

# Interpretable CNNs for Object Classification

Quanshi Zhang, Xin Wang, Ying Nian Wu, Huilin Zhou, and Song-Chun Zhu, *Fellow, IEEE*

**Abstract**—This paper proposes a generic method to learn interpretable convolutional filters in a deep convolutional neural network (CNN) for object classification, where each interpretable filter encodes features of a specific object part. Our method does not require additional annotations of object parts or textures for supervision. Instead, we use the same training data as traditional CNNs. Our method automatically assigns each interpretable filter in a high conv-layer with an object part of a certain category during the learning process. Such explicit knowledge representations in conv-layers of CNN help people clarify the logic encoded in the CNN, *i.e.* answering what patterns the CNN extracts from an input image and uses for prediction. We have tested our method using different benchmark CNNs with various structures to demonstrate the broad applicability of our method. Experiments have shown that our interpretable filters are much more semantically meaningful than traditional filters.

**Index Terms**—Convolutional Neural Networks, Interpretable Deep Learning

## 1 INTRODUCTION

In recent years, convolutional neural networks (CNNs) [8], [13], [16] have achieved superior performance in many visual tasks, such as object classification and detection. In spite of the good performance, a deep CNN has been considered a black-box model with weak feature interpretability for decades. Boosting the feature interpretability of a deep model gradually attracts increasing attention recently, but it presents significant challenges for state-of-the-art algorithms.

In this paper, we focus on a new task, *i.e.* without any additional annotations for supervision, revising a CNN to make its high conv-layers (*e.g.* the top two conv-layers) encode interpretable object-part knowledge. The revised CNN is termed an *interpretable CNN*.

More specifically, we propose a generic interpretable layer to ensure each filter in the proposed interpretable layer learns specific, discriminative object-part features. Filters in the interpretable layer are supposed to have some introspection of their feature representations and regularize their features towards object parts during the end-to-end learning. We trained interpretable CNNs on several benchmark datasets, and experimental results show that filters in the interpretable layer consistently represented the same object part across input images.

Note that our task of improving feature interpretability of a CNN is essentially different from the conventional visualization [5], [6], [19], [24], [28], [40] and diagnosis [2], [11], [18], [21] of pre-trained CNNs. The interpretable CNN learns more interpretable features, whereas previous methods mainly explain pre-

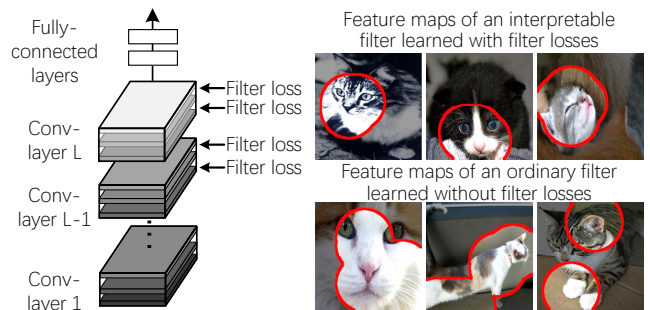


Fig. 1. Comparison of an interpretable filter's feature maps with a filter's feature maps in a traditional CNN.

trained neural networks, instead of learning new and more interpretable features.

In addition, as discussed in [2], filters in low conv-layers usually describe textural patterns, while filters in high conv-layers are more likely to represent part patterns. Therefore, we focus on part-based interpretability and propose a method to ensure each filter in a high conv-layer to represent an object part.

Fig. 1 visualizes the difference between a traditional filter and our interpretable filter. In a traditional CNN, a filter usually describes a mixture of patterns. For example, the filter may be activated by both the head part and the leg part of a cat. In contrast, the filter in our interpretable CNN is expected to be activated by a certain part.

Thus, the goal of this study can be summarized as follows. We propose a generic interpretable conv-layer to construct the interpretable CNN. Feature representations of the interpretable conv-layers are interpretable, *i.e.* each filter in the interpretable conv-layer learns to consistently represent the same object part across different images. In addition, the interpretable conv-layer needs to satisfy the following properties:

- The interpretable CNN needs to be learned with-

• Quanshi Zhang, Xin Wang, and Huilin Zhou is with the John Hopcroft Center and the MoE Key Lab of Artificial Intelligence, AI Institute, at the Shanghai Jiao Tong University, China. Ying Nian Wu and Song-Chun Zhu are with the University of California, Los Angeles, USA.

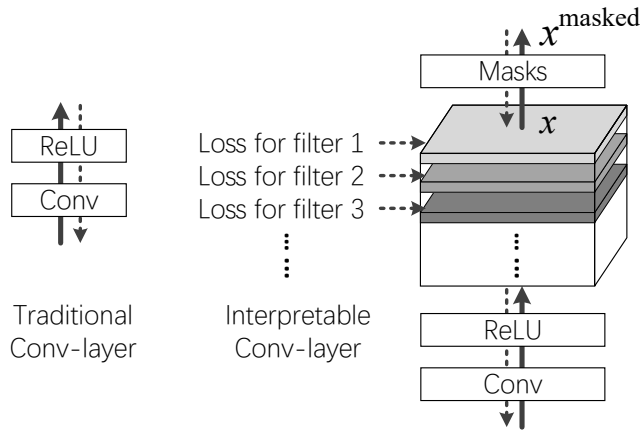


Fig. 2. Structures of an ordinary conv-layer and an interpretable conv-layer. Solid and dashed lines indicate the forward and backward propagations, respectively. During the forward propagation, our CNN assigns each interpretable filter with a specific mask *w.r.t.* each input image during the learning process.

out any additional annotations of object parts for supervision. We use the same training samples as the original CNN for learning.

- The interpretable CNN does not change the loss function of the classification task, and it can be broadly applied to different benchmark CNNs with various structures.
- As an exploratory research, learning strict representations of object parts may hurt a bit the discrimination power. However, we need to control the decrease within a small range.

**Method:** As shown in Fig. 2, we propose a simple yet effective loss. We simply add the loss to the feature map of each filter in a high conv-layer, so as to construct an interpretable conv-layer. The filter loss is proposed based the assumption that only a single target object is contained in the input image. The filter loss pushes the filter towards the representation of a specific object part.

Theoretically, we can prove that the loss encourages a low entropy of inter-category activations and a low entropy of spatial distributions of neural activations. In other words, this loss ensures that (i) each filter must encode an object part of a single object category, instead of representing multiple categories; (ii) The feature must consistently be triggered by a single specific part across multiple images, rather than be simultaneously triggered by different object regions in each input image. It is assumed that repetitive patterns on different object regions are more likely to describe low-level textures, instead of high-level parts.

**Value of feature interpretability:** Such explicit object-part representations in conv-layers of CNN can help people clarify the decision-making logic encoded in the CNN at the object-part level. Given an input image, the interpretable conv-layer enable people to

explicitly identify which object parts are memorized and used by the CNN for classification without ambiguity. Note that the automatically learned object part may not have an explicit name, *e.g.* a filter in the interpretable conv-layer may describe a partial region of a semantic part or the joint of two parts.

In critical applications, clear disentanglement of visual concepts in high conv-layers helps people trust a network’s prediction. As analyzed in [45], a good performance on testing images cannot always ensure correct feature representations considering potential dataset bias. For example, in [45], a CNN used an unreliable context—eye features—to identify the “lip-stick” attribute of a face image. Therefore, people need to semantically and visually explain what patterns are learned by the CNN.

**Contributions:** In this paper, we focus on a new task, *i.e.* end-to-end learning an interpretable CNN without any part annotations, where filters of high conv-layers represent specific object parts. We propose a simple yet effective method to learn interpretable filters, and the method can be broadly applied to different benchmark CNNs. Experiments show that our approach has significantly boosted feature interpretability of CNNs.

A preliminary version of this paper appeared in [46].

## 2 RELATED WORK

The interpretability and the discrimination power are two crucial aspects of a CNN [2]. In recent years, different methods are developed to explore the semantics hidden inside a CNN. Our previous paper [48] provides a comprehensive survey of recent studies in exploring visual interpretability of neural networks, including (i) the visualization and diagnosis of CNN representations, (ii) approaches for disentangling CNN representations into graphs or trees, (iii) the learning of CNNs with disentangled and interpretable representations, and (iv) middle-to-end learning based on model interpretability.

### 2.1 Interpretation of pre-trained neural networks

**Network visualization:** Visualization of filters in a CNN is the most direct way of exploring the pattern that is encoded by the filter. Gradient-based visualization [19], [28], [40] showed the appearance that maximized the score of a given unit. Furthermore, Bau *et al.* [2] defined and analyzed the interpretability of each filter. They classified all potential semantics into the following six types, *objects*, *parts*, *scenes*, *textures*, *materials*, and *colors*. We can further summarize the semantics of *objects* and *parts* as part patterns with specific contours and consider the other four semantics as textural patterns without explicit shapes. Recently, [20] provided tools to visualize filters of a CNN. Dosovitskiy *et al.* [5] proposed up-convolutional

nets to invert feature maps of conv-layers to images. However, up-convolutional nets cannot mathematically ensure the visualization result reflects actual neural representations.

Although above studies can produce clear visualization results, theoretically, gradient-based visualization of a filter usually selectively visualizes the strongest activations of a filter in a high conv-layer, instead of illustrating knowledge hidden behind all activations of the filter; otherwise, the visualization result will be chaotic. Similarly, [2] selectively analyzed the semantics of the highest 0.5% activations of each filter. In comparisons, we aim to purify the semantic meaning of each filter in a high conv-layer, *i.e.* letting most activations of a filter be explainable, instead of extracting meaningful neural activations for visualization.

**Pattern retrieval:** Unlike passive visualization, some methods actively retrieve certain units with certain meanings from CNNs. Just like mid-level features [30] of images, pattern retrieval mainly focuses on mid-level representations in conv-layers. For example, Zhou *et al.* [49], [50] selected units from feature maps to describe “scenes”. Simon *et al.* discovered objects from feature maps of conv-layers [26], and selected certain filters to represent object parts [27]. Zhang *et al.* [42] extracted certain neural activations of a filter to represent object parts in a weakly-supervised manner. They also disentangled CNN representations via active question-answering and summarized the disentangled knowledge using an And-Or graph [43]. [44] used human interactions to refine the AOG representation of CNN knowledge. [7] used a gradient-based method to explain visual question-answering. Other studies [12], [17], [34], [36] selected filters or neural activations with specific meanings from CNNs for various applications. Unlike the retrieval of meaningful neural activations from noisy features, our method aims to substantially boost the interpretability of features in intermediate conv-layers.

**Model diagnosis:** Many approaches have been proposed to diagnose CNN features, including exploring semantic meanings of convolutional filters [32], evaluating the transferability of filters [39], and the analysis of feature distributions of different categories [1]. The LIME [21] and the SHAP [18] are general methods to extract input units of a neural network that are used for the prediction score. For CNNs oriented to visual tasks, gradient-based visualization methods [6], [24] and [14] extracted image regions that are responsible for the network output, in order to clarify the logic of network prediction. These methods require people to manually check image regions accountable for the label prediction for each testing image. [10] extracted relationships between representations of various categories from a CNN. In contrast, given an interpretable CNN, people can

directly identify object parts or filters that are used for prediction.

As discussed by Zhang *et al.* [45], knowledge representations of a CNN may be significantly biased due to dataset bias, even though the CNN sometimes exhibits good performance. For example, a CNN may extract unreliable contextual features for prediction. Network-attack methods [11], [31], [32] diagnosed network representation flaws using adversarial samples of a CNN. For example, influence functions [11] can be used to generate adversarial samples, in order to fix the training set and further debug representations of a CNN. [15] discovered blind spots of knowledge representation of a pre-trained CNN in a weakly-supervised manner.

**Distilling neural networks into explainable models:** Furthermore, some method distilled CNN knowledge into another model with interpretable features for explanations. [33] distilled knowledge of a neural network into an additive model to explain the knowledge inside the network. [47] roughly represented the rationale of each CNN prediction using a semantic tree structure. Each node in the tree represented a decision-making mode of the CNN. Similarly, [41] used a semantic graph to summarize and explain all part knowledge hidden inside conv-layers of a CNN.

## 2.2 Learning interpretable feature representations

Unlike the diagnosis and visualization of pre-trained CNNs, some approaches were developed to learn meaningful feature representations in recent years. Automatically learning interpretable feature representations without additional human annotations proposes new challenges to state-of-the-art algorithms. For example, [22] required people to label dimensions of the input that were related to each output, in order to learn a better model. Hu *et al.* [9] designed logic rules to regularize network outputs during the learning process. Sabour *et al.* [23] proposed a capsule model, where each feature dimension of a capsule may represent a specific meaning. Similarly, we invent a generic filter loss to regularize the representation of a filter to improve its interpretability.

In addition, unlike the visualization methods (*e.g.* the Grad-CAM method [24]) using a single saliency map for visualization, our interpretable CNN disentangles feature representations and uses different filters to represent the different object parts.

## 3 ALGORITHM

Given a target conv-layer of a CNN, we expect each filter in the conv-layer to be activated by a certain

object part of a certain category, while remain inactivated on images of other categories<sup>1</sup>. Let  $\mathbf{I}$  denote a set of training images, where  $\mathbf{I}_c \subset \mathbf{I}$  represents the subset that belongs to category  $c$ , ( $c = 1, 2, \dots, C$ ). Theoretically, we can use different types of losses to learn CNNs for multi-class classification and binary classification of a single class (*i.e.*  $c = 1$  for images of a category and  $c = 2$  for random images).

In the following paragraphs, we focus on the learning of a single filter  $f$  in a conv-layer. Fig. 2 shows the structure of our interpretable conv-layer. We add a loss to the feature map  $x$  of the filter  $f$  after the ReLU operation. The filter loss  $Loss_f$  pushes the filter  $f$  to represent a specific object part of the category  $c$  and keep silent on images of other categories. Please see Section 3.2 for the determination of the category  $c$  for the filter  $f$ . Let  $\mathbf{X} = \{x | x = f(I) \in \mathbb{R}^{n \times n}, I \in \mathbf{I}\}$  denote a set of feature maps of  $f$  after an ReLU operation *w.r.t.* different images. Given an input image  $I \in \mathbf{I}_c$ , the feature map in an intermediate layer  $x = f(I)$  is an  $n \times n$  matrix,  $x_{ij} \geq 0$ . If the target part appears, we expect the feature map  $x = f(I)$  to exclusively activate at the target part's location; otherwise, the feature map should keep inactivated.

Therefore, a high interpretability of the filter  $f$  requires a high mutual information between the feature map  $x = f(I)$  and the part location, *i.e.* the part location can roughly determine activations on the feature map  $x$ .

Accordingly, we formulate the filter loss as the minus mutual information, as follows.

$$Loss_f = -MI(\mathbf{X}; \Omega) = - \sum_{\mu \in \Omega} p(\mu) \sum_x p(x|\mu) \log \frac{p(x|\mu)}{p(x)} \quad (1)$$

where  $MI(\cdot)$  denotes the mutual information;  $\Omega = \{\mu_1, \mu_2, \dots, \mu_{n^2}\} \cup \{\mu^-\}$ . We use  $\mu_1, \mu_2, \dots, \mu_{n^2}$  to denote the  $n^2$  neural units on the feature map  $x$ , each  $\mu = [i, j] \in \Omega$ ,  $1 \leq i, j \leq n$ , corresponding to a location candidate for the target part.  $\mu^-$  denotes a dummy location for the case when the target part does not appear on the image.

Given an input image, the above loss forces each filter to match and only match one of the templates, *i.e.* making the feature map of the filter contain a single significant activation peak at most. This ensures each filter to represent a specific object part.

•  $p(\mu)$  measures the probability of the target part appearing at the location  $\mu$ . If annotations of part locations are given, then the computation of  $p(\mu)$  is simple. People can manually assign a semantic part with the filter  $f$ , and then  $p(\mu)$  can be determined using part annotations.

However, in our study, the target part of filter  $f$  is not pre-defined before the learning process. Instead,

1. To avoid ambiguity, we evaluate or visualize the semantic meaning of each filter by using the feature map after the ReLU and mask operations.

the part corresponding to  $f$  needs to be determined during the learning process. More crucially, we do not have any ground-truth annotations of the target part, which boosts the difficulty of calculating  $p(\mu)$ .

• The conditional likelihood  $p(x|\mu)$  measures the fitness between a feature map  $x$  and the part location  $\mu \in \Omega$ . In order to simplify the computation of  $p(x|\mu)$ , we design  $n^2$  templates for  $f$ ,  $\{T_{\mu_1}, T_{\mu_2}, \dots, T_{\mu_{n^2}}\}$ . As shown in Fig. 3, each template  $T_{\mu_i}$  is an  $n \times n$  matrix.  $T_{\mu_i}$  describes the ideal distribution of activations for the feature map  $x$  when the target part mainly triggers the  $i$ -th unit in  $x$ . In addition, we also design a negative template  $T^-$  corresponding to the dummy location  $\mu^-$ . The feature map can match to  $T^-$ , when the target part does not appear on the input image. In this study, the prior probability is given as  $p(\mu_i) = \frac{\alpha}{n^2}$ ,  $p(\mu^-) = 1 - \alpha$ , where  $\alpha$  is a constant prior likelihood.

Note that in Equation (1), we do not manually assign filters with different categories. Instead, we use the negative template  $\mu^-$  to help the assignment of filters. *I.e.* the negative template ensures that each filter represents a specific object part (if the input image does not belong to the target part, then the input image is supposed to match  $\mu^-$ ), which also ensures a clear assignment of filters to categories. Here, we assume two categories do not share object parts, *e.g.* eyes of dogs and those of cats do not have similar contextual appearance.

We define  $p(x|\mu)$  below, which follows a standard form widely used in [25], [38].

$$p(x|\mu) \approx p(x|T_\mu) = \frac{1}{Z_\mu} \exp [\text{tr}(x \cdot T_\mu)] \quad (2)$$

where  $Z_\mu = \sum_{x \in \mathbf{X}} \exp[\text{tr}(x \cdot T_\mu)]$ .  $\text{tr}(\cdot)$  indicates the trace of a matrix, and  $\text{tr}(x \cdot T_\mu) = \sum_{ij} x_{ij} t_{ji}$ ,  $x, T_\mu \in \mathbb{R}^{n \times n}$ .  $p(x) = \sum_{\mu} p(\mu) p(x|\mu)$ .

**Part templates:** As shown in Fig. 3, a negative template is given as  $T^- = (t_{ij}^-)$ ,  $t_{ij}^- = -\tau < 0$ , where  $\tau$  is a positive constant. A positive template corresponding to  $\mu$  is given as  $T_\mu = (t_{ij}^+)$ ,  $t_{ij}^+ = \tau \cdot \max(1 - \beta \frac{\| [i, j] - \mu \|_1}{n}, -1)$ , where  $\| \cdot \|_1$  denotes the L-1 norm distance. Note that the lowest value in a positive template is -1 instead of 0. It is because that the negative value in the template penalizes neural activations outside the domain of the highest activation peak, which ensures each filter mainly has at most a single significant activation peak.

### 3.1 Part localization & the mask layer

Given an input image  $I$ , the filter  $f$  computes a feature map  $x$  after the ReLU operation. Without ground-truth annotations of the target part for  $f$ , in this study, we determine the part location on  $x$  during the learning process. We consider the neural unit with the strongest activation  $\hat{\mu} = \arg\max_{\mu=[i, j]} x_{ij}$ ,  $1 \leq i, j \leq n$  as the target part location.

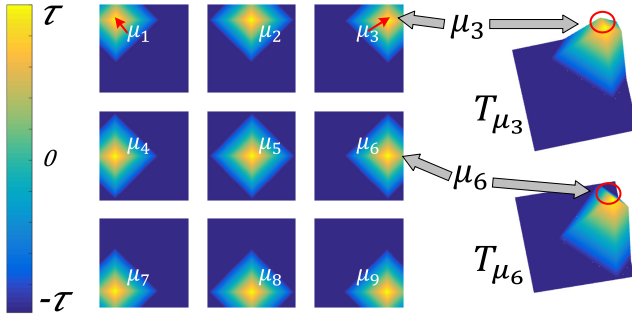


Fig. 3. Templates of  $T_{\mu_i}$ . We show a toy example of  $n = 3$ . Each template  $T_{\mu_i}$  matches to a feature map  $x$  when the target part mainly triggers the  $i$ -th unit in  $x$ . In fact, the algorithm also supports a round template based on the L-2 norm distance. Here, we use the L-1 norm distance instead to speed up the computation.

As shown in Fig. 2, we add a mask layer above the interpretable conv-layer. Based on the estimated part position  $\hat{\mu}$ , the mask layer assigns a specific mask with  $x$  to filter out noisy activations. The mask operation is separate from the filter loss in Equation (1). Our method selects the template  $T_{\hat{\mu}}$  w.r.t. the part location  $\hat{\mu}$  as the mask. We compute  $x^{\text{masked}} = \max\{x \circ T_{\hat{\mu}}, 0\}$  as the output masked feature map, where  $\circ$  denotes the Hadamard (element-wise) product. The mask operation supports gradient back-propagations.

Fig. 4 visualizes the masks  $T_{\hat{\mu}}$  chosen for different images, and compares the original and masked feature maps. The CNN selects different templates for different images.

Note that although a filter usually has much stronger neural activations on the target category than on other categories, the magnitude of neural activations is still not discriminative enough for classification. Moreover, during the testing process, people do not have ground-truth class labels of input images. Thus, to ensure stable feature extraction, our method only selects masks from the  $n^2$  positive templates  $\{T_{\mu_i}\}$  and omits the negative template  $T^-$  for all images, no matter whether or not input images contain the target part. Such an operation is conducted during the forward process for both training and testing processes.

### 3.2 Learning

We train the interpretable CNN in an end-to-end manner. During the forward-propagation process, each filter in the CNN passes its information in a bottom-up manner, just like traditional CNNs. During the back-propagation, each filter in an interpretable conv-layer receives gradients w.r.t. its feature map  $x$  from both the final task loss  $\mathbf{L}(\hat{y}_k, y_k^*)$  on the  $k$ -th sample

and the filter loss,  $\text{Loss}_f$ , as follows:

$$\frac{\partial \text{Loss}}{\partial x_{ij}} = \lambda \sum_f \frac{\partial \text{Loss}_f}{\partial x_{ij}} + \frac{1}{N} \sum_{k=1}^N \frac{\partial \mathbf{L}(\hat{y}_k, y_k^*)}{\partial x_{ij}} \quad (3)$$

where  $\lambda$  is a weight. Then, we back propagate  $\frac{\partial \text{Loss}}{\partial x_{ij}}$  to lower layers and compute gradients w.r.t. feature maps and gradients w.r.t. parameters in lower layers to update the CNN.

For implementation, gradients of  $\text{Loss}_f$  w.r.t. each element  $x_{ij}$  of feature map  $x$  are computed as follows.

$$\begin{aligned} \frac{\partial \text{Loss}_f}{\partial x_{ij}} &= \frac{1}{Z_{\hat{\mu}}} \sum_{\mu} p(\mu) t_{ij} e^{tr(x \cdot T_{\mu})} \left\{ tr(x \cdot T_{\mu}) - \log [Z_{\mu} p(x)] \right\} \\ &\approx \frac{p(\hat{\mu}) \hat{t}_{ij}}{Z_{\hat{\mu}}} e^{tr(x \cdot T_{\hat{\mu}})} \left\{ tr(x \cdot T_{\hat{\mu}}) - \log Z_{\hat{\mu}} - \log p(x) \right\} \end{aligned} \quad (4)$$

where  $T_{\hat{\mu}}$  is the target template for feature map  $x$ . If the input image  $I$  belongs to the target category of filter  $f$ , then  $\hat{\mu} = \arg\max_{\mu=[i,j]} x_{ij}$ . If image  $I$  belongs to other categories, then  $\hat{\mu} = \mu^-$ . Please see the appendix for the proof of the above equation.

Considering  $\forall \mu \in \Omega \setminus \{\hat{\mu}\}$ ,  $e^{tr(x \cdot T_{\hat{\mu}})} \gg e^{tr(x \cdot T_{\mu})}$  and  $p(\hat{\mu}) \gg p(\mu)$  after initial learning episodes, we make the above approximation to simplify the computation. Because  $Z_{\hat{\mu}}$  is computed using numerous feature maps, we can roughly treat  $Z_{\hat{\mu}}$  as a constant to compute gradients in the above equation. We gradually update the value of  $Z_{\hat{\mu}}$  during the training process. More specifically, we can use a subset of feature maps to approximate the value of  $Z_{\mu}$ , and continue to update  $Z_{\mu}$  when we receive more feature maps during the training process. Similarly, we can approximate  $p(x)$  using a subset of feature maps. We can also approximate  $p(x) = \sum_{\mu} p(\mu) p(x|\mu) = \sum_{\mu} p(\mu) \frac{\exp[tr(x \cdot T_{\mu})]}{Z_{\mu}} \approx \sum_{\mu} p(\mu) \mathbb{E}_x \frac{\exp[tr(x \cdot T_{\mu})]}{Z_{\mu}}$  without huge computation.

**Determining the target category for each filter:** We need to assign each filter  $f$  with a target category  $\hat{c}$  to approximate gradients in Equation (4). We simply assign the filter  $f$  with the category  $\hat{c}$  whose images activate  $f$  the most, i.e.  $\hat{c} = \arg\max_c \mathbb{E}_{x=f(I): I \in \mathbf{I}_c} \sum_{ij} x_{ij}$ .

### 3.3 Understanding the filter loss

The filter loss in Equation (1) can be re-written as

$$\text{Loss}_f = -H(\Omega) + H(\Omega'|\mathbf{X}) + \sum_x p(\Omega^+, x) H(\Omega^+|X=x) \quad (5)$$

where  $\Omega' = \{\mu^-, \Omega^+\}$ .  $H(\Omega) = -\sum_{\mu \in \Omega} p(\mu) \log p(\mu)$  is a constant prior entropy of part locations. Thus, the filter loss minimizes two conditional entropies,  $H(\Omega'|\mathbf{X})$  and  $H(\Omega^+|X=x)$ . Please see the appendix for the proof of the above equation.

**Low inter-category entropy:** The second term  $H(\Omega' = \{\mu^-, \Omega^+\}|\mathbf{X})$  is computed as

$$H(\Omega' = \{\mu^-, \Omega^+\}|\mathbf{X}) = - \sum_x p(x) \sum_{\mu \in \{\mu^-, \Omega^+\}} p(\mu|x) \log p(\mu|x) \quad (6)$$



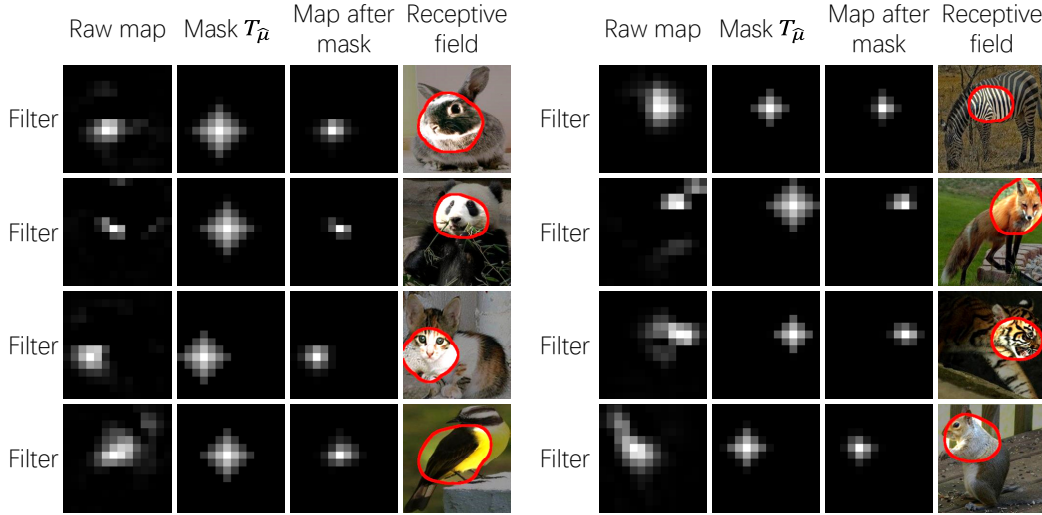


Fig. 4. Given an input image  $I$ , from left to right, we consequently show the feature map of a filter after the ReLU layer  $x$ , the assigned mask  $T_{\hat{\mu}}$ , the masked feature map  $x^{\text{masked}}$ , and the image-resolution RF of activations in  $x^{\text{masked}}$  computed by [49].

where  $\Omega^+ = \{\mu_1, \dots, \mu_n\} \subset \Omega$ ,  $p(\Omega^+|x) = \sum_{\mu \in \Omega^+} p(\mu|x)$ . We define the set of all real locations  $\Omega^+$  as a single label to represent category  $c$ . We use the dummy location  $\mu^-$  to roughly indicate matches to other categories.

This term encourages a low conditional entropy of inter-category activations, *i.e.* a well-learned filter  $f$  needs to be exclusively activated by a certain category  $c$  and keep silent on other categories. We can use a feature map  $x$  of  $f$  to identify whether or not the input image belongs to category  $c$ , *i.e.*  $x$  fitting to either  $T_{\hat{\mu}}$  or  $T^-$ , without significant uncertainty.

•**Low spatial entropy:** The third term in Equation (5) is given as

$$H(\Omega^+|X=x) = \sum_{\mu \in \Omega^+} \tilde{p}(\mu|x) \log \tilde{p}(\mu|x) \quad (7)$$

where  $\tilde{p}(\mu|x) = \frac{p(\mu|x)}{p(\Omega^+|x)}$ . This term encourages a low conditional entropy of the spatial distribution of  $x$ 's activations. *I.e.* given an image  $I \in \mathbf{I}_c$ , a well-learned filter should only be activated in a single region  $\hat{\mu}$  of the feature map  $x$ , instead of being repetitively triggered at different locations.

## 4 EXPERIMENTS

In experiments, we applied our method to modify four types of CNNs with various structures into interpretable CNNs and learned interpretable CNNs based on three benchmark datasets, in order to demonstrate the broad applicability. We learned interpretable CNNs for binary classification of a single category and multi-category classification. We used different techniques to visualize the knowledge encoded in interpretable filters, in order to qualitatively illustrate semantic meanings of these filters. Furthermore, we

used two types of evaluation metrics, *i.e.* the object-part interpretability and the location instability, to measure the clarity of the meaning of a filter.

Our experiments showed that an interpretable filter in our interpretable CNN usually consistently represented the same part through different input images, while a filter in an ordinary CNN mainly described a mixture of semantics.

We chose three benchmark datasets with part annotations for training and testing, including the ILSVRC 2013 DET Animal-Part dataset [42], the CUB200-2011 dataset [35], and the VOC Part dataset [4]. These datasets provide ground-truth bounding boxes of entire objects. For landmark annotations, the ILSVRC 2013 DET Animal-Part dataset [42] contains ground-truth bounding boxes of heads and legs of 30 animal categories. The CUB200-2011 dataset [35] contains a total of 11.8K bird images of 200 species, and the dataset provides center positions of 15 bird landmarks. The VOC Part dataset [4] contains ground-truth part segmentations of 107 object landmarks in six animal categories.

We used these datasets, because they contain ground-truth annotations of object landmarks<sup>2</sup> (parts) to evaluate the semantic clarity of each filter. As mentioned in [4], [42], animals usually consist of non-rigid parts, which present considerable challenges for part localization. As in [4], [42], we selected animal categories in the three datasets for testing.

We learned interpretable filters based on structures of four typical CNNs for evaluation, including the AlexNet [13], the VGG-M [29], the VGG-S [29], the

2. To avoid ambiguity, a landmark is referred to as the *central position* of a semantic part (a part with an explicit name, *e.g.* a head, a tail). In contrast, the part corresponding to a filter does not have an explicit name.

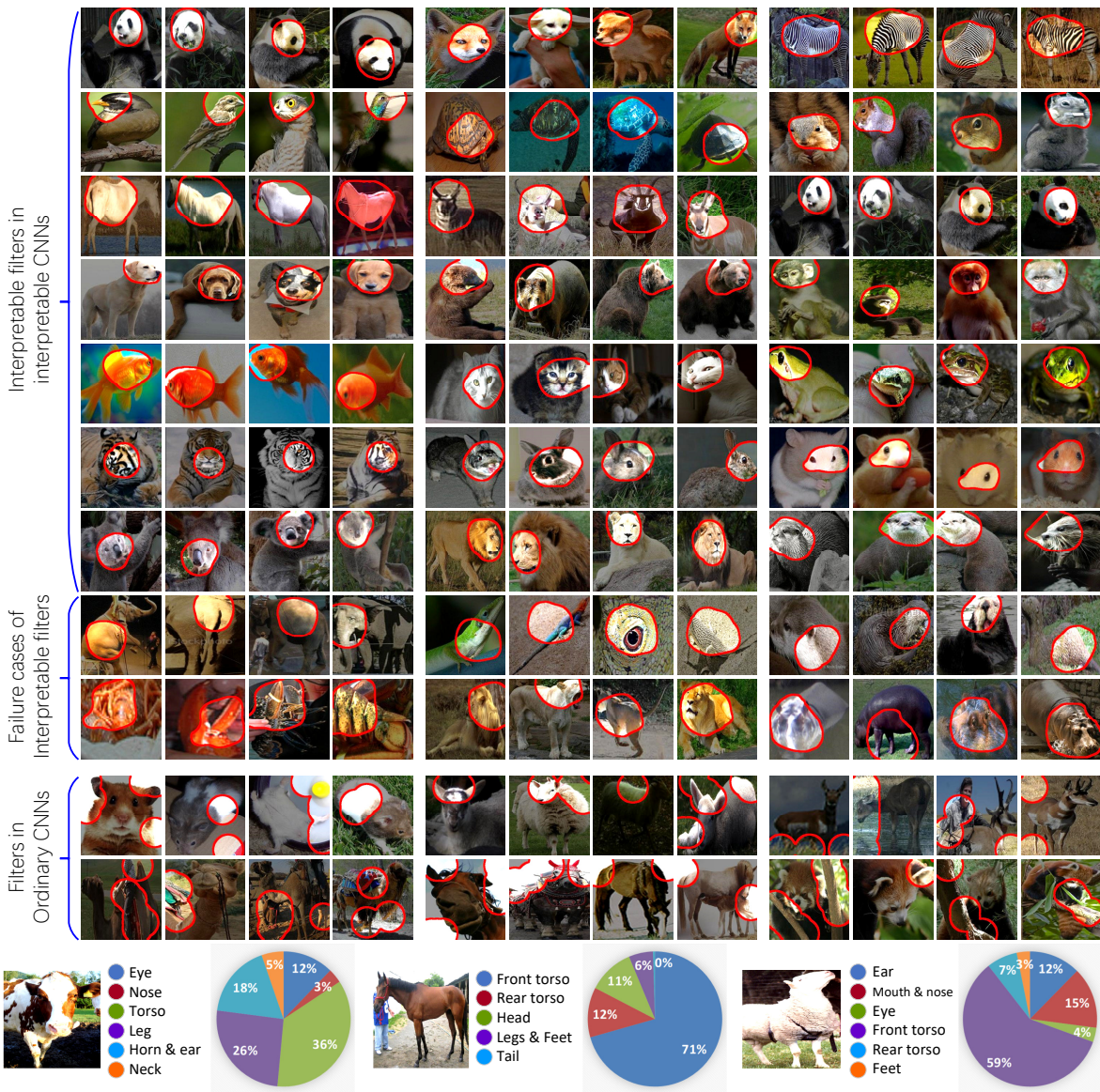


Fig. 5. Visualization of filters in top conv-layers (top) and quantitative contribution of object parts to the prediction (bottom). (top) We used [49] to estimate the image-resolution receptive field of activations in a feature map to visualize a filter's semantics. Each group of four feature maps for a category are computed using the same interpretable filter. These images show that each interpretable filter is consistently activated by the same object part through different images. Four rows visualize filters in interpretable CNNs, and two rows correspond to filters in ordinary CNNs. (bottom) The clear disentanglement of object-part representations help people to quantify the contribution of different object parts to the network prediction. We show the explanation for part contribution, which was generated by the method of [3].

VGG-16 [29]. Note that skip connections in residual networks [8] make a single feature map contain patterns of different filters. Thus, we did not use residual networks for testing to simplify the evaluation. Given a CNN, all filters in the top conv-layer were set as interpretable filters. Then, we inserted another conv-layer with  $M$  filters above the top conv-layer, which did not change the size of output feature maps. *I.e.* we set  $M = 512$  for the VGG-16, VGG-M, and VGG-S networks, and  $M = 256$  for the AlexNet. Filters in the

new conv-layer were also interpretable filters. Each filter was a  $3 \times 3 \times M$  tensor with a bias term.

*Implementation details:* We set parameters as  $\tau = \frac{0.5}{n^2}$ ,  $\alpha = \frac{n^2}{1+n^2}$ , and  $\beta \approx 4$ .  $\beta$  was updated during the learning process. We set a decreasing weight for filter losses, *i.e.*  $\lambda \propto \frac{1}{t} \mathbb{E}_{x \in \mathbf{X}} \max_{i,j} x_{ij}$  for the  $t$ -th epoch. We initialized fully-connected (FC) layers and the new conv-layer, but we loaded parameters of the lower conv-layers from a CNN that was pre-trained using [13], [29]. We then fine-tuned parameters of all layers



Image Heatmap Image Heatmap Image Heatmap Image Heatmap Image Heatmap Image Heatmap

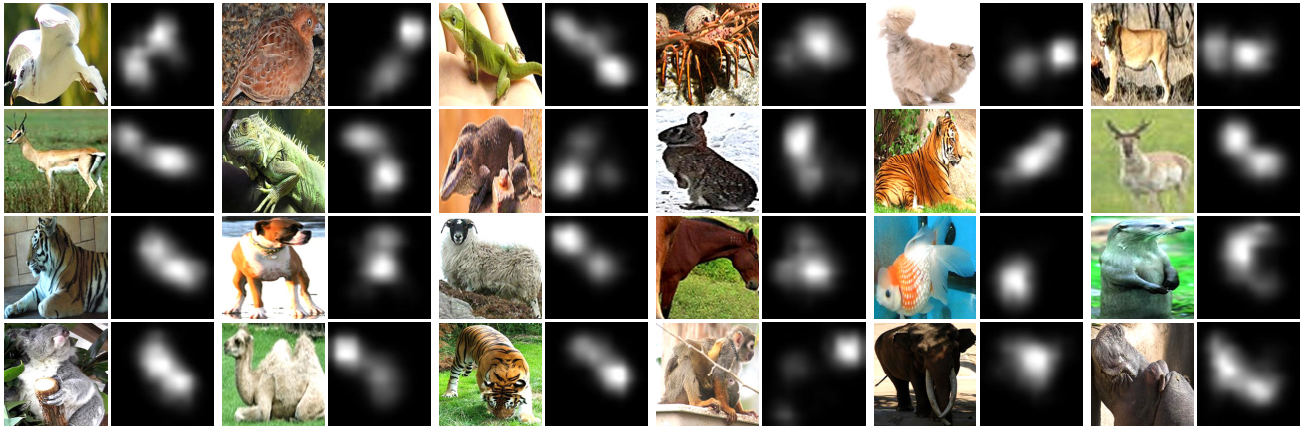


Fig. 6. Heatmaps for distributions of object parts that are encoded in interpretable filters. We use all filters in the top conv-layer to compute the heatmap. Interpretable filters usually selectively modeled distinct object parts of a category and ignored other parts.

in the interpretable CNN using training images in the dataset. To enable a fair comparison, when we learned the traditional CNN as a baseline, we also initialized FC layers of the traditional CNN, used pre-trained parameters in conv-layers, and then fine-tuned the CNN.

#### 4.1 Experiments

*Binary classification of a single category:* We learned interpretable CNNs based on above four types of network structures to classify each animal category in above three datasets. We also learned ordinary CNNs using the same data for comparison. We used the logistic log loss for binary classification of a single category from random images. We followed experimental settings in [41], [42] to crop objects of the target category as positive samples. Images of other categories were regarded as negative samples.

*Multi-category classification:* We learned interpretable CNNs to classify the six animal categories in the VOC Part dataset [4] and also learned interpretable CNNs to classify the thirty categories in the ILSVRC 2013 DET Animal-Part dataset [42]. In experiments, we tried both the softmax log loss and the logistic log loss<sup>3</sup> for multi-category classification.

#### 4.2 Qualitative Visualization of filters

We followed the method proposed by Zhou *et al.* [49] to compute the receptive fields (RFs) of neural activations of a filter. We used neural activations after ReLU and mask operations and scaled up RFs to the image resolution. As discussed in [2], the traditional

idea of directly propagating the theoretical receptive field of a neural unit in a feature map back to the image plane cannot accurately reflect the real image-resolution RF of the neural unit (*i.e.* the image region that contributes most to the score of the neural unit). Therefore, we used the method of [49] to compute real RFs.

Studies in both [49] and [2] have introduced methods to compute real RFs of neural activations on a given feature map. For ordinary CNNs, we simply used a round RF for each neural activation. We overlapped all activated RFs in a feature map to compute the final RF of the feature map.

Fig. 5 shows RFs<sup>4</sup> of filters in top conv-layers of CNNs, which were trained for binary classification of a single category. Filters in interpretable CNNs were mainly activated by a certain object part, whereas feature maps of ordinary CNNs after ReLU operations usually represented various object parts and textures. The clear disentanglement of object-part representations can help people to quantify the contribution of different object parts to the network prediction. Fig. 5 shows the explanation for part contribution, which was generated by the method of [3].

We found that interpretable CNNs usually encoded head patterns of animals in its top conv-layer for classification, although no part annotations were used to train the CNN. We can understand such results from the perspective of the information bottleneck [37] as follows. (i) Our interpretable filters selectively encode the most distinct parts of each category (*i.e.* the head for most categories), which minimizes the conditional entropy of the final classification given feature maps of a conv-layer. (ii) Each interpretable filter represents a specific part of an object, which minimizes the mutual information between the input image and middle-layer feature maps. The interpretable CNN

3. We considered the output  $y_c$  for each category  $c$  independent to outputs for other categories, thereby a CNN making multiple independent binary classifications of different categories for each image. Table 7 reported the average accuracy of the multiple classification outputs of an image.



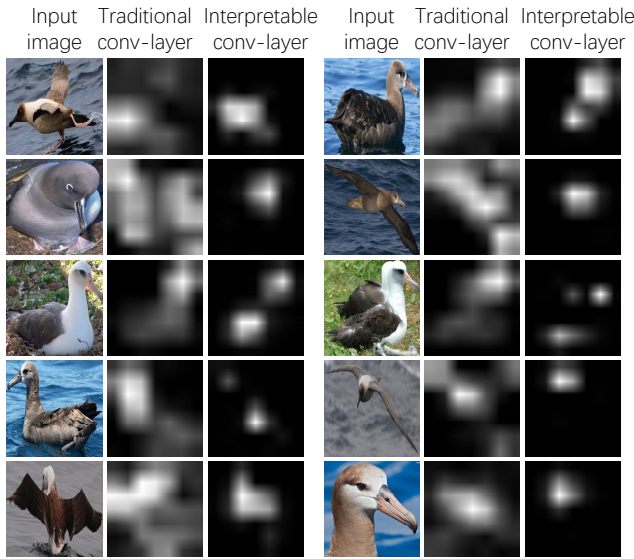


Fig. 7. Grad-CAM visualizations [24] of the traditional conv-layer and the interpretable conv-layer. Unlike the traditional conv-layer, the interpretable conv-layer usually selectively modeled distinct object parts of a category and ignored other parts.

“forgets” as much irrelevant information as possible.

In addition to the visualization of RFs, we also visualized heatmaps for part distributions and the grad-CAM attention map of an interpretable conv-layer. Fig. 6 shows heatmaps for distributions of object parts that were encoded in interpretable filters. Fig. 7 compares grad-CAM visualizations [24] of an interpretable conv-layer and those of a traditional conv-layer. We chose the top conv-layer of the traditional VGG-16 net and the top conv-layer of the interpretable VGG-16 net for visualization. Interpretable filters usually selectively modeled distinct object parts of a category and ignored other parts.

### 4.3 Quantitative evaluation of part interpretability

Filters in low conv-layers usually represent simple patterns or object details, whereas those in high conv-layers are more likely to describe large-scale parts. Therefore, in experiments, we used the following two metrics to evaluate the clarity of part semantics of the top conv-layer of a CNN.

#### 4.3.1 Evaluation metric: part interpretability

The metric was originally proposed by Bau *et al.* [2] to measure the object-part interpretability of filters. For each filter  $f$ ,  $\mathbf{X}$  denotes a set of feature maps after ReLU/mask operations on different input images. Then, the distribution of activation scores over all positions in all feature maps was computed. [2] set a threshold  $T_f$  such that  $p(x_{ij} > T_f) = 0.005$  to select strongest activations from all positions  $[i, j]$  from  $x \in \mathbf{X}$  as valid activations for  $f$ 's semantics.

Then, image-resolution RFs of valid neural activations of each input image  $I$  were computed<sup>4</sup>. The RFs on image  $I$ , termed  $S_f^I$ , corresponded to part regions of  $f$ .

The fitness between the filter  $f$  and the  $k$ -th part on image  $I$  was reported as the intersection-over-union score  $IoU_{f,k}^I = \frac{\|S_f^I \cap S_k^I\|}{\|S_f^I \cup S_k^I\|}$ , where  $S_k^I$  represents the ground-truth mask of the  $k$ -th part on image  $I$ . Given an image  $I$ , the filter  $f$  was associated with the  $k$ -th part if  $IoU_{f,k}^I > 0.2$ . The criterion  $IoU_{f,k}^I > 0.2$  was stricter than  $IoU_{f,k}^I > 0.04$  in [2], because object-part semantics usually needs a stricter criterion than textural semantics and color semantics in [2]. The average probability of the  $k$ -th part being associating with the filter  $f$  was reported as  $P_{f,k} = \mathbb{E}_{I:\text{with } k\text{-th part}} \mathbf{1}(IoU_{f,k}^I > 0.2)$ . Note that a single filter may be associated with multiple object parts in an image. The highest probability of part association for each filter was used as the interpretability of filter  $f$ , i.e.  $P_f = \max_k P_{f,k}$ .

For the binary classification of a single category, we used testing images of the target category to evaluate the feature interpretability. In the VOC Part dataset [4], four parts were chosen for the *bird* category. We merged segments of the head, beak, and 1/r-eyes as the head part, merged segments of the torso, neck, and 1/r-wings as the torso part, merged segments of 1/r-legs/feet as the leg part, and used the tail segment as the fourth part. We used five parts for both the *cat* category and the *dog* category. We merged segments of the head, 1/r-eyes, 1/r-ears, and nose as the head part, merged segments of the torso and neck as the torso part, merged segments of frontal 1/r-legs/paws as the frontal legs, merged segments of back 1/r-legs/paws as the back legs, and used the tail as the fifth part. Part definitions for the *cow*, *horse*, and *sheep* category were similar those for the *cat* category, except for that we omitted the tail part of these categories. In particular, we added 1/r-horn segments of the horse to the head part. The average part interpretability  $P_f$  over all filters was computed for evaluation.

For the multi-category classification, we first determined the target category  $\hat{c}$  for each filter  $f$  i.e.  $\hat{c} = \operatorname{argmax}_c \mathbb{E}_{x=f(I): I \in \mathbf{I}_c} \sum_{i,j} x_{ij}$ . Then, we computed  $f$ 's object-part interpretability using images of the target category  $\hat{c}$  by following above instructions.

4. [49] computes the RF when the filter represents an object part. Fig. 5 used RFs computed by [49] to visualize filters. However, when a filter in an ordinary CNN does not have consistent contours, it is difficult for [49] to align different images to compute an average RF. Thus, for ordinary CNNs, we simply used a round RF for each valid activation. We overlapped all activated RFs in a feature map to compute the final RF as mentioned in [2]. For a fair comparison, in Section 4.3.1, we uniformly applied these RFs to both interpretable CNNs and ordinary CNNs.

	bird	cat	cow	dog	horse	sheep	Avg.
AlexNet	0.332	0.363	0.340	0.374	0.308	0.373	0.348
AlexNet, interpretable	<b>0.770</b>	<b>0.565</b>	<b>0.618</b>	<b>0.571</b>	<b>0.729</b>	<b>0.669</b>	<b>0.654</b>
VGG-16	0.519	0.458	0.479	0.534	0.440	0.542	0.495
VGG-16, interpretable	<b>0.818</b>	<b>0.653</b>	<b>0.683</b>	<b>0.900</b>	<b>0.795</b>	<b>0.772</b>	<b>0.770</b>
VGG-M	0.357	0.365	0.347	0.368	0.331	0.373	0.357
VGG-M, interpretable	<b>0.821</b>	<b>0.632</b>	<b>0.634</b>	<b>0.669</b>	<b>0.736</b>	<b>0.756</b>	<b>0.708</b>
VGG-S	0.251	0.269	0.235	0.275	0.223	0.287	0.257
VGG-S, interpretable	<b>0.526</b>	<b>0.366</b>	<b>0.291</b>	<b>0.432</b>	<b>0.478</b>	0.251	<b>0.390</b>

TABLE 1

Part interpretability of filters in CNNs for binary classification of a single category based on the VOC Part dataset [4].

Network	Logistic log loss <sup>3</sup>	Softmax log loss
VGG-16	0.710	0.723
interpretable	<b>0.938</b>	<b>0.897</b>
VGG-M	0.478	0.502
interpretable	<b>0.770</b>	<b>0.734</b>
VGG-S	0.479	0.435
interpretable	<b>0.572</b>	<b>0.601</b>

TABLE 2

Part interpretability of filters in CNNs that are trained for multi-category classification based on the VOC Part dataset [4]. Filters in our interpretable CNNs exhibited significantly better part interpretability than ordinary CNNs in all comparisons.

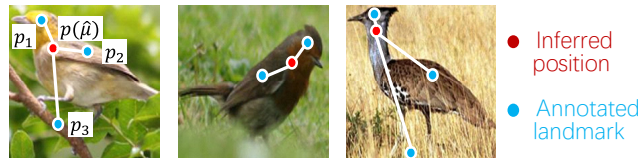


Fig. 8. Notation for computing the location instability.

#### 4.3.2 Evaluation metric: location instability

The second metric measures the instability of part locations, which was used in [41], [46]. It is assumed that if  $f$  consistently represented the same object part through different objects, then distances between the inferred part  $\hat{\mu}$  and some ground-truth landmarks<sup>2</sup> should keep stable among different objects. For example, if  $f$  represented the shoulder part without ambiguity, then the distance between the inferred position and the head will not change a lot among different objects.

Therefore, the deviation of the distance between the inferred position  $\hat{\mu}$  and a specific ground-truth landmark among different images was computed. The location  $\hat{\mu}$  was inferred as the neural unit with the highest activation on  $f$ 's feature map. We reported the average deviation *w.r.t.* different landmarks as the location instability of  $f$ .

Please see Fig. 8. Given an input image  $I$ ,  $d_I(p_k, \hat{\mu}) = \frac{\|p_k - \hat{\mu}\|}{\sqrt{w^2 + h^2}}$  denotes the normalized distance

Neural network	Avg. location instability
AlexNet	0.150
AlexNet+ordinary layer	0.118
AlexNet, interpretable	<b>0.070</b>
VGG-16	0.137
VGG-16+ordinary layer	0.097
VGG-16, interpretable	<b>0.076</b>
VGG-M	0.148
VGG-M+ordinary layer	0.107
VGG-M, interpretable	<b>0.065</b>
VGG-S	0.148
VGG-S+ordinary layer	0.103
VGG-S, interpretable	<b>0.073</b>

TABLE 5

Location instability of filters ( $\mathbb{E}_{f,k}[D_{f,k}]$ ) in CNNs for binary classification of a single category using the CUB200-2011 dataset.

between the inferred part and the  $k$ -th landmark  $p_k$ , where  $p(\hat{\mu})$  is referred to as the center of the unit  $\hat{\mu}$ 's RF.  $\sqrt{w^2 + h^2}$  measures the diagonal length of  $I$ .  $D_{f,k} = \sqrt{\text{var}_I[d_I(p_k, \hat{\mu})]}$  is termed as the *relative location deviation* of filter  $f$  *w.r.t.* the  $k$ -th landmark, where  $\text{var}_I[d_I(p_k, \hat{\mu})]$  is the variation of  $d_I(p_k, \hat{\mu})$ . Because each landmark could not appear in all testing images, for each filter  $f$ , the metric only used inference results on top-ranked 100 images with the highest inference scores to compute  $D_{f,k}$ . In this way, the average of relative location deviations of all the filters in a conv-layer *w.r.t.* all  $K$  landmarks, *i.e.*  $\mathbb{E}_f \mathbb{E}_{k=1}^K D_{f,k}$ , was reported as the location instability of  $f$ .

We used the most frequent object parts as landmarks to measure the location instability. For the ILSVRC 2013 DET Animal-Part dataset [42], we used the *head* and *frontal legs* of each category as landmarks for evaluation. For the VOC Part dataset [4], we selected the *head*, *neck*, and *torso* of each category as landmarks. For the CUB200-2011 dataset [35], we used the *head*, *back*, *tail* of birds as landmarks.

In particular, for multi-category classification, we first determined the target category of for each filter  $f$  and then computed the relative location deviation  $D_{f,k}$  using landmarks of  $f$ 's target category. Because

	gold.	bird	frog	turt.	liza.	koala	lobs.	dog	fox	cat	lion
AlexNet	0.161	0.167	0.152	0.153	0.175	0.128	0.123	0.144	0.143	0.148	0.137
AlexNet+ordinary layer	0.154	0.157	0.143	0.146	0.170	0.120	0.118	0.127	0.117	0.136	0.120
AlexNet, interpretable	<b>0.084</b>	<b>0.095</b>	<b>0.090</b>	<b>0.107</b>	<b>0.097</b>	<b>0.079</b>	<b>0.077</b>	<b>0.093</b>	<b>0.087</b>	<b>0.095</b>	<b>0.084</b>
VGG-16	0.153	0.156	0.144	0.150	0.170	0.127	0.126	0.143	0.137	0.148	0.139
VGG-16+ordinary layer	0.136	0.127	0.120	0.136	0.147	0.108	0.111	0.111	0.097	0.134	0.102
VGG-16, interpretable	<b>0.076</b>	<b>0.099</b>	<b>0.086</b>	<b>0.115</b>	<b>0.113</b>	<b>0.070</b>	<b>0.084</b>	<b>0.077</b>	<b>0.069</b>	<b>0.086</b>	<b>0.067</b>
VGG-M	0.161	0.166	0.151	0.153	0.176	0.128	0.125	0.145	0.145	0.150	0.140
VGG-M+ordinary layer	0.147	0.144	0.135	0.142	0.159	0.114	0.115	0.119	0.111	0.128	0.114
VGG-M, interpretable	<b>0.088</b>	<b>0.088</b>	<b>0.089</b>	<b>0.108</b>	<b>0.099</b>	<b>0.080</b>	<b>0.074</b>	<b>0.090</b>	<b>0.082</b>	<b>0.103</b>	<b>0.079</b>
VGG-S	0.158	0.166	0.149	0.151	0.173	0.127	0.124	0.143	0.142	0.148	0.138
VGG-S+ordinary layer	0.150	0.132	0.133	0.138	0.156	0.113	0.111	0.110	0.104	0.125	0.112
VGG-S, interpretable	<b>0.087</b>	<b>0.101</b>	<b>0.093</b>	<b>0.107</b>	<b>0.096</b>	<b>0.084</b>	<b>0.078</b>	<b>0.091</b>	<b>0.082</b>	<b>0.101</b>	<b>0.082</b>
	tiger	bear	rabb.	hams.	squi.	horse	zebra	swine	hippo.	catt.	sheep
AlexNet	0.142	0.144	0.148	0.128	0.149	0.152	0.154	0.141	0.141	0.144	0.155
AlexNet+ordinary layer	0.123	0.133	0.136	0.112	0.145	0.149	0.142	0.137	0.139	0.141	0.149
AlexNet, interpretable	<b>0.090</b>	<b>0.095</b>	<b>0.095</b>	<b>0.077</b>	<b>0.095</b>	<b>0.098</b>	<b>0.084</b>	<b>0.091</b>	<b>0.089</b>	<b>0.097</b>	<b>0.101</b>
VGG-16	0.144	0.143	0.146	0.125	0.150	0.150	0.153	0.141	0.140	0.140	0.150
VGG-16+ordinary layer	0.127	0.112	0.119	0.100	0.112	0.134	0.140	0.126	0.126	0.131	0.135
VGG-16, interpretable	<b>0.097</b>	<b>0.081</b>	<b>0.079</b>	<b>0.066</b>	<b>0.065</b>	<b>0.106</b>	<b>0.077</b>	<b>0.094</b>	<b>0.083</b>	<b>0.102</b>	<b>0.097</b>
VGG-M	0.145	0.144	0.150	0.128	0.150	0.151	0.158	0.140	0.140	0.143	0.155
VGG-M+ordinary layer	0.124	0.131	0.134	0.108	0.132	0.138	0.141	0.133	0.131	0.135	0.142
VGG-M, interpretable	<b>0.089</b>	<b>0.101</b>	<b>0.097</b>	<b>0.082</b>	<b>0.095</b>	<b>0.095</b>	<b>0.080</b>	<b>0.095</b>	<b>0.084</b>	<b>0.092</b>	<b>0.094</b>
VGG-S	0.142	0.143	0.148	0.128	0.146	0.149	0.155	0.139	0.140	0.141	0.155
VGG-S+ordinary layer	0.117	0.127	0.127	0.105	0.122	0.136	0.137	0.133	0.131	0.130	0.143
VGG-S, interpretable	<b>0.089</b>	<b>0.097</b>	<b>0.091</b>	<b>0.076</b>	<b>0.098</b>	<b>0.096</b>	<b>0.080</b>	<b>0.092</b>	<b>0.088</b>	<b>0.094</b>	<b>0.101</b>
	ante.	camel	otter	arma.	monk.	elep.	red pa.	gia.pa.			<b>Av.</b>
AlexNet	0.147	0.153	0.159	0.160	0.139	0.125	0.140	0.125			<b>0.146</b>
AlexNet+ordinary layer	0.148	0.143	0.145	0.151	0.125	0.116	0.127	0.102			<b>0.136</b>
AlexNet, interpretable	<b>0.085</b>	<b>0.102</b>	<b>0.104</b>	<b>0.095</b>	<b>0.090</b>	<b>0.085</b>	<b>0.084</b>	<b>0.073</b>			<b>0.091</b>
VGG-16	0.144	0.149	0.154	0.163	0.136	0.129	0.143	0.125			<b>0.144</b>
VGG-16+ordinary layer	0.122	0.121	0.134	0.143	0.108	0.110	0.115	0.102			<b>0.121</b>
VGG-16, interpretable	<b>0.091</b>	<b>0.105</b>	<b>0.093</b>	<b>0.100</b>	<b>0.074</b>	<b>0.084</b>	<b>0.067</b>	<b>0.063</b>			<b>0.085</b>
VGG-M	0.146	0.154	0.160	0.161	0.140	0.126	0.142	0.127			<b>0.147</b>
VGG-M+ordinary layer	0.130	0.135	0.140	0.150	0.120	0.112	0.120	0.106			<b>0.130</b>
VGG-M, interpretable	<b>0.077</b>	<b>0.104</b>	<b>0.102</b>	<b>0.093</b>	<b>0.086</b>	<b>0.087</b>	<b>0.089</b>	<b>0.068</b>			<b>0.090</b>
VGG-S	0.143	0.154	0.158	0.157	0.140	0.125	0.139	0.125			<b>0.145</b>
VGG-S+ordinary layer	0.125	0.133	0.135	0.147	0.119	0.111	0.118	0.100			<b>0.126</b>
VGG-S, interpretable	<b>0.077</b>	<b>0.102</b>	<b>0.105</b>	<b>0.094</b>	<b>0.090</b>	<b>0.086</b>	<b>0.078</b>	<b>0.072</b>			<b>0.090</b>

TABLE 3

Location instability of filters ( $\mathbb{E}_{f,k}[D_{f,k}]$ ) in CNNs that are trained for the binary classification of a single category using the ILSVRC 2013 DET Animal-Part dataset [42]. Filters in our interpretable CNNs exhibited significantly lower localization instability than ordinary CNNs.

filters in baseline CNNs did not exclusively represent a single category, we simply assigned filter  $f$  with the category whose landmarks can achieve the lowest location deviation to simplify the computation. *I.e.* for a baseline CNN, we used  $\mathbb{E}_f \min_c \mathbb{E}_{k \in Part_c} D_{f,k}$  to evaluate the location instability, where  $Part_c$  denotes the set of part indexes belonging to category  $c$ .

#### 4.3.3 Comparisons between metrics of filter interpretability and location instability

Although the filter interpretability [2] and the location instability [46] are the two most state-of-the-art met-

rics to evaluate the interpretability of a convolution filter, these metrics still have some limitations.

Firstly, the filter interpretability [2] assumes that the feature map of an automatically learned filter should well match the ground-truth segment of a semantic part (with an explicit part name), an object, or a texture. For example, it assumes that a filter may represent the exact segment of the head part. However, without ground-truth annotations of object parts or textures for supervision, there is no mechanism to assign explicit semantic meanings with filters during the learning process. In most cases, filters in



	bird	cat	cow	dog	horse	sheep	Avg.
AlexNet	0.153	0.131	0.141	0.128	0.145	0.140	<b>0.140</b>
AlexNet+ordinary layer	0.147	0.125	0.139	0.112	0.146	0.143	<b>0.136</b>
AlexNet, interpretable w/o filter loss	0.091	0.090	0.091	0.089	0.086	0.088	<b>0.089</b>
AlexNet, interpretable	0.090	0.089	0.090	0.088	0.087	0.088	<b>0.088</b>
VGG-16	0.145	0.133	0.146	0.127	0.143	0.143	<b>0.139</b>
VGG-16+ordinary layer	0.125	0.121	0.137	0.102	0.131	0.137	<b>0.125</b>
VGG-16, interpretable w/o filter loss	0.099	0.087	0.102	0.078	0.096	0.101	<b>0.094</b>
VGG-16, interpretable	0.101	0.098	0.105	0.074	0.097	0.100	<b>0.096</b>
VGG-M	0.152	0.132	0.143	0.130	0.145	0.141	<b>0.141</b>
VGG-M+ordinary layer	0.142	0.120	0.139	0.115	0.141	0.142	<b>0.133</b>
VGG-M, interpretable w/o filter loss	0.089	0.095	0.091	0.086	0.086	0.091	<b>0.090</b>
VGG-M, interpretable	0.086	0.094	0.090	0.087	0.084	0.084	<b>0.088</b>
VGG-S	0.152	0.131	0.141	0.128	0.144	0.141	<b>0.139</b>
VGG-S+ordinary layer	0.137	0.115	0.133	0.107	0.133	0.138	<b>0.127</b>
VGG-S, interpretable w/o filter loss	0.093	0.096	0.093	0.086	0.083	0.090	<b>0.090</b>
VGG-S, interpretable	0.089	0.092	0.092	0.087	0.086	0.088	<b>0.089</b>

TABLE 4

Location instability of filters ( $\mathbb{E}_{f,k}[D_{f,k}]$ ) in CNNs that are trained for binary classification of a single category using the VOC Part dataset [4]. Filters in our interpretable CNNs exhibited significantly lower localization instability than ordinary CNNs.

	ILSVRC Part [42]	VOC Part [4]	
	Logistic log loss <sup>3</sup>	Logistic log loss <sup>3</sup>	Softmax log loss
VGG-16	–	0.128	0.142
ordinary layer	–	0.096	0.099
interpretable	–	<b>0.073</b>	<b>0.075</b>
VGG-M	0.167	0.135	0.137
ordinary layer	–	0.117	0.107
interpretable	<b>0.096</b>	<b>0.083</b>	<b>0.087</b>
VGG-S	0.131	0.138	0.138
ordinary layer	–	0.127	0.099
interpretable	<b>0.083</b>	<b>0.078</b>	<b>0.082</b>

TABLE 6

Location instability of filters ( $\mathbb{E}_{f,k}[D_{f,k}]$ ) in CNNs that are trained for multi-category classification. Filters in our interpretable CNNs exhibited significantly lower localization instability than ordinary CNNs in all comparisons.

an interpretable CNN (as well as a few filters in traditional CNNs) may describe a specific object part without explicit names, *e.g.* the region of both the head and neck or the region connecting the torso and the tail. Therefore, in both [2] and [46], people did not require the inferred object region to describe the exact segment of a semantic part, and simply set a relatively loose criterion  $IoU_{f,k}^I > 0.04$  or  $0.2$  to compute the filter interpretability.

Secondly, the location instability was proposed in [46]. The location instability of a filter is evaluated using the average deviation of distances between the inferred position and some ground-truth landmarks. There is also an assumption for this evaluation metric, *i.e.* the distance between an inferred part and a specific landmark should not change a lot through different images. As a result, people cannot set landmarks as the head and the tail of a snake, because the distance between different parts of a snake continuously

changes when the snake moves.

Generally speaking, there are two advantages to use the location instability for evaluation:

- The computation of the location instability [46] is independent to the size of the receptive field (RF) of a neural activation. This solves a big problem with the evaluation of filter interpretability, *i.e.* state-of-the-art methods of computing a neural activation’s image-resolution RFs (*e.g.* [49]) can only provide an approximate scale of the RF. The metric of location instability only uses central positions of part inferences of a filter, rather than use the entire inferred part segment, for evaluation. Thus, the location instability is a robust metric to evaluate the object-part interpretability of a filter.
- The location instability allows a filter to represent an object part without an explicit name (a half of the head).

Nevertheless, the evaluation metric for filter interpretability is still an open problem.

#### 4.3.4 Robustness to adversarial attacks

In this experiment, we applied adversarial attacks [32] to both original CNNs and interpretable CNNs. The CNNs were learned to classify birds in the CUB200-2011 dataset and random images. Table 15 compares

the average adversarial distortion  $\sqrt{\frac{\sum(I'_i - I_i)^2}{n}}$  of the adversarial signal among all images between original CNNs and interpretable CNNs, where  $I$  represents the input image while  $I'$  denotes the adversarial counterpart. Because interpretable CNNs exclusively encoded object-part patterns and ignored textures, original CNNs usually exhibited stronger robustness to adversarial attacks than interpretable CNNs.

#### 4.3.5 Experimental results and analysis

Feature interpretability of different CNNs is evaluated in Tables 1, 2, 3, 4, 5, and 6. Tables 1 and 2 show

results based on the metric in [2]. Tables 3, 4, and 5 list location instability of CNNs for binary classification of a single category. Table 6 reports location instability of CNNs that were learned for multi-category classification.

We compared our interpretable CNNs with two types of CNNs, *i.e.* the original CNN, the CNN with an additional conv-layer on the top (termed *AlexNet/VGG-16/VGG-M/VGG-S+ordinary layer*). To construct the CNN with a new conv-layer, we put a new conv-layer on the top of conv-layer. The filter size of the new conv-layer was  $3 \times 3 \times \text{channel number}$ , and output feature maps of the new conv-layer were in the same size of input feature maps. Because our interpretable CNN had an additional interpretable conv-layer, we designed the baseline CNN with a new conv-layer to enable fair comparisons. Our interpretable filters exhibited significantly higher part interpretability and lower location instability than traditional filters in baseline CNNs over almost all comparisons. Table 7 reports the classification accuracy of different CNNs. Ordinary CNNs exhibited better performance in binary classification, while interpretable CNNs outperformed baseline CNNs in multi-category classification.

In addition, to prove the discrimination power of the learned filter, we further tested the average accuracy when we used the maximum activation score in a single filter's feature map as a metric for binary classification between birds in the CUB200-2011 dataset [35] and random images. In the scenario of classifying birds from random images, filters in the CNN was expected to learn the common appearance of birds, instead of summarizing knowledge from random images. Thus, we chose filters in the top conv-layer. If the maximum activation score of a filter exceeded a threshold, then we classified the input image as a bird; otherwise not. The threshold was set to the one that maximized the classification accuracy. Table 8 reports the average classification accuracy over all filters. Our interpretable filters outperformed ordinary filters.

Given a CNN for binary classification of an animal category in the VOC Part dataset [4], we manually annotated the part name corresponding to the learned filters in the CNN. Table 9 reports the ratio of interpretable filters that corresponds to each object part.

Besides, we also analyzed samples that were incorrectly classified by the interpretable CNN. We used VGG-16 networks for the binary classification of an animal category in the VOC Part dataset [4]. We annotated the object-part name corresponding to each interpretable filter in the top interpretable layer. For each false positive sample without the target category, Fig. 9 localized the image regions that were incorrectly detected as specific object parts by interpretable filters. This figure helped people understand the reason for misclassification.

Multi-category classification			
	ILSVRC Part	VOC Part	
	logistic <sup>3</sup>	logistic <sup>3</sup>	softmax
VGG-M	96.73	93.88	81.93
interpretable	<b>97.99</b>	<b>96.19</b>	<b>88.03</b>
VGG-S	96.98	94.05	78.15
interpretable	<b>98.72</b>	<b>96.78</b>	<b>86.13</b>
VGG-16	–	97.97	89.71
interpretable	–	<b>98.50</b>	<b>91.60</b>

Binary classification of a single category			
	ILSVRC Part	VOC Part	CUB200
AlexNet	<b>96.28</b>	<b>95.40</b>	<b>95.59</b>
interpretable w/o filter loss	–	93.98	–
interpretable	95.38	93.93	95.35
VGG-M	<b>97.34</b>	<b>96.82</b>	<b>97.34</b>
interpretable w/o filter loss	–	93.13	–
interpretable	95.77	94.17	96.03
VGG-S	<b>97.62</b>	<b>97.74</b>	<b>97.24</b>
interpretable w/o filter loss	–	93.83	–
interpretable	95.64	95.47	95.82
VGG-16	<b>98.58</b>	<b>98.66</b>	<b>98.91</b>
interpretable w/o filter loss	–	97.02	–
interpretable	96.67	95.39	96.51

TABLE 7

Classification accuracy based on different datasets. In the binary classification of a single category, ordinary CNNs performed better, while in multi-category classification, interpretable CNNs exhibited superior performance.

	filters in ordinary CNNs	filters in interpretable CNNs
AlexNet	68.7	<b>75.1</b>
VGG-M	69.9	<b>80.2</b>
VGG-16	72.1	<b>82.4</b>

TABLE 8

Classification accuracy based on a single filter. We reported the average accuracy to demonstrate the discrimination power of individual filters.

#### 4.4 Effects of the filter loss

In this section, we evaluated effects of the filter loss. We compared the interpretable CNN learned with the filter loss with that without the filter loss (*i.e.* only using the mask layer without the filter loss).

##### 4.4.1 Semantic purity of neural activations

We proposed a metric to measure the semantic purity of neural activations of a filter. If a filter was activated at multiple locations besides the highest peak (*i.e.* the one corresponding to the target part), we considered this filter to have low semantic purity.

The semantic purity of a filter was measured as the ratio of neural activations within the range of the mask to all neural activations of the filter. In other words, the purity of a filter indicated that whether a filter was learned to represent a single part or represent multiple parts.

	bird	cat	cow	dog	horse	sheep	Avg.
VGG-16, interpretable	<b>0.568</b>	<b>0.656</b>	<b>0.534</b>	<b>0.573</b>	<b>0.570</b>	<b>0.492</b>	<b>0.566</b>
VGG-16 + mask layer, w/o filter loss	0.385	0.373	0.342	0.435	0.382	0.303	0.370
VGG-M, interpretable	<b>0.444</b>	<b>0.616</b>	<b>0.395</b>	<b>0.540</b>	<b>0.408</b>	<b>0.387</b>	<b>0.465</b>
VGG-M + mask layer, w/o filter loss	0.230	0.341	0.185	0.295	0.250	0.215	0.253
VGG-S, interpretable	<b>0.418</b>	<b>0.437</b>	<b>0.390</b>	<b>0.398</b>	<b>0.421</b>	<b>0.369</b>	<b>0.406</b>
VGG-S + mask layer, w/o filter loss	0.224	0.312	0.165	0.234	0.208	0.161	0.217

TABLE 10

Semantic purity of neural activations of interpretable filters learned with and without the filter loss from the VOC Part dataset. Filters learned with the filter loss exhibited significantly higher semantic purity than those learned without filter loss.

	gold.	bird	frog	turt.	liza.	koala	lobs.	dog	fox	cat	lion
VGG-16, interpretable	<b>0.614</b>	<b>0.624</b>	<b>0.585</b>	<b>0.527</b>	<b>0.531</b>	<b>0.607</b>	<b>0.581</b>	<b>0.606</b>	<b>0.660</b>	<b>0.609</b>	<b>0.572</b>
VGG-16 + mask layer, w/o filter loss	0.331	0.441	0.342	0.352	0.362	0.339	0.344	0.482	0.449	0.393	0.344
VGG-M, interpretable	<b>0.533</b>	<b>0.442</b>	<b>0.385</b>	<b>0.444</b>	<b>0.430</b>	<b>0.408</b>	<b>0.430</b>	<b>0.587</b>	<b>0.564</b>	<b>0.466</b>	<b>0.572</b>
VGG-M + mask layer, w/o filter loss	0.307	0.251	0.249	0.198	0.266	0.255	0.212	0.384	0.393	0.389	0.293
VGG-S, interpretable	<b>0.454</b>	<b>0.453</b>	<b>0.389</b>	<b>0.409</b>	<b>0.389</b>	<b>0.414</b>	<b>0.398</b>	<b>0.431</b>	<b>0.456</b>	<b>0.442</b>	<b>0.395</b>
VGG-S + mask layer, w/o filter loss	0.245	0.265	0.186	0.189	0.216	0.212	0.209	0.316	0.330	0.284	0.197
	tiger	bear	rabb.	hams.	squi.	horse	zebra	swine	hippo.	catt.	sheep
VGG-16, interpretable	<b>0.583</b>	<b>0.614</b>	<b>0.535</b>	<b>0.679</b>	<b>0.677</b>	<b>0.559</b>	<b>0.532</b>	<b>0.530</b>	<b>0.564</b>	<b>0.466</b>	<b>0.572</b>
VGG-16 + mask layer, w/o filter loss	0.368	0.369	0.364	0.356	0.419	0.407	0.343	0.363	0.393	0.389	0.293
VGG-M, interpretable	<b>0.455</b>	<b>0.442</b>	<b>0.472</b>	<b>0.452</b>	<b>0.434</b>	<b>0.397</b>	<b>0.421</b>	<b>0.412</b>	<b>0.420</b>	<b>0.425</b>	<b>0.420</b>
VGG-M + mask layer, w/o filter loss	0.345	0.345	0.316	0.245	0.320	0.283	0.224	0.278	0.290	0.397	0.284
VGG-S, interpretable	<b>0.448</b>	<b>0.419</b>	<b>0.411</b>	<b>0.405</b>	<b>0.416</b>	<b>0.430</b>	<b>0.469</b>	<b>0.384</b>	<b>0.403</b>	<b>0.439</b>	<b>0.404</b>
VGG-S + mask layer, w/o filter loss	0.215	0.197	0.212	0.179	0.217	0.210	0.215	0.177	0.200	0.323	0.196
	ante.	camel	otter	arma.	monk.	elep.	red pa.	gia.pa.			Avg.
VGG-16, interpretable	<b>0.573</b>	<b>0.544</b>	<b>0.565</b>	<b>0.855</b>	<b>0.657</b>	<b>0.562</b>	<b>0.718</b>	<b>0.697</b>			<b>0.595</b>
VGG-16 + mask layer, w/o filter loss	0.440	0.367	0.315	0.321	0.377	0.333	0.391	0.383			0.372
VGG-M, interpretable	<b>0.584</b>	<b>0.435</b>	<b>0.441</b>	<b>0.419</b>	<b>0.400</b>	<b>0.399</b>	<b>0.541</b>	<b>0.468</b>			<b>0.459</b>
VGG-M + mask layer, w/o filter loss	0.411	0.315	0.232	0.175	0.225	0.173	0.334	0.343			0.285
VGG-S, interpretable	<b>0.471</b>	<b>0.386</b>	<b>0.386</b>	<b>0.436</b>	<b>0.394</b>	<b>0.408</b>	<b>0.459</b>	<b>0.471</b>			<b>0.420</b>
VGG-S + mask layer, w/o filter loss	0.388	0.202	0.184	0.183	0.188	0.185	0.262	0.252			0.226

TABLE 11

Semantic purity of neural activations of interpretable filters learned the ILSVRC 2013 DET Animal-Part dataset with and without the filter loss. Filters learned with the filter loss exhibited significantly higher semantic purity than those learned without the filter loss.

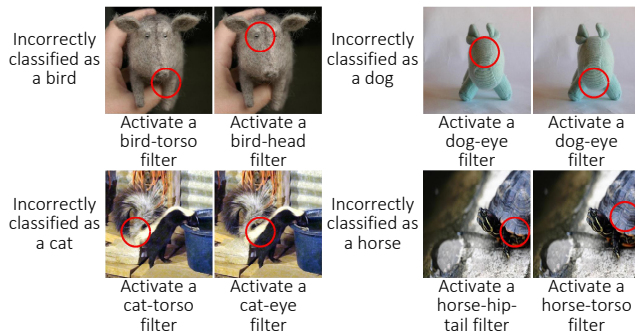


Fig. 9. Examples that were incorrectly classified by the interpretable CNN.

Let  $x \in \mathbb{R}^{n \times n}$  be neural activations of a filter (after the ReLU layer and before the mask layer). The corresponding mask was given as  $T_{\hat{\mu}} \in \mathbb{R}^{n \times n}$ . The purity of neural activations was defined as  $purity = \frac{\sum_f \sum_{i,j} \max(0, x_{ij}) \cdot \mathbf{1}(T_{\hat{\mu}, ij} > 0)}{\sum_f \sum_{i,j} \max(0, x_{ij})}$ .  $\mathbf{1}(\cdot)$  was the indicator func-

tion, which returns 1 if the condition in the braces was satisfied, and returns 0 otherwise. The purity was supposed to be higher if neural activations were more concentrated.

We compared the purity of neural activations between the interpretable CNN and the CNN trained with the mask layer but without filter loss. We constructed these CNNs based on the architectures of VGG-16, VGG-M and VGG-S, and learned the CNNs for binary classification on an animal category in the VOC Part dataset and the ILSVRC 2013 DET Animal-Part dataset. Experimental results are shown in Table 10 and Table 11. It demonstrated that filters learned with the filter loss exhibited higher semantic purity than those learned without the filter loss. The filter loss forced each filter to exclusively represented a single object part during the training process.



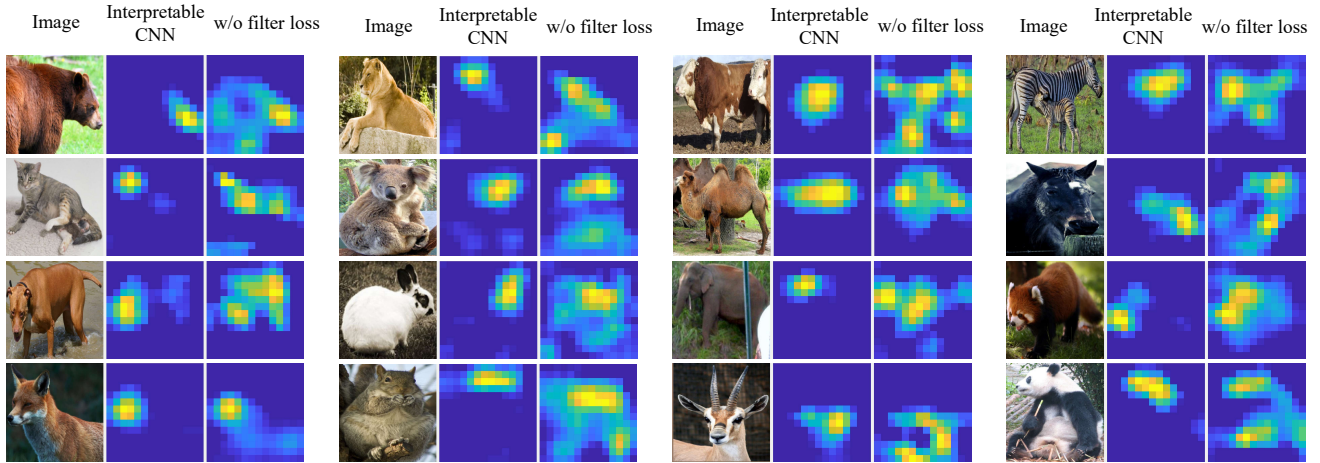


Fig. 10. Visualization of neural activations (before the mask layer) of interpretable filters learned with and without the filter loss. We visualized neural activations of the first interpretable conv-layer before the mask layer in the CNN. In comparison, visualization results in Fig. 5 correspond to feature maps after the mask layer. Filters trained with the filter loss tended to generate more concentrated neural activations and have higher semantic purity than filters learned without the filter loss.

	bird	cat	cow	dog	horse	sheep	Avg.
VGG-16	0.144	0.134	0.146	0.127	0.141	0.142	0.139
VGG-16 + mask layer, w/o filter loss	0.141	<b>0.126</b>	0.139	0.124	0.138	0.136	0.134
VGG-16, interpretable	<b>0.130</b>	0.127	<b>0.134</b>	<b>0.109</b>	<b>0.131</b>	<b>0.126</b>	<b>0.126</b>

TABLE 12

Location instability of filters ( $\mathbb{E}_{f,k}[D_{f,k}]$ ) in the first conv-layer. The CNN was trained for the binary classification of a single category using the VOC Part dataset [4]. For the baseline, we added a conv-layer to the ordinary VGG-16 network, and selected the corresponding conv-layer in the network to enable fair comparisons. Interpretable filters learned with both the filter loss and the mask layer exhibited much lower localization instability than those learned with the mask layer but without the filter loss.

	bird	cow	cat	dog	horse	sheep
head	19.2	—	—	—	15.4	—
neck	21.2	—	—	5.8	—	—
torso	36.5	32.7	3.8	38.5	75.0	52.0
hip & tail	5.8	—	—	—	—	—
foot	5.8	—	—	—	9.6	3.8
wing	11.5	—	—	—	—	—
eye	—	30.8	55.8	11.5	7.7	—
nose & mouth	—	15.4	32.7	3.8	19.2	—
side face	—	9.6	—	—	—	—
leg	—	11.5	5.8	23.1	—	—
ear & horn	—	—	1.9	17.3	17.3	—

TABLE 9

Statistics of semantic meanings of interpretable filters. “—” indicates that the part is not selected as a label to describe the filter in a CNN. Except for CNNs for the bird and the horse, CNNs for other animals paid attention to detailed structures of the head. Thus, we annotated fine-grained parts inside the head for these CNNs.

#### 4.4.2 Visualization of filters

Besides the quantitative analysis of neural activation purity, we visualized filter activations to compare the interpretable CNN and the CNN trained without

Model	Original CNN	Interpretable CNN
VGG-M	0.00302±0.00123	0.00243±0.00120
VGG-S	0.00305±0.00120	0.00266±0.00133
VGG-16	0.00293±0.00128	0.00280±0.00152

TABLE 15

Average adversarial distortion of the original CNN and the interpretable CNN.

the filter loss. Fig. 10 visualized neural activations of the filter in the first interpretable conv-layer of VGG-16 before the mask layer. Visualization results demonstrated that filters trained with the filter loss could generate more concentrated neural activations.

#### 4.4.3 Location instability

We compared the location instability among interpretable filters learned with the filter loss, those learned without the filter loss, and ordinary filters. Here, we used filters in the first interpretable conv-layer of the interpretable CNN and filters in the corresponding conv-layer of the traditional CNNs for comparison.

We constructed the competing CNNs based on the

	gold.	bird	frog	turt.	liza.	koala	lobs.	dog	fox	cat	lion
VGG-16	0.152	0.152	0.142	0.150	0.167	0.125	0.126	0.139	0.137	0.149	0.134
VGG-16 + mask layer, w/o filter loss	0.143	0.148	0.141	0.146	0.162	0.121	0.120	0.137	0.130	0.142	0.130
VGG-16, interpretable	<b>0.105</b>	<b>0.127</b>	<b>0.131</b>	<b>0.139</b>	<b>0.157</b>	<b>0.102</b>	<b>0.118</b>	<b>0.123</b>	<b>0.105</b>	<b>0.129</b>	<b>0.106</b>
	tiger	bear	rabb.	hams.	squi.	horse	zebra	swine	hippo.	catt.	sheep
VGG-16	0.143	0.139	0.143	0.124	0.145	0.147	0.155	0.138	0.139	0.142	0.148
VGG-16 + mask layer, w/o filter loss	0.137	0.131	0.140	0.119	0.134	0.143	0.142	0.133	0.131	0.134	0.145
VGG-16, interpretable	<b>0.112</b>	<b>0.118</b>	<b>0.122</b>	<b>0.098</b>	<b>0.106</b>	<b>0.136</b>	<b>0.095</b>	<b>0.125</b>	<b>0.120</b>	<b>0.126</b>	<b>0.128</b>
	ante.	camel	otter	arma.	monk.	elep.	red pa.	gia.pa.			<b>Avg.</b>
VGG-16	0.141	0.144	0.155	0.156	0.134	0.127	0.139	0.126			0.142
VGG-16 + mask layer, w/o filter loss	0.134	0.142	0.147	0.153	0.129	0.124	0.125	0.116			0.136
VGG-16, interpretable	<b>0.119</b>	<b>0.131</b>	<b>0.134</b>	<b>0.118</b>	<b>0.112</b>	<b>0.109</b>	<b>0.098</b>	<b>0.088</b>			<b>0.118</b>

TABLE 13

Location instability of filters ( $\mathbb{E}_{f,k}[D_{f,k}]$ ) in the first interpretable conv-layer. The CNN was trained for the binary classification of a single category using the ILSVRC 2013 DET Animal-Part dataset [42]. For the baseline, we added a conv-layer to the ordinary VGG-16 network, and selected the corresponding conv-layer in the network to enable fair comparisons. Interpretable filters learned with both the filter loss and the mask layer exhibited much lower localization instability than those learned with the mask layer but without the filter loss.

Model	Logistic log loss		Softmax log loss	
	on target categories	on other categories	on target categories	on other categories
VGG-M, interpretable	107.9	7.5	6.2	1.1
w/o filter loss	18.6	6.3	4.7	1.1
VGG-S, interpretable	32.1	10.5	18.8	3.1
w/o filter loss	18.9	8.2	11.5	2.3
VGG-16, interpretable	948.9	81.5	106.7	5.1
w/o filter loss	40.5	12.9	67.0	6.0

TABLE 14

The mean value of neural activations on the target categories and those on other categories. Filters learned with the filter loss exhibited were usually more discriminative than those learned without the filter loss.

VGG-16 architecture, and these CNNs were trained for single-category classification based on the VOC Part dataset and the ILSVRC 2013 DET Animal-Part dataset. As shown in Table 12 and Table 13, the filter loss forced each filter to focus on a specific object part and reduced the location instability.

#### 4.4.4 Activation magnitudes

We further tested effects of the interpretable loss on neural activations among different categories. We used the VGG-M, VGG-S, and VGG-16 networks with either the logistic log loss or the softmax loss, which was trained to classify animal categories in the VOC Part dataset [4]. For each interpretable filter, given images of its target category, we recorded their neural activations (*i.e.* recording the maximal activation value in each of their feature maps). At the same time, we also recorded neural activations on other categories of the filter. The interpretable filter was supposed to activate much more strongly on its target category than on other (unrelated) categories. We collected all activation records on corresponding categories of all filters, and their mean value is reported in Table 14. In comparison, we also computed the mean value of

all activations on unrelated categories of all filters in Table 14. This table shows that the interpretable filter was usually activated more strongly on the target category than on other categories. Furthermore, we also compared the proposed interpretable CNN with the ablation baseline *w/o filter loss*, in which the CNN was learned without the filter loss. Table 14 shows that the filter loss made each filter more prone to being triggered by a single category, *i.e.* boosting the feature interpretability.

## 5 CONCLUSION AND DISCUSSIONS

In this paper, we have proposed a general method to enhance feature interpretability of CNNs. We design a loss to push a filter in high conv-layers towards the representation of an object part during the learning process without any part annotations. Experiments have shown that each interpretable filter consistently represents a certain object part of a category through different input images. In comparison, each filter in the traditional CNN usually represents a mixture of parts and textures.

Meanwhile, the interpretable CNN still has some drawbacks. First, in the scenario of multi-category

classification, filters in a conv-layer are assigned with different categories. In this way, when we need to classify a large number of categories, theoretically, each category can only obtain a few filters, which will decrease a bit the classification performance. Otherwise, the interpretable conv-layer must contain lots of filters to enable the classification of a large number of categories. Second, the learning of the interpretable CNN has a strong assumption, *i.e.* each input image must contain a single object, which limits the applicability of the interpretable CNN. Third, the filter loss is only suitable to learn high conv-layers, because low conv-layers usually represent textures, instead of object parts. Finally, the interpretable CNN is not suitable to encode textural patterns.

## ACKNOWLEDGMENTS

This work is partially supported by National Natural Science Foundation of China (U19B2043 and 61906120), DARPA XAI Award N66001-17-2-4029, NSF IIS 1423305, and ARO project W911NF1810296.

## APPENDIX

### PROOF OF EQUATION (4)

$$\begin{aligned}
 \frac{\partial \text{Loss}}{\partial x_{ij}} &= - \sum_{\mu \in \Omega} p(\mu) \left\{ \frac{\partial p(x|\mu)}{\partial x_{ij}} [\log p(x|\mu) - \log p(x) + 1] \right. \\
 &\quad \left. - p(x|\mu) \frac{\partial \log p(x)}{\partial x_{ij}} \right\} \\
 &= - \sum_{\mu \in \Omega} p(\mu) \left\{ \frac{\partial p(x|\mu)}{\partial x_{ij}} [\log p(x|\mu) - \log p(x) + 1] \right. \\
 &\quad \left. - p(x|\mu) \frac{1}{p(x)} \frac{\partial p(x)}{\partial x_{ij}} \right\} \\
 &= - \sum_{\mu \in \Omega} p(\mu) \left\{ \frac{\partial p(x|\mu)}{\partial x_{ij}} [\log p(x|\mu) - \log p(x) + 1] \right. \\
 &\quad \left. - p(x|\mu) \frac{1}{p(x)} \sum_{\mu'} \left[ p(\mu') \frac{\partial p(x|\mu')}{\partial x_{ij}} \right] \right\} \\
 &= - \sum_{\mu \in \Omega} p(\mu) \left\{ \frac{\partial p(x|\mu)}{\partial x_{ij}} [\log p(x|\mu) - \log p(x) + 1] \right\} \\
 &\quad + \sum_{\mu \in \Omega} p(\mu) \frac{\partial p(x|\mu)}{\partial x_{ij}} \frac{\sum_{\mu'} p(\mu') p(x|\mu')}{p(x)} \\
 &= - \sum_{\mu \in \Omega} p(\mu) \left\{ \frac{\partial p(x|\mu)}{\partial x_{ij}} [\log p(x|\mu) - \log p(x) + 1] \right\} \\
 &\quad + \sum_{\mu \in \Omega} p(\mu) \frac{\partial p(x|\mu)}{\partial x_{ij}} \\
 &= - \sum_{\mu \in \Omega} \frac{\partial p(x|\mu)}{\partial x_{ij}} p(\mu) [\log p(x|\mu) - \log p(x)] \\
 &= - \sum_{\mu \in \Omega} \frac{t_{ij} p(\mu) e^{tr(x \cdot T)}}{Z_{\mu}} \left\{ tr(x \cdot T) - \log [Z_{\mu} p(x)] \right\}
 \end{aligned} \tag{8}$$

### PROOF OF EQUATION (5)

$$\begin{aligned}
 \text{Loss} &= -MI(\mathbf{X}; \Omega) \quad // \quad \Omega = \{\mu^-, \mu_1, \mu_2, \dots, \mu_{n^2}\} \\
 &= -H(\Omega) + H(\Omega|\mathbf{X}) \\
 &= -H(\Omega) - \sum_x p(x) \sum_{\mu \in \Omega} p(\mu|x) \log p(\mu|x) \\
 &= -H(\Omega) - \sum_x p(x) \left\{ p(\mu^-|x) \log p(\mu^-|x) \right. \\
 &\quad \left. + \sum_{\mu \in \Omega^+} p(\mu|x) \log p(\mu|x) \right\} \\
 &= -H(\Omega) - \sum_x p(x) \left\{ p(\mu^-|x) \log p(\mu^-|x) \right. \\
 &\quad \left. + \sum_{\mu \in \Omega^+} p(\mu|x) \log \left[ \frac{p(\mu|x)}{p(\Omega^+|x)} p(\Omega^+|x) \right] \right\} \\
 &= -H(\Omega) - \sum_x p(x) \left\{ p(\mu^-|x) \log p(\mu^-|x) \right. \\
 &\quad \left. + p(\Omega^+|x) \log p(\Omega^+|x) + \sum_{\mu \in \Omega^+} p(\mu|x) \log \frac{p(\mu|x)}{p(\Omega^+|x)} \right\} \\
 &= -H(\Omega) + H(\Omega' = \{\mu^-, \Omega^+\}|\mathbf{X}) \\
 &\quad + \sum_x p(\Omega^+, x) H(\Omega^+|X = x)
 \end{aligned}$$

where  $p(\Omega^+|x) = \sum_{\mu \in \Omega^+} p(\mu|x)$ ,  $H(\Omega^+|X = x) = \sum_{\mu \in \Omega^+} \tilde{p}(\mu|X = x) \log \tilde{p}(\mu|X = x)$ ,  $\tilde{p}(\mu|X = x) = \frac{p(\mu|x)}{p(\Omega^+|x)}$ .

## REFERENCES

- [1] M. Aubry and B. C. Russell. Understanding deep features with computer-generated imagery. *In ICCV*, 2015.
- [2] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba. Network dissection: Quantifying interpretability of deep visual representations. *In CVPR*, 2017.
- [3] R. Chen, H. Chen, G. Huang, J. Ren, and Q. Zhang. Explaining neural networks semantically and quantitatively. *In arXiv:1812.07169*, 2018.
- [4] X. Chen, R. Mottaghi, X. Liu, S. Fidler, R. Urtasun, and A. Yuille. Detect what you can: Detecting and representing objects using holistic models and body parts. *In CVPR*, 2014.
- [5] A. Dosovitskiy and T. Brox. Inverting visual representations with convolutional networks. *In CVPR*, 2016.
- [6] R. C. Fong and A. Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. *In arXiv:1704.03296v1*, 2017.
- [7] Y. Goyal, A. Mohapatra, D. Parikh, and D. Batra. Towards transparent ai systems: Interpreting visual question answering models. *In arXiv:1608.08974v2*, 2016.
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *In CVPR*, 2016.
- [9] Z. Hu, X. Ma, Z. Liu, E. Hovy, and E. P. Xing. Harnessing deep neural networks with logic rules. *In arXiv:1603.06318v2*, 2016.
- [10] V. K. Ithapu. Decoding the deep: Exploring class hierarchies of deep representations using multiresolution matrix factorization. *In CVPR Workshop on Explainable Computer Vision and Job Candidate Screening Competition*, 2017.
- [11] P. Koh and P. Liang. Understanding black-box predictions via influence functions. *In ICML*, 2017.
- [12] S. Kolouri, C. E. Martin, and H. Hoffmann. Explaining distributed neural activations via unsupervised learning. *In CVPR Workshop on Explainable Computer Vision and Job Candidate Screening Competition*, 2017.
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *In NIPS*, 2012.



- [14] D. Kumar, A. Wong, and G. W. Taylor. Explaining the unexplained: A class-enhanced attentive response (clear) approach to understanding deep neural networks. In *CVPR Workshop on Explainable Computer Vision and Job Candidate Screening Competition*, 2017.
- [15] H. Lakkaraju, E. Kamar, R. Caruana, and E. Horvitz. Identifying unknown unknowns in the open world: Representations and policies for guided exploration. In *AAAI*, 2017.
- [16] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, 1998.
- [17] B. J. Lengerich, S. Konam, E. P. Xing, S. Rosenthal, and M. Veloso. Visual explanations for convolutional neural networks via input resampling. In *ICML Workshop on Visualization for Deep Learning*, 2017.
- [18] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *NIPS*, 2017.
- [19] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. In *CVPR*, 2015.
- [20] C. Olah, A. Mordvintsev, and L. Schubert. Feature visualization. *Distill*, 2017. <https://distill.pub/2017/feature-visualization>.
- [21] M. T. Ribeiro, S. Singh, and C. Guestrin. “why should i trust you?” explaining the predictions of any classifier. In *KDD*, 2016.
- [22] A. S. Ross, M. C. Hughes, and F. Doshi-Velez. Right for the right reasons: Training differentiable models by constraining their explanations. In *arXiv:1703.03717v1*, 2017.
- [23] S. Sabour, N. Frosst, and G. E. Hinton. Dynamic routing between capsules. In *NIPS*, 2017.
- [24] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, 2017.
- [25] Z. Si and S.-C. Zhu. Learning and-or templates for object recognition and detection. In *PAMI*, 2013.
- [26] M. Simon and E. Rodner. Neural activation constellations: Unsupervised part model discovery with convolutional networks. In *ICCV*, 2015.
- [27] M. Simon, E. Rodner, and J. Denzler. Part detector discovery in deep convolutional neural networks. In *ACCV*, 2014.
- [28] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: visualising image classification models and saliency maps. In *arXiv:1312.6034*, 2013.
- [29] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [30] S. Singh, A. Gupta, and A. A. Efros. Unsupervised discovery of mid-level discriminative patches. In *ECCV*, 2012.
- [31] J. Su, D. V. Vargas, and S. Kouichi. One pixel attack for fooling deep neural networks. In *arXiv:1710.08864*, 2017.
- [32] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. In *ICLR*, 2014.
- [33] J. Vaughan, A. Sudjianto, E. Brahimi, J. Chen, and V. N. Nair. Explainable neural networks based on additive index models. In *arXiv:1806.01933*, 2018.
- [34] C. Ventura, D. Masip, and A. Lapedriza. Interpreting cnn models for apparent personality trait regression. In *CVPR Workshop on Explainable Computer Vision and Job Candidate Screening Competition*, 2017.
- [35] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. Technical report, In California Institute of Technology, 2011.
- [36] A. S. Wicaksana and C. C. S. Liem. Human-explainable features for job candidate screening prediction. In *CVPR Workshop on Explainable Computer Vision and Job Candidate Screening Competition*, 2017.
- [37] N. Wolchover. New theory cracks open the black box of deep learning. In *Quanta Magazine*, 2017.
- [38] J. Xie, W. Hu, S.-C. Zhu, and Y. N. Wu. Learning inhomogeneous frame models for object patterns. In *CVPR*, 2014.
- [39] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *NIPS*, 2014.
- [40] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *ECCV*, 2014.
- [41] Q. Zhang, R. Cao, F. Shi, Y. Wu, and S.-C. Zhu. Interpreting cnn knowledge via an explanatory graph. In *AAAI*, 2018.
- [42] Q. Zhang, R. Cao, Y. N. Wu, and S.-C. Zhu. Growing interpretable part graphs on convnets via multi-shot learning. In *AAAI*, 2016.
- [43] Q. Zhang, R. Cao, Y. N. Wu, and S.-C. Zhu. Mining object parts from cnns via active question-answering. In *CVPR*, 2017.
- [44] Q. Zhang, R. Cao, S. Zhang, M. Edmonds, Y. N. Wu, and S.-C. Zhu. Interactively transferring cnn patterns for part localization. In *arXiv:1708.01783*, 2017.
- [45] Q. Zhang, W. Wang, and S.-C. Zhu. Examining cnn representations with respect to dataset bias. In *AAAI*, 2018.
- [46] Q. Zhang, Y. N. Wu, and S.-C. Zhu. Interpretable convolutional neural networks. In *CVPR*, 2018.
- [47] Q. Zhang, Y. Yang, H. Ma, and Y. N. Wu. Interpreting cnns via decision trees. In *CVPR*, 2019.
- [48] Q. Zhang and S.-C. Zhu. Visual interpretability for deep learning: a survey. in *Frontiers of Information Technology & Electronic Engineering*, 19(1):27–39, 2018.
- [49] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Object detectors emerge in deep scene cnns. In *ICRL*, 2015.
- [50] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016.



**Quanshi Zhang** received the B.S. degree in machine intelligence from the Peking University, China, in 2009 and M.S. and Ph.D. degrees in center for spatial information science from the University of Tokyo, Japan, in 2011 and 2014, respectively. In 2014, he went to the University of California, Los Angeles, as a post-doctoral associate. Now, he is an associate professor at the Shanghai Jiao Tong University. His research interests include computer vision, machine learning,

and robotics.



**Xin Wang** is a Ph.D. student an internship student at the Shanghai Jiao Tong University. His research mainly focuses on machine learning and computer vision.



**Ying Nian Wu** received a Ph.D. degree from the Harvard University in 1996. He was an Assistant Professor at the University of Michigan between 1997 and 1999 and an Assistant Professor at the University of California, Los Angeles between 1999 and 2001. He became an Associate Professor at the University of California, Los Angeles in 2001. From 2006 to now, he is a professor at the University of California, Los Angeles. His research interests include statistics, machine learning, and computer vision.



**Huilin Zhou** received a B.S. degree in mathematics from the University of Electronic Science and Technology of China in 2019. Now, she is a Ph.D. candidate at the Shanghai Jiao Tong University. Her research interests include computer vision and machine learning.



**Song-Chun Zhu** received a Ph.D. degree from Harvard University, and is a professor with the Department of Statistics and the Department of Computer Science at UCLA. His research interests include computer vision, statistical modeling and learning, cognition and AI, and visual arts. He received a number of honors, including the Marr Prize in 2003 with Z. Tu et. al. on image parsing, the Aggarwal prize from the Int'l Association of Pattern Recognition in 2008, twice Marr Prize honorary nominations in 1999 for texture modeling and 2007 for object modeling with Y.N. Wu et al., a Sloan Fellowship in 2001, the US NSF Career Award in 2001, and the US ONR Young Investigator Award in 2001. He is a Fellow of IEEE.