-trained on MSCOCO without fine-tuning through the visual representation obtained from Inception [54] architecture. In order to obtain Grad-CAM map for individual words in the ground-truth caption we one-hot encode each of the visual words at the corresponding time-steps and compute the neuron importance score using Eq. (1) and combine with the convolution feature maps using Eq. (2).

Comparison to Human Attention We manually created an object category to word mapping that maps object categories like <person> to a list of potential fine-grained labels like [“child”, “man”, “woman”, ...]. We map a total of 830 visual words existing in COCO captions to 80 COCO categories. We then use the segmentation annotations for the 80 categories as human attention for this subset of matching words.

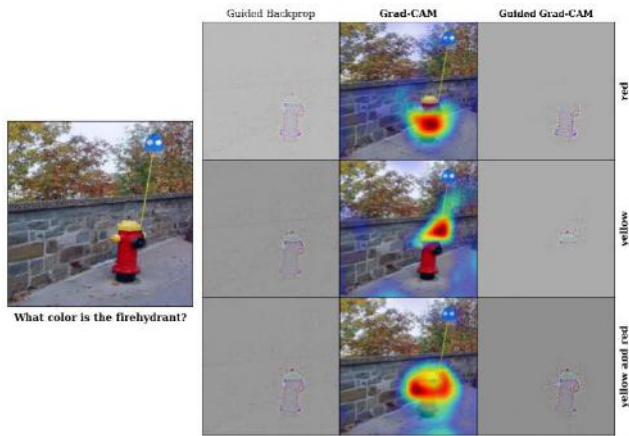
We then use the pointing evaluation from [58]. For each visual word from the caption, we generate the Grad-CAM map and then extract the maximally activated point. We then evaluate if the point lies within the human attention map segmentation for the corresponding COCO category, thereby counting it as a hit or a miss. The pointing accuracy is then calculated as

$Acc = \frac{\#Hits}{\#Hits + \#Misses}$. We perform this experiment on 1000 randomly sampled images from COCO dataset and obtain an accuracy of 30.0%. Some qualitative examples can be found in Fig. 11.

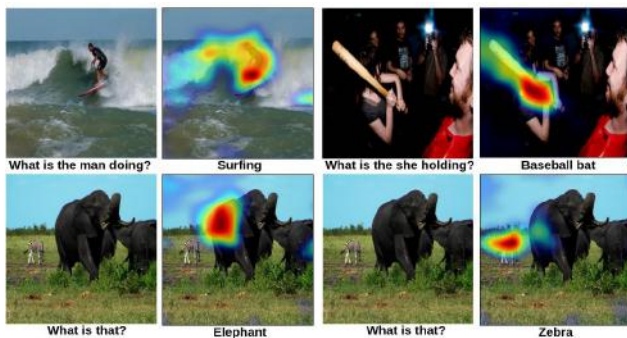
8.2 Visual Question Answering

Typical VQA pipelines [3, 20, 42, 46] consist of a CNN to process 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 (1000 being the size of the answer space). Since this is a classification problem, we pick an answer (the score y^c in (3)) and use its score to compute Grad-CAM visualizations over the image to explain the answer. Despite the complexity of the task, involving both visual and textual components, the explanations (of the VQA model from Lu *et al.* [38]) described in Fig. 12 are surprisingly intuitive and informative. We quantify the performance of Grad-CAM via correlation with occlusion maps, as in Sec. 5.3. Grad-CAM achieves a rank correlation (with occlusion maps) of 0.60 ± 0.038 whereas Guided Backpropagation achieves 0.42 ± 0.038 , indicating higher faithfulness of our Grad-CAM visualization.

Comparison to Human Attention. Das *et al.* [9] 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 VQA model from [38] on 1374 val question-image (QI) pairs from [3] using the rank correlation evaluation protocol as in [9]. Grad-CAM and human attention maps have a correlation of 0.136, which is higher than chance or random attention maps (zero correlation). This shows that despite not being trained on grounded image-text pairs, even non-attention



(a) Visualizing VQA model from [38]



(b) Visualizing ResNet based Hierarchical co-attention VQA model from [39]

Fig. 12: 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 visualizations are highly interpretable and help explain any target prediction – for “red”, the model focuses on the bottom red part of the firehydrant; when forced to answer “yellow”, the model concentrates on it’s top yellow cap, and when forced to answer “yellow and red”, it looks at the whole firehydrant! (b) Our approach is capable of providing interpretable explanations even for complex models.

based CNN + LSTM based VQA models are surprisingly good at localizing regions for predicting a particular answer.

Visualizing ResNet-based VQA model with co-attention.

Lu *et al.* [39] use a 200 layer ResNet [24] to encode the image, and jointly learn a hierarchical attention mechanism on the question and image. Fig. 12b shows Grad-CAM visualizations for this network. As we visualize deeper layers of the ResNet, we see small changes in Grad-CAM for most adjacent layers and larger changes between layers that involve dimensionality reduction. More visualizations for ResNets can be found in Sec. G. To the best of our knowledge, we are the first to visualize decisions from ResNet-based models.

9 Conclusion

In this work, we proposed a novel class-discriminative localization technique – Gradient-weighted Class Activation Mapping (Grad-CAM) – for making *any* CNN-based model more transparent by producing visual explanations. Further, we combined Grad-CAM localizations with existing high-resolution visualization techniques to obtain the best of both worlds – high-resolution and class-discriminative Guided Grad-CAM visualizations. Our visualizations outperform existing approaches on both axes – interpretability and faithfulness to original model. Extensive human studies reveal that our visualizations can discriminate between classes more accurately, better expose the trustworthiness of a classifier, and help identify biases in datasets. Further, we devise a way to identify important neurons through Grad-CAM and provide a way to obtain textual explanations for model decisions. Finally, we show the broad applicability of Grad-CAM to various off-the-shelf architectures for tasks such as image classification, image captioning and visual question answering. 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 and use it. Future work includes explaining decisions made by deep networks in domains such as reinforcement learning, natural language processing and video applications.

10 Acknowledgements

This work was funded in part by NSF CAREER awards to DB and DP, DARPA XAI grant to DB and DP, ONR YIP awards to DP and DB, ONR Grant N00014-14-1-0679 to DB, a Sloan Fellowship to DP, ARO YIP awards to DB and DP, an Allen Distinguished Investigator award to DP from the Paul G. Allen Family Foundation, ICTAS Junior Faculty awards to DB and DP, Google Faculty Research Awards to DP and DB, Amazon Academic Research Awards to DP and DB, AWS in Education Research grant to DB, and NVIDIA GPU donations to DB. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the U.S. Government, or any sponsor.

Appendix

A Appendix Overview

In the appendix, we provide:

- I - Ablation studies evaluating our design choices
- II - More qualitative examples for image classification, captioning and VQA

- III - More details of Pointing Game evaluation technique
- IV - Qualitative comparison to existing visualization techniques
- V - More qualitative examples of textual explanations

B Ablation studies

We perform several ablation studies to explore and validate our design choices for computing Grad-CAM visualizations. This includes visualizing different layers in the network, understanding importance of ReLU in (2), analyzing different types of gradients (for ReLU backward pass), and different gradient pooling strategies.

1. Grad-CAM for different layers

We show Grad-CAM visualizations for the “tiger-cat” class at different convolutional layers in AlexNet and VGG-16. As expected, the results from Fig. 13 show that localization becomes progressively worse as we move to earlier convolutional layers. This is because later convolutional layers better capture high-level semantic information while retaining spatial information than earlier layers, that have smaller receptive fields and only focus on local features.

2. Design choices

Method	Top-1 Loc 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 3: Localization results on ILSVRC-15 val for the ablations. Note that this evaluation is over 10 crops, while visualizations are single crop.

We evaluate different design choices via top-1 localization errors on the ILSVRC-15 val set [14]. See Table. 3.

2.1. Importance of ReLU in (3)

Removing ReLU ((3)) increases error by 15.3%. Negative values in Grad-CAM indicate confusion between multiple occurring classes.

2.2. 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). We observe that using GMP lowers the localization ability of Grad-CAM. An example can be found in Fig. 15 below. This may be due to the fact that max is statistically less robust to noise compared to the averaged gradient.

2.3. Effect of different ReLU on Grad-CAM

We experiment with Guided-ReLU [53] and Deconv-ReLU [57] as modifications to the backward pass of ReLU.

Guided-ReLU: Springenberg *et al.* [53] introduced Guided Backprop, where the backward pass of ReLU is modified to only pass positive gradients to regions of positive activations. Applying this change to the computation of Grad-CAM introduces a drop in the class-discriminative ability as can be seen in Fig. 16, but it marginally improves localization performance as can be seen in Table. 3.

Deconv-ReLU: In Deconvolution [57], Zeiler and Fergus introduced a modification to the backward pass of ReLU to only pass positive gradients. Applying this modification to the computation of Grad-CAM leads to worse results (Fig. 16). This indicates that negative gradients also carry important information for class-discriminateness.

C Qualitative results for vision and language tasks

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.

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. 17 correspond to the VGG-16 [52] network trained on ImageNet.

Fig. 17 shows randomly sampled examples from COCO [35] 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. 17.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. 17.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.

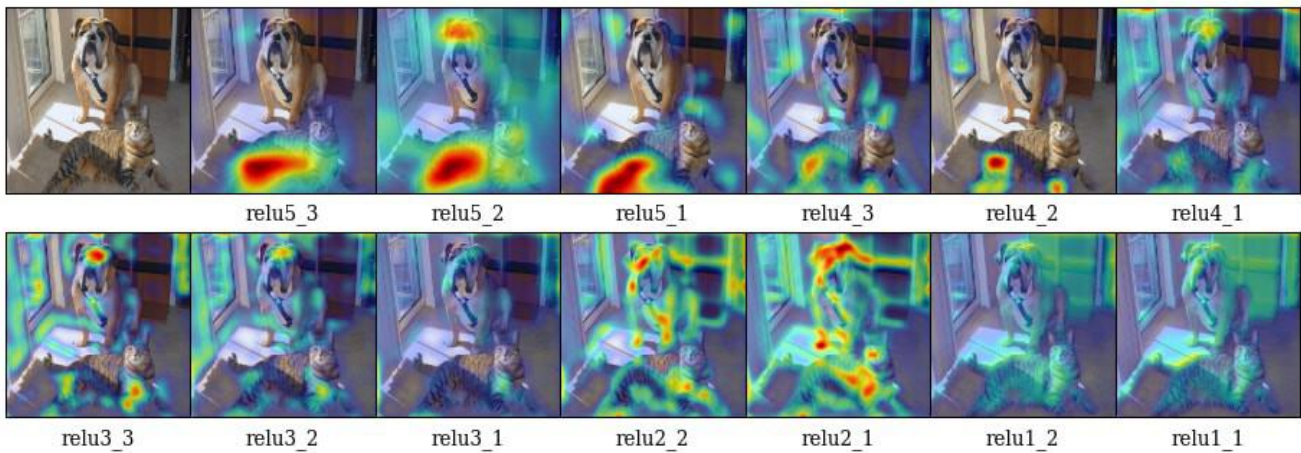


Fig. 13: Grad-CAM at different convolutional layers for the ‘tiger cat’ class. This figure analyzes how localizations change qualitatively as we perform Grad-CAM with respect to different feature maps in a CNN (VGG16 [52]). We find that the best looking visualizations are often obtained after the deepest 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, that deeper convolutional layer capture more semantic concepts.

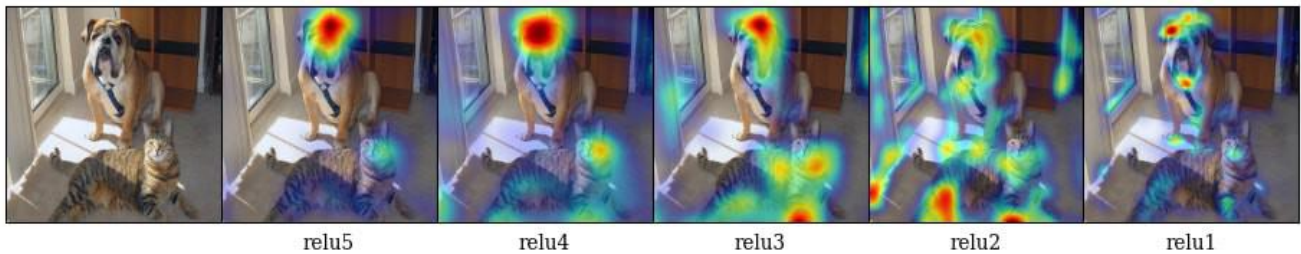


Fig. 14: Grad-CAM localizations for “tiger cat” category for different rectified convolutional layer feature maps for AlexNet.

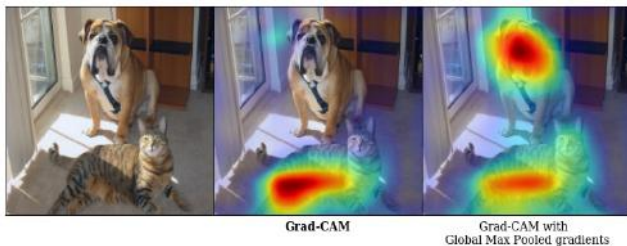


Fig. 15: Grad-CAM visualizations for “tiger cat” category with Global Average Pooling and Global Max Pooling.

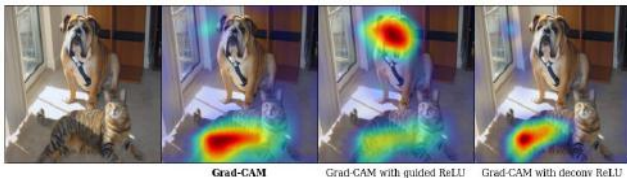


Fig. 16: Grad-CAM visualizations for “tiger cat” category for different modifications to the ReLU backward pass. The best results are obtained when we use the actual gradients during the computation of Grad-CAM.

“computer keyboard, keypad” class:

<http://i.imgur.com/QMhsRzf.jpg>

“sunglasses, dark glasses, shades” class:

<http://i.imgur.com/a1C7DGh.jpg>

2. Image Captioning

We use the publicly available Neuraltalk2 code and model⁹ 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 [35] 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.

3. Visual Question Answering (VQA)

We use Grad-CAM and Guided Grad-CAM to explain why a publicly available VQA model [38] 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

⁹ <https://github.com/karpathy/neuraltalk2>

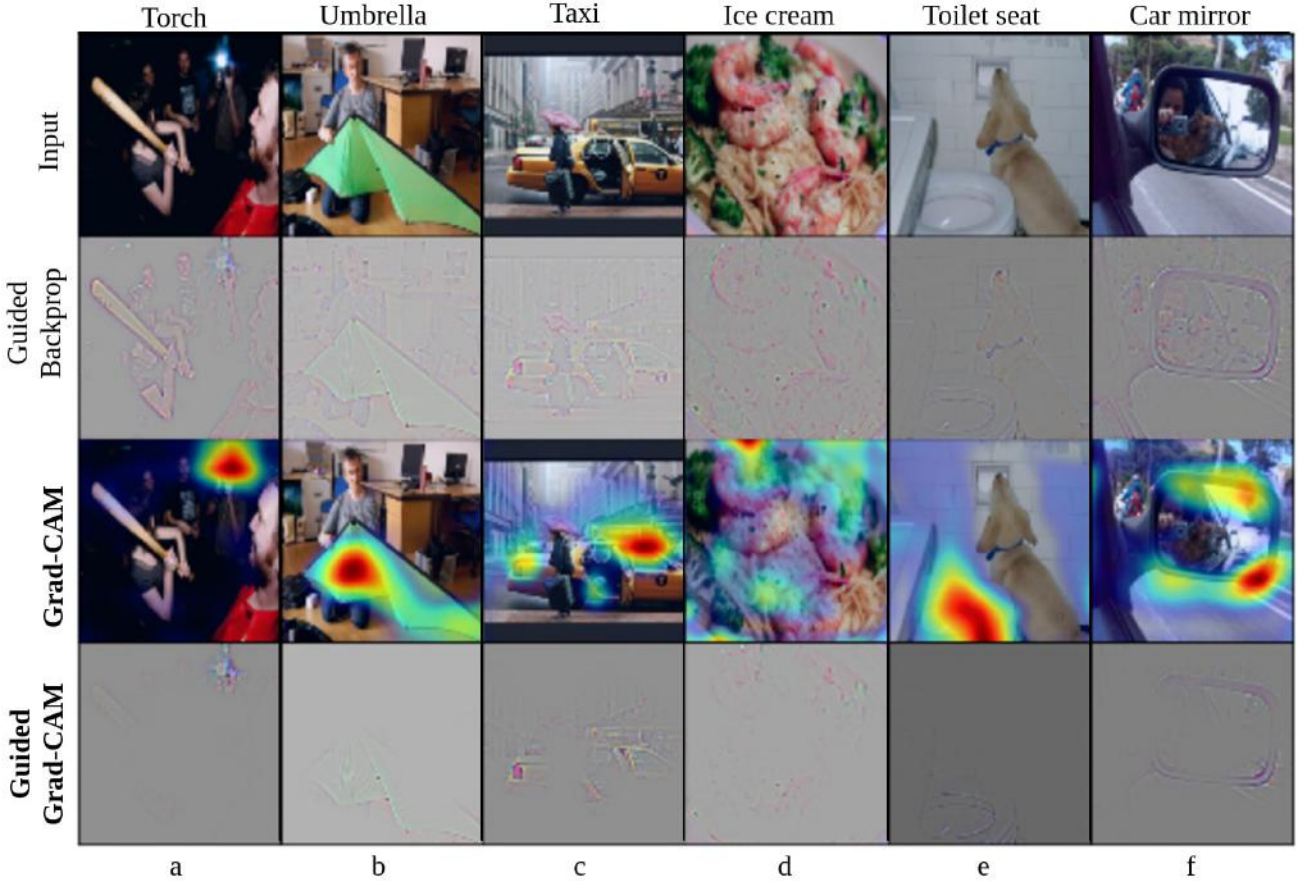


Fig. 17: Visualizations for randomly sampled images from the COCO validation dataset. Predicted classes are mentioned at the top of each column.

multiplied 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 [33], VGG-16 and VGG-19 [52]. 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. 19, 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. 19, for the question, “*Is this a whole orange?*”, the model looks for regions around the orange to answer “No”.

D More details of Pointing Game

In [58], 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.

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,

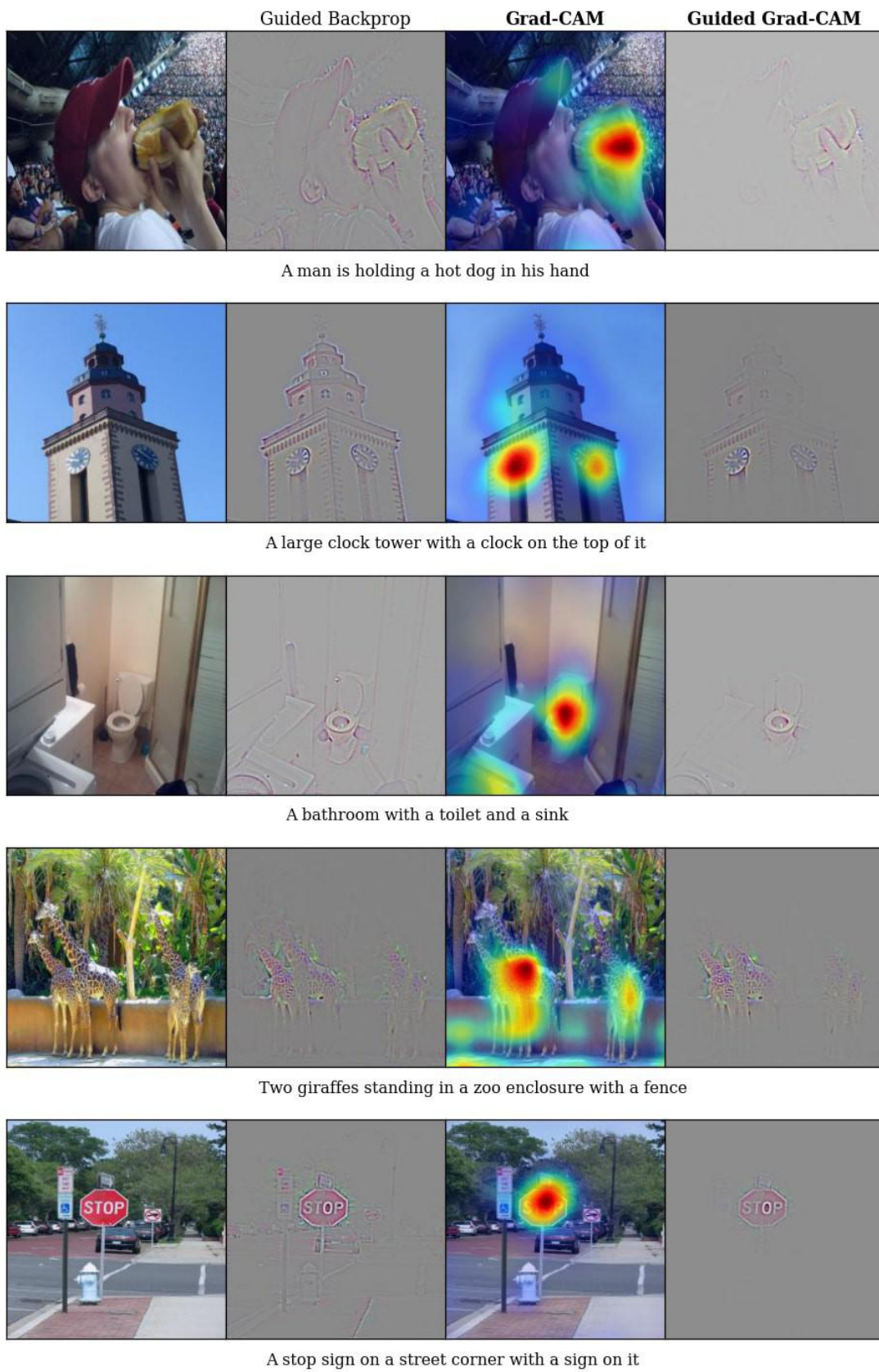


Fig. 18: Guided Backpropagation, Grad-CAM and Guided Grad-CAM visualizations for the captions produced by the Neuraltalk2 image captioning model.

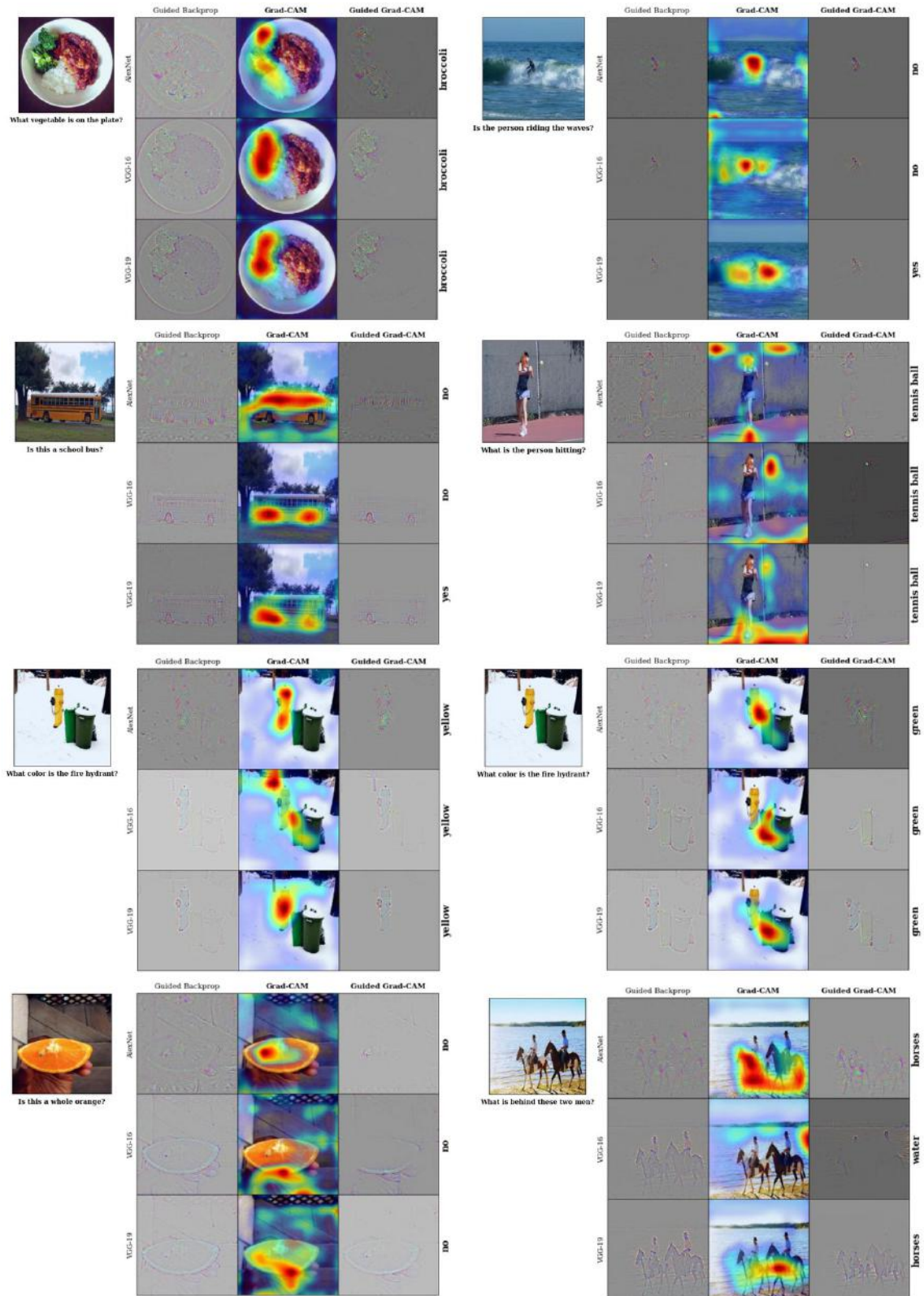


Fig. 19: Guided Backpropagation, Grad-CAM and Guided Grad-CAM visualizations for the answers from a VQA model. For each image-question pair, we show visualizations for AlexNet, VGG-16 and VGG-19. Notice how the attention changes in row 3, as we change the answer from *Yellow* to *Green*.

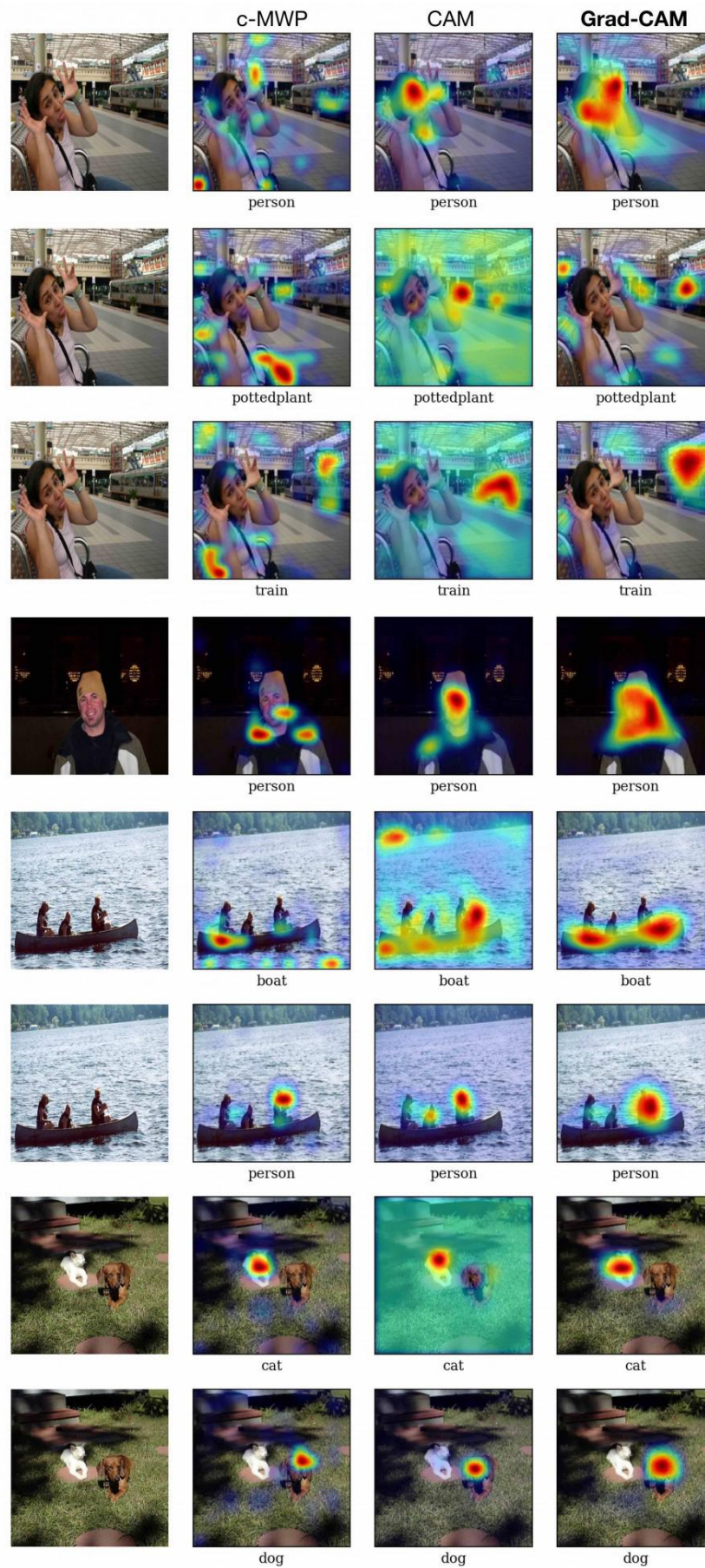


Fig. 20: Visualizations for ground-truth categories (shown below each image) for images sampled from the PASCAL [17] validation set.

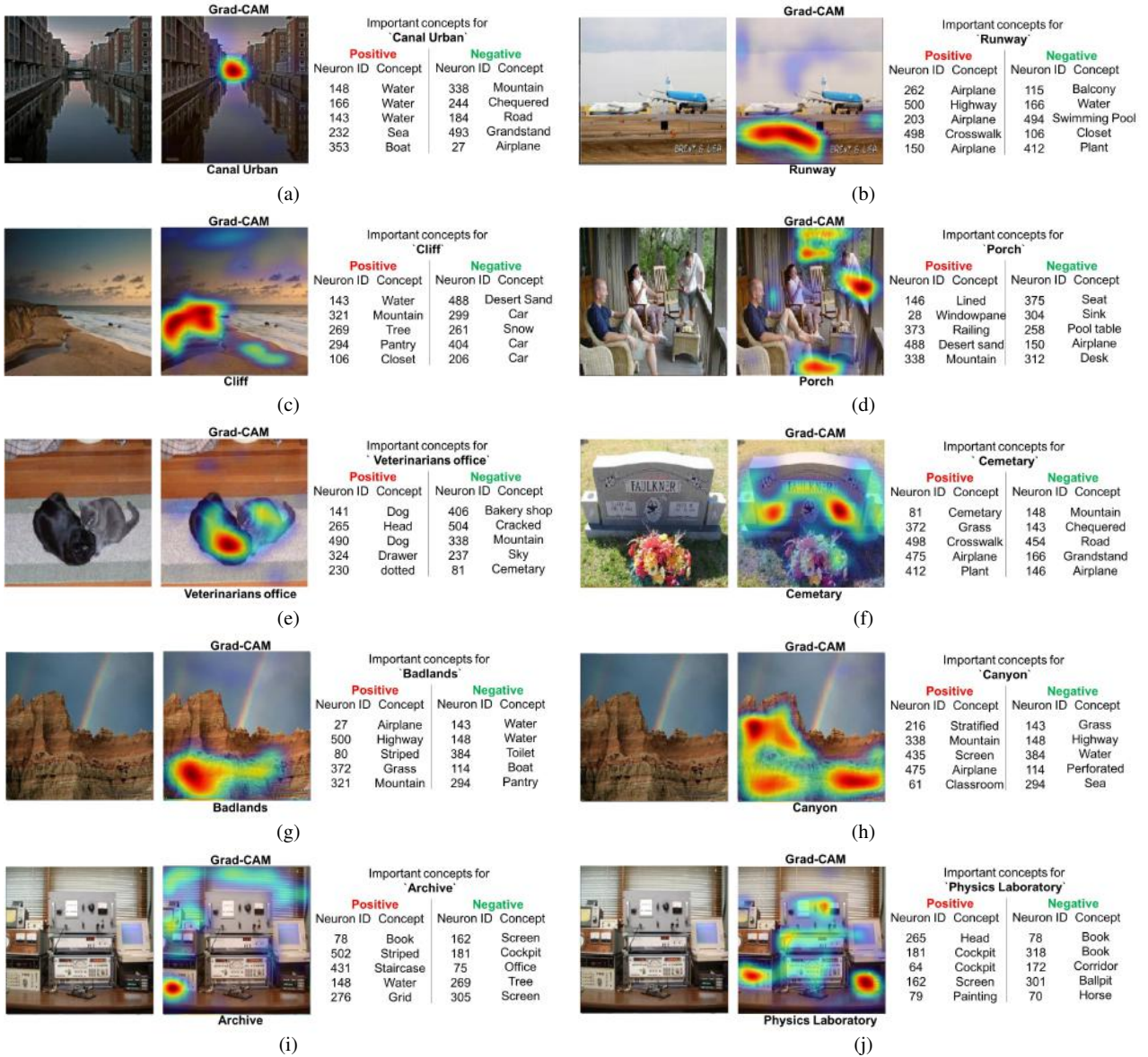


Fig. 21: More Qualitative examples showing visual explanations and textual explanations for VGG-16 trained on Places365 dataset ([61]). For textual explanations we provide the most important neurons for the predicted class along with their names. Important neurons can be either be persuasive (positively important) or inhibitive (negatively important). The first 3 rows show positive examples, and the last 2 rows show failure cases.

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% on VGG-16).

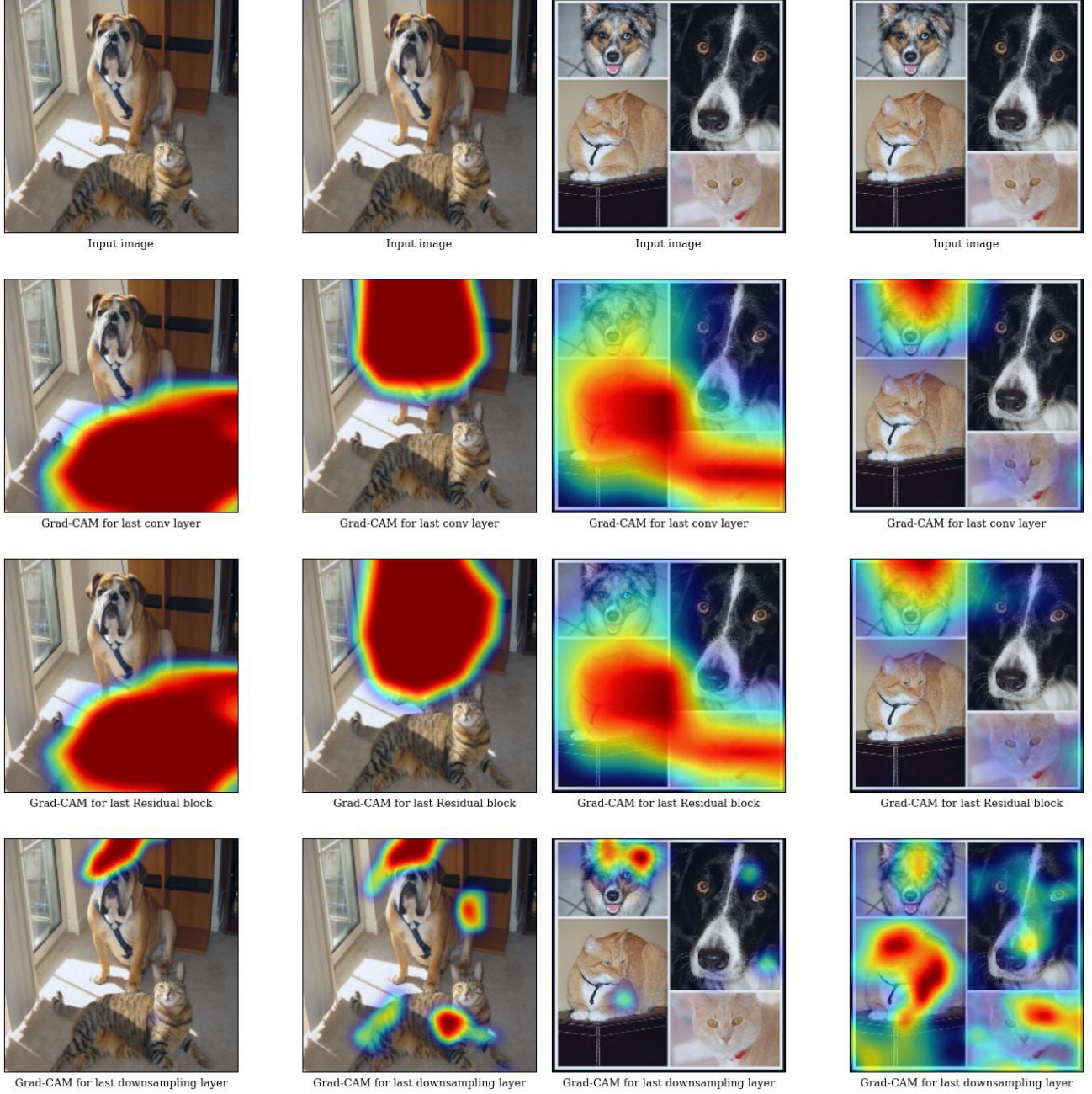
E Qualitative comparison to Excitation Backprop (c-MWP) and CAM

In this section we provide more qualitative results comparing Grad-CAM with CAM [59] and c-MWP [58] on Pascal [17].

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. 20.

F Visual and Textual explanations for Places dataset

Fig. 21 shows more examples of visual and textual explanations (Sec. 7) for the image classification model (VGG-16) trained on Places365 dataset ([61]).



(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.

Fig. 22: We observe that the discriminative ability of Grad-CAM significantly reduces as we encounter the downsampling layer.

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 ImageNet¹⁰.

Current ResNets [24] typically consist of residual blocks. One set of blocks use identity skip connections (shortcut connections between two layers having identical output di-

mensions). These sets of residual blocks are interspersed with downsampling modules that alter dimensions of propagating signal. As can be seen in Fig. 22 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. 22). We observe similar trends for other ResNet architectures (18 and 50-layer).

¹⁰ We use the 200-layer ResNet architecture from <https://github.com/facebook/fb.resnet.torch>.

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. 1, 2, 10, 12
4. D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Computer Vision and Pattern Recognition*, 2017. 1, 3, 10
5. L. Bazzani, A. Bergamo, D. Anguelov, and L. Torresani. Self-taught object localization with deep networks. In *WACV*, 2016. 4
6. 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. 4
7. 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. 1, 10
8. 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
9. 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. 12
10. A. Das, S. Datta, G. Gkioxari, S. Lee, D. Parikh, and D. Batra. Embodied Question Answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
11. A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. Moura, D. Parikh, and D. Batra. Visual Dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1
12. A. Das, S. Kottur, J. M. Moura, S. Lee, and D. Batra. Learning cooperative visual dialog agents with deep reinforcement learning. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017. 1
13. H. de Vries, F. Strub, S. Chandar, O. Pietquin, H. Larochelle, and A. C. Courville. Guesswhat?! visual object discovery through multi-modal dialogue. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1
14. 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, 6, 14
15. A. Dosovitskiy and T. Brox. Inverting Convolutional Networks with Convolutional Networks. In *CVPR*, 2015. 3
16. D. Erhan, Y. Bengio, A. Courville, and P. Vincent. Visualizing Higher-layer Features of a Deep Network. *University of Montreal*, 1341, 2009. 3
17. M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>, 2009. 19, 20
18. 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. 1
19. 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
20. 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. 1, 10, 12
21. R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In *CVPR*, 2014. 1
22. I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *stat*, 2015. 9
23. D. Gordon, A. Kembhavi, M. Rastegari, J. Redmon, D. Fox, and A. Farhadi. Iqa: Visual question answering in interactive environments. *arXiv preprint arXiv:1712.03316*, 2017. 1
24. K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 2, 13, 21
25. D. Hoiem, Y. Chodpathumwan, and Q. Dai. Diagnosing Error in Object Detectors. In *ECCV*, 2012. 2
26. P. Jackson. *Introduction to Expert Systems*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 3rd edition, 1998. 2
27. 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. 6
28. E. Johns, O. Mac Aodha, and G. J. Brostow. Becoming the Expert - Interactive Multi-Class Machine Teaching. In *CVPR*, 2015. 2
29. J. Johnson, A. Karpathy, and L. Fei-Fei. DenseCap: Fully Convolutional Localization Networks for Dense Captioning. In *CVPR*, 2016. 1, 10, 11
30. 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
31. A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015. 10, 11
32. A. Kolesnikov and C. H. Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *ECCV*, 2016. 7
33. A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1, 5, 6, 16
34. M. Lin, Q. Chen, and S. Yan. Network in network. In *ICLR*, 2014. 3
35. 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. 14, 15
36. Z. C. Lipton. The Mythos of Model Interpretability. *ArXiv e-prints*, June 2016. 2, 3
37. J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 1
38. 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. 12, 13, 15
39. J. Lu, J. Yang, D. Batra, and D. Parikh. Hierarchical question-image co-attention for visual question answering. In *NIPS*, 2016. 13
40. A. Mahendran and A. Vedaldi. Salient deconvolutional networks. In *European Conference on Computer Vision*, 2016. 3
41. A. Mahendran and A. Vedaldi. Visualizing deep convolutional neural networks using natural pre-images. *International Journal of Computer Vision*, pages 1–23, 2016. 3, 4
42. M. Malinowski, M. Rohrbach, and M. Fritz. Ask your neurons: A neural-based approach to answering questions about images. In *ICCV*, 2015. 1, 10, 12
43. 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
44. 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
45. P. O. Pinheiro and R. Collobert. From image-level to pixel-level labeling with convolutional networks. In *CVPR*, 2015. 3
46. M. Ren, R. Kiros, and R. Zemel. Exploring models and data for image question answering. In *NIPS*, 2015. 1, 10, 12
47. M. T. Ribeiro, S. Singh, and C. Guestrin. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *SIGKDD*, 2016. 3, 8
48. R. R. Selvaraju, P. Chattopadhyay, M. Elhoseiny, T. Sharma, D. Batra, D. Parikh, and S. Lee. Choose your neuron: Incorporating domain knowledge through neuron-importance. In *Proceedings*

- of the European Conference on Computer Vision (ECCV), pages 526–541, 2018. 3
49. R. R. Selvaraju, S. Lee, Y. Shen, H. Jin, S. Ghosh, L. Heck, D. Batra, and D. Parikh. Taking a hint: Leveraging explanations to make vision and language models more grounded. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019. 3
 50. 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
 51. K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *CoRR*, abs/1312.6034, 2013. 3, 6
 52. K. Simonyan and A. Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *ICLR*, 2015. 5, 6, 14, 15, 16
 53. J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. A. Riedmiller. Striving for Simplicity: The All Convolutional Net. *CoRR*, abs/1412.6806, 2014. 2, 3, 6, 14
 54. C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 6, 12
 55. O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015. 1, 10, 12
 56. C. Vondrick, A. Khosla, T. Malisiewicz, and A. Torralba. HOGgles: Visualizing Object Detection Features. *ICCV*, 2013. 9
 57. M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014. 2, 3, 4, 6, 8, 10, 14
 58. J. Zhang, Z. Lin, J. Brandt, X. Shen, and S. Sclaroff. Top-down Neural Attention by Excitation Backprop. In *ECCV*, 2016. 4, 6, 7, 12, 16, 20
 59. B. Zhou, A. Khosla, L. A., A. Oliva, and A. Torralba. Learning Deep Features for Discriminative Localization. In *CVPR*, 2016. 2, 3, 5, 6, 20
 60. B. Zhou, A. Khosla, À. Lapedriza, A. Oliva, and A. Torralba. Object detectors emerge in deep scene cnns. *CoRR*, abs/1412.6856, 2014. 10
 61. B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. 10, 11, 20