



# Choose Your Neuron: Incorporating Domain Knowledge Through Neuron-Importance

Ramprasaath R. Selvaraju<sup>1</sup>(✉), Prithvijit Chattopadhyay<sup>1</sup>,  
Mohamed Elhoseiny<sup>2</sup>, Tilak Sharma<sup>2</sup>, Dhruv Batra<sup>1,2</sup>, Devi Parikh<sup>1,2</sup>,  
and Stefan Lee<sup>1</sup>

<sup>1</sup> Georgia Institute of Technology, Atlanta, USA

{ramprs,prithvijit3,dbatra,parikh,steflee}@gatech.edu

<sup>2</sup> Facebook, Menlo Park, USA

{elhoseiny,tilaksharma,dbatra,dparikh}@fb.com

**Abstract.** Individual neurons in convolutional neural networks supervised for image-level classification tasks have been shown to implicitly learn semantically meaningful concepts ranging from simple textures and shapes to whole or partial objects – forming a “dictionary” of concepts acquired through the learning process. In this work we introduce a simple, efficient zero-shot learning approach based on this observation. Our approach, which we call Neuron Importance-Aware Weight Transfer (NIWT), learns to map domain knowledge about novel “*unseen*” classes onto this dictionary of learned concepts and then optimizes for network parameters that can effectively combine these concepts – essentially learning classifiers by discovering and composing learned semantic concepts in deep networks. Our approach shows improvements over previous approaches on the CUBirds and AWA2 generalized zero-shot learning benchmarks. We demonstrate our approach on a diverse set of semantic inputs as external domain knowledge including attributes and natural language captions. Moreover by learning inverse mappings, NIWT can provide visual and textual explanations for the predictions made by the newly learned classifiers and provide neuron names. Our code is available at <https://github.com/ramprs/neuron-importance-zsl>.

**Keywords:** Zero Shot Learning · Interpretability · Grad-CAM

## 1 Introduction

Deep neural networks have pushed the boundaries of standard classification tasks in the past few years, with performance on many challenging benchmarks reaching near human-level accuracies. One caveat however is that these deep models require massive labeled datasets – failing to generalize from few examples

---

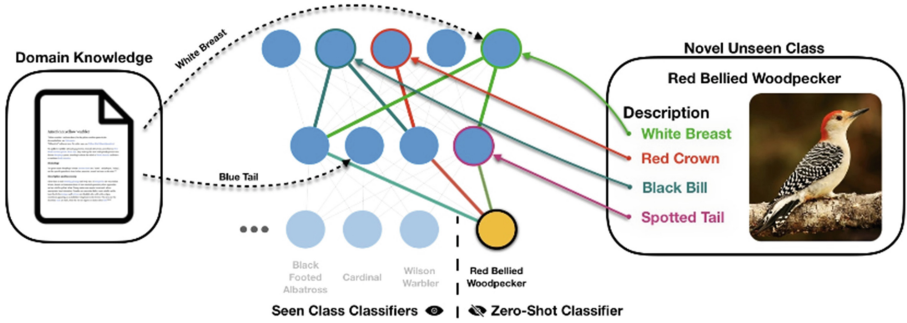
R. R. Selvaraju and P. Chattopadhyay—Equal Contribution.

R. R. Selvaraju—Work done partly at Facebook.

© Springer Nature Switzerland AG 2018

V. Ferrari et al. (Eds.): ECCV 2018, LNCS 11217, pp. 540–556, 2018.

[https://doi.org/10.1007/978-3-030-01261-8\\_32](https://doi.org/10.1007/978-3-030-01261-8_32)



**Fig. 1.** We present our Neuron Importance-Aware Weight Transfer (NIWT) approach which maps free-form domain knowledge about unseen classes to relevant concept-sensitive neurons within a pretrained deep network. We then optimize the weights of a novel classifier such that the activation of this set of neurons results in high output scores for the unseen classes in the generalized zero-shot learning setting.

or descriptions of unseen classes like humans can. To close this gap, the task of learning deep classifiers for unseen classes from external domain knowledge alone – termed zero-shot learning (ZSL) – has been the topic of increased interest within the community [2, 3, 5, 10, 11, 14, 16, 17, 21, 25, 28, 30, 35].

As humans, much of the way we acquire and transfer knowledge about novel concepts is in reference to or via composition of concepts which are already known. For instance, upon hearing that “A Red Bellied Woodpecker is a small, round bird with a white breast, red crown, and spotted wings.”, we can compose our understanding of colors and birds to imagine how we might distinguish such an animal from other birds. However, applying a similar compositional learning strategy for deep neural networks has proven challenging.

While individual neurons in deep networks have been shown to learn localized, semantic concepts, these units lack referable groundings – *e.g.* even if a network contains units sensitive to “white breast” and “red crown”, there is no explicit mapping of these neurons to the relevant language name or description. This observation encouraged prior work in interpretability to crowd-source “neuron names” to discover these groundings [4]. However, this annotation process is model dependent and needs to be re-executed for each model trained, which makes it expensive and impractical. Moreover, even if given perfect “neuron names”, it is an open question how to leverage this neuron-level descriptive supervision to train novel classifiers. This question is at the heart of our approach.

Many existing zero-shot learning approaches make use of deep features (*e.g.* vectors of activations from some late layer in a network pretrained on some large-scale task) to learn joint embeddings with class descriptions [1, 3, 5, 7–9, 23, 31]. These higher-level features collapse many underlying concepts in the pursuit of class discrimination; consequentially, accessing lower-level concepts and recombining them in new ways to represent novel classes is difficult with

these features. Mapping class descriptions to lower-level activations directly on the other hand is complicated by the high intra-class variance of activations due to both spatial and visual differences within instances of a class. Our goal is to address these challenges by grounding class descriptions (including attributes and free-form text) to the *importance* of lower-layer neurons to final network decisions [26].

In our approach, which we call Neuron Importance-based Weight Transfer (NIWT), we learn a mapping between class-specific domain knowledge and the importances of individual neurons within a deep network. This mapping is learnt using images (to compute neuron-importance) and corresponding domain knowledge representation(s) of training classes. We then use this learned mapping to predict neuron importances from knowledge about unseen classes and optimize classification weights such that the resulting network aligns with the predicted importances. In other words, based on domain-knowledge of the unseen categories, we can predict which low-level neurons should matter in the final classification decision. We can then learn network weights such that the neurons predicted to matter actually do contribute to the final decision. In this way, we connect the description of a previous unseen category to weights of a classifier that can predict this category at test time – all without having seen a single image from this category. To the best of our knowledge, this is the first zero-shot learning approach to align domain knowledge to intermediate neurons within a deep network. As an additional benefit, the learned mapping from domain knowledge to neuron importances grounds the neurons in interpretable semantics; automatically performing neuron naming.

We focus on the challenging generalized zero-shot (GZSL) learning setting. Unlike standard ZSL settings which evaluate performance only on unseen classes, GZSL considers both unseen and seen classes to measure the performance. In effect, GZSL is made more challenging by dropping the unrealistic assumption that test instances are known a priori to be from unseen classes in standard ZSL. We validate our approach across two standard datasets - Caltech-UCSD Birds (CUB) [29] and Animals with Attributes 2 (AWA2) [31] - showing improved performance over existing methods. Moreover, we examine the quality of our grounded explanations for classifier decisions through textual and visual examples.

**Contributions.** Concretely, we make the following contributions in this work:

- We introduce a zero-shot learning approach based on mapping unseen class descriptions to neuron importance within a deep network and then optimizing unseen classifier weights to effectively combine these concepts. We demonstrate the effectiveness of our approach by reporting improvements on the generalized zero-shot benchmark on CUB and AWA2. We also show our approach can handle arbitrary forms of domain knowledge including attributes and captions.
- In contrast to existing approaches, our method is capable of explaining its zero-shot predictions with human-interpretable semantics from attributes. We show how inverse mappings from neuron importance to domain knowledge can

also be learned to provide interpretable visual and textual explanations for the decisions made by newly learned classifiers for seen and unseen classes.

## 2 Related Work

**Model Interpretability.** Our method aligns human interpretable domain knowledge to neurons within deep neural networks, instilling these neurons with understandable semantic meanings. There has been significant recent interest in building machine learning models that are transparent and interpretable in their decision making process. For deep networks, several works propose explanations based on internal states or structures of the network [12, 26, 33, 36]. Most related to our work is the approach of Selvaraju *et al.* [26] which computes neuron importance as part of a visual explanation pipeline. In this work, we leverage these importance scores to embed free-form domain knowledge to individual neurons in a deep network and train new classifiers based on this information. In contrast, Grad-CAM [26] simply visualizes the importance of input regions.

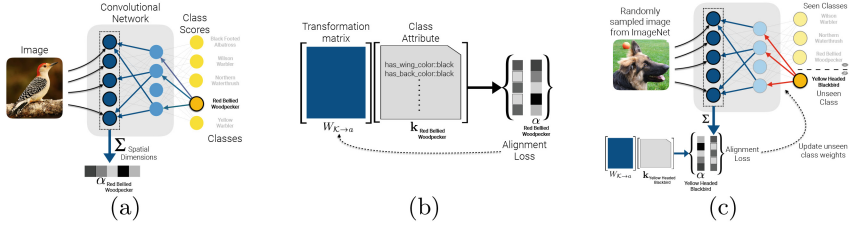
**Attribute-based Zero-Shot Learning.** One long-pursued approach for zero-shot learning is to leverage knowledge about common attributes and shared parts (e.g., furry, in addition to being simpler and more efficient [2, 3, 25, 31]).

**Text-Based Zero-Shot Learning (ZSL).** In parallel research, pure text articles extracted from the web have been leveraged instead of attributes to design zero-shot visual classifiers [8]. The description of a new category is purely textual (avoiding the use of attributes) and could be extracted easily by just mining article(s) about the class of interest from the web (e.g., Wikipedia). Recent approaches have adopted deep neural network based classifiers, leading to a noticeable improvement on zero-shot accuracy (Bo *et al.* [18]). The proposed approaches mainly rely on learning a similarity function between text descriptions and images (either linearly [8, 25] or non-linearly via deep neural networks [18] or kernels [7]). At test-time, classification is performed by associating the image to the class with the highest similarity to the corresponding class-level text. Recently, Reed *et al.* [24] showed that by collecting 10 sentences per-image, their sentence-based approach can outperform attribute-based alternatives on CUB.

In contrast to these approaches, we directly map external domain knowledge (text-based or otherwise) to internal components (neurons) of deep neural networks rather than learning associative mappings between images and text – providing interpretability for our novel classifiers.

## 3 Neuron Importance-Aware Weight Transfer (NIWT)

In this section, we describe our Neuron Importance-Aware Weight Transfer (NIWT) approach to zero-shot learning. At a high level, NIWT maps free-form domain knowledge to neurons within a deep network and then learns classifiers based on novel class descriptions which respect these groundings. Concretely,



**Fig. 2.** Our Neuron Importance-Aware Weight Transfer (NIWT) approach can be broken down into three stages. (a) class-specific neuron importances are extracted for seen classes at a fixed layer, (b) a linear transform is learned to project free-form domain knowledge to these extracted importances, and (c) weights for new classifiers are optimized such that neuron importances match those predicted by this mapping for unseen classes.

NIWT consists of three steps: (1) estimating the importance of individual neuron(s) at a fixed layer w.r.t. the decisions made by the network for the seen classes (see Fig. 2a), (2) learning a mapping between domain knowledge and these neuron-importances (see Fig. 2b), and (3) optimizing classifier weights with respect to predicted neuron-importances for unseen classes (see Fig. 3c). We discuss each stage in the following sections.

### 3.1 Preliminaries: Generalized Zero-Shot Learning (GZSL)

Consider a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  comprised of example input-output pairs from a set of *seen classes*  $\mathcal{S} = \{1, \dots, s\}$  and *unseen classes*  $\mathcal{U} = \{s+1, \dots, s+u\}$ . For convenience, we use the subscripts  $\mathcal{S}$  and  $\mathcal{U}$  to indicate subsets corresponding to seen and unseen classes respectively, e.g.  $\mathcal{D}_{\mathcal{S}} = \{(x_i, y_i) \mid y_i \in \mathcal{S}\}$ . Further, assume there exists domain knowledge  $\mathcal{K} = \{k_1, \dots, k_{s+u}\}$  corresponding to each class (e.g. class level attributes or natural language descriptions). Concisely, the goal of generalized zero-shot learning is then to learn a mapping  $f : \mathcal{X} \rightarrow \mathcal{S} \cup \mathcal{U}$  from the input space  $\mathcal{X}$  to the combined set of seen and unseen class labels using only the domain knowledge  $\mathcal{K}$  and instances  $\mathcal{D}_{\mathcal{S}}$  belonging to the seen classes.

### 3.2 Class-Dependent Neuron Importance

Class descriptions capture salient concepts about the content of corresponding images – for example, describing the coloration and shape of a bird’s head. Similarly, a classifier must also learn discriminative visual concepts in order to succeed; however, these concepts are not grounded in human interpretable language. In this stage, we identify neurons corresponding to these discriminative concepts before aligning them with domain knowledge in Sect. 3.3.

Consider a deep neural network  $\text{NET}_{\mathcal{S}}(\cdot)$  trained for classification which predicts scores  $\{o_c \mid c \in \mathcal{S}\}$  for seen classes  $\mathcal{S}$ . One intuitive measure of a neuron

$n$ 's importance to the final score  $o_c$  is simply the gradient of  $o_c$  with respect to the neuron's activation  $a^n$  (where  $n$  indexes the channel dimension). For networks containing convolutional units (which are replicated spatially), we follow [26] and simply compute importance as the mean gradient (along spatial dimensions), writing the neuron importance  $\alpha_c^n$  as

$$\alpha_c^n = \overbrace{\frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W \frac{\partial o_c}{\partial a_{ij}^n}}^{\text{global average pooling}} \underbrace{\hspace{10em}}_{\text{gradients via backprop}} \quad (1)$$

where  $a_{i,j}^n$  is the activation of neuron  $n$  at spatial position  $i, j$ . For a given input, the importance of every neuron in the network can be computed for a given class via a single backward pass followed by a global average pooling operation for convolutional units. In practice, we focus on  $\alpha$ 's from single layers in the network in our experiments. We note that other measures of neuron importance have been proposed [15,32] in various contexts; however, this simple gradient-based importance measure has some notable properties which we leverage.

Firstly, we find gradient-based importance scores to be quite consistent across images of the same class despite the visual variation between instances, and likewise to correlate poorly across classes. To assess this quantitatively, we computed  $\alpha$ 's for neurons in the final convolutional layer of a convolutional neural network trained on a fine-grained multi-class task (conv5-3 of VGG-16 [27] trained on AWA2 [31]) for 10,000 randomly selected images. We observed an average rank correlation of 0.817 for instances within the same class and 0.076 across pairs of classes. This relative invariance of  $\alpha$ 's to intra-class input variation may be due in part to the piece-wise linear decision boundaries in networks using ReLU [20] activations. As shown in [22], transitions between these linear regions are much less frequent between same-class inputs than across classes. Within the same linear region, activation gradients (and hence  $\alpha$ 's) are trivially identical.

Secondly, this measure is fully differentiable with respect to model parameters which we use to learn novel classifiers with gradient methods (see Sect. 3.4).

### 3.3 Mapping Domain Knowledge to Neurons

Without loss of generality, consider a single layer  $L$  within  $\text{NET}_{\mathcal{S}}(\cdot)$ . Given an instance  $(x_i, y_i) \in \mathcal{D}_{\mathcal{S}}$ , let  $\mathbf{a}_c = \{\alpha_c^n \mid n \in L\}$  be a vector of importances computed for neurons in  $L$  with respect to class  $c$  when  $x_i$  is passed through the network. In this section, we learn a simple linear mapping from domain knowledge to these importance vectors – aligning interpretable semantics with individual neurons.

We first compute the importance vector  $\mathbf{a}_{y_i}$  for each seen class instance  $(x_i, y_i)$  and match it with the domain knowledge representation  $k_{y_i}$  of the corresponding class. Given this dataset of  $(\mathbf{a}_{y_i}, k_{y_i})$  pairs, we learn a linear transform  $W_{\mathcal{K} \rightarrow a}$  to map domain knowledge to importances. As importances are gradient

based, we penalize errors in the predicted importances based on cosine distance – emphasizing alignment over magnitude. We minimize the cosine distance loss as

$$\mathcal{L}(\mathbf{a}_{y_i}, \mathbf{k}_{y_i}) = 1 - \frac{(W_{\mathcal{K} \rightarrow a} \cdot \mathbf{k}_{y_i}) \cdot \mathbf{a}_{y_i}}{\|W_{\mathcal{K} \rightarrow a} \cdot \mathbf{k}_{y_i}\| \|\mathbf{a}_{y_i}\|}, \quad (2)$$

via gradient descent to estimate  $W_{\mathcal{K} \rightarrow a}$ . We stop training when average rank-correlation of predicted and true importance vectors stabilizes for a set of held out validation classes from  $\mathcal{S}$ .

Notably, this is a many-to-one mapping with the domain knowledge of one class needing to predict many different importance vectors. Despite this, this mapping achieves average rank correlations of 0.2 to 0.5 for validation class instances. We explore the impact of error in importance vector prediction on weight optimization in Sect. 3.4. We also note that this simple linear mapping can also be learned in an inverse fashion, mapping neuron importances back to semantic concepts within the domain knowledge (which we explore in Sect. 6).

### 3.4 Neuron Importance to Classifier Weights

In this section, we use predicted importances to learn classifiers for the unseen classes. As these new classifiers will be built atop the trained seen-class network  $\text{NET}_{\mathcal{S}}$ , we modify  $\text{NET}_{\mathcal{S}}$  to extend the output space to include the unseen class – expanding the final fully-connected layer to include additional neurons with weight vectors  $\mathbf{w}^1, \dots, \mathbf{w}^u$  for the unseen classes such that the network now additionally outputs scores  $\{o_c \mid c \in \mathcal{U}\}$ . We refer to this expanded network as  $\text{NET}_{\mathcal{S} \cup \mathcal{U}}$ . At this stage, the weights for the unseen classes are sampled randomly from a multivariate normal distribution with parameters estimated from the seen class weights and as such the output scores are uncalibrated and uninformative.

Given the learned mapping  $W_{\mathcal{K} \rightarrow A}$  and unseen class domain knowledge  $\mathcal{K}_{\mathcal{U}}$ , we can predict unseen class importances  $A_{\mathcal{U}} = \{\mathbf{a}_1, \dots, \mathbf{a}_u\}$  with the importance vector for unseen class  $c$  predicted as  $\mathbf{a}_c = W_{\mathcal{K} \rightarrow a} \mathbf{k}_c$ . For a given input, we can compute importance vectors  $\hat{\mathbf{a}}_c$  for each unseen class  $c$ . As  $\hat{\mathbf{a}}^c$  is a function of the weight parameters  $\mathbf{w}_c$ , we can simply supervise  $\hat{\mathbf{a}}_c$  with the predicted importances  $\mathbf{a}_c$  and optimize  $w^c$  with gradient descent – minimizing the cosine distance loss between predicted and observed importance vectors. However, the cosine distance loss does not account for scale and without regularization the scale of weights (and as consequence the outputs) of seen and unseen classes might vary drastically, resulting in bias towards one set or the other.

To address this problem, we introduce a  $L_2$  regularization term which constrains the learned unseen weights to be a similar scale as the mean of seen weights  $\bar{\mathbf{w}}_{\mathcal{S}}$ . We write the final objective as

$$\mathcal{L}(\hat{\mathbf{a}}_c, \mathbf{a}_c) = 1 - \frac{\hat{\mathbf{a}}_c \cdot \mathbf{a}_c}{\|\hat{\mathbf{a}}_c\| \|\mathbf{a}_c\|} + \lambda \|\mathbf{w}_c - \bar{\mathbf{w}}_{\mathcal{S}}\|, \quad (3)$$

where  $\lambda$  is controls the strength of this regularization. We examine the effect of this trade-off in Sect. 5.1, finding training to be robust to a wide range of

$\lambda$  values. We note that as observed importances  $\mathbf{a}^c$  are themselves computed from network gradients, updating weights based on this loss requires computing a Hessian-vector product; however, this is relatively efficient as the number of weights for each unseen class is small and independent of those for other classes.

**Training Images.** Note that to perform the optimization described above, we need to pass images through the network to compute importance vectors. We observe importances to be only weakly correlated with image features and find they can be computed for any of the unseen classes irrespective of the input image class – as such, we find simply inputting images with natural statistics to be sufficient. Specifically, we pair random images from ImageNet [6] with random tuples  $(\hat{\mathbf{a}}_c, \mathbf{k}_c)$  to perform the importance to weight optimization.

## 4 Experiments

In this section, we evaluate our approach on generalized zero-shot learning (GZSL) (Sect. 4.1) and present analysis for each stage of NIWT (Sect. 5).

### 4.1 Experimental Setting

**Datasets and Metrics.** We conduct our GZSL experiments on the

- **Animals with Attributes 2 (AWA2)** [31] – The AWA2 dataset consists of 37,322 images of 50 animal species (on average 764 per class but with a wide range). Each class is labeled with 85 binary and continuous attributes.
- **Caltech-UCSD Birds 200 (CUB)** [29] – The CUB dataset consists of 11788 images corresponding to 200 species of birds. Each image and each species has been annotated with 312 binary and continuous attribute labels respectively. These attributes describe fine-grained physical bird features such as the color and shape of specific body parts. Additionally, each image is associated with 10 human captions [24].

For both datasets, we use the GZSL splits proposed in [31] which ensure that no unseen class occurs within the ImageNet [6] dataset which is commonly used for training classification networks for feature extraction. As in [30], we evaluate our approach using class-normalized accuracy computed over both seen and unseen classes (*e.g.* 200-way for CUB) – breaking the results down into unseen accuracy  $\text{Acc}_U$ , seen accuracy  $\text{Acc}_S$ , and the harmonic mean between them  $H$ .

**Models.** We experiment with ResNet101 [13] and VGG16 [27] models pretrained on ImageNet [6] and fine-tuned on the seen classes. For each, we train a version by finetuning all layers and another by updating only the final classification weights. Compared to ResNet, where we see sharp declines for fixed models (60.6% finetuned vs 28.26% fixed for CUB and 90.10% vs 70.7% for AWA2), VGG achieves similar accuracies for both finetuned and fixed settings (74.84% finetuned vs 66.8% fixed for CUB and 92.32% vs 91.44% for AWA2). We provide more training details in the Appendix.



**NIWT Settings.** To train the domain knowledge to importance mapping we hold out five seen classes and stop optimization when rank correlation between observed and predicted importances is highest. For attribute vectors, we use the class level attributes directly and for captions on CUB we use average word2vec embeddings [19] for each class. When optimizing for weights given importances, we stop when the loss fails to improve by 1% over 40 iterations. We choose values of  $\lambda$  (between  $1e^{-5}$  to  $1e^{-2}$ ), learning rate ( $1e^{-5}$  to  $1e^{-2}$ ) and the batch size ( $\{16, 32, 64\}$ ) by grid search on H for a disjoint set of validation classes sampled from the seen classes of the proposed splits [31] based (see Table 1).

**Baselines.** We compare NIWT with a number of well-performing zero-shot learning approaches based on learning joint embeddings of image features and class information. Methods like ALE [2] focus on learning compatibility functions for class labels and visual features using some form of ranking loss. In addition to comparing with ALE as reported in [31], we also compare with settings where the hyper-parameters have been directly tuned on the test-set.

We also compare against the recent Deep Embedding approach of [34] which also leverages deep networks, jointly aligning domain knowledge with deep features end-to-end. For both of the mentioned baselines, we utilize code provided by the authors and report results by directly tuning hyper-parameters on the test-set so as to convey an upper-bound of performance.

## 4.2 Results

We show results in Table 1 for AWA2 and CUB using all model settings. There are a number of interesting trends to observe:

1. **NIWT sets the state of the art in generalized zero-shot learning.** For both datasets, NIWT-Attributes based on VGG establishes a new state of the art for harmonic mean (48.1% for AWA2 and 37.0% for CUB). For AWA2, this corresponds to a  $\sim 10\%$  improvement over prior state-of-the-art which is based on deep feature embeddings. These results imply that mapping domain knowledge to internal neurons can lead to improved results.
2. **Seen-class finetuning yields improved harmonic mean H.** For CUB and AWA2, finetuning the VGG network on seen class images offers significant gains for NIWT (26.7% $\rightarrow$ 37.0% H and 36.1% $\rightarrow$ 48.1% H respectively); finetuning ResNet sees similar gains (17.3% $\rightarrow$ 27.7% H on CUB and 27.5% $\rightarrow$ 40.5% H on AWA2). Notably, these trends seem inconsistent for the compared methods.
3. **NIWT effectively grounds both attributes and free-form language.** We see strong performance both for attributes and captions across both networks (37.0% and 23.6% H for VGG and 27.7% and 23.8% H for ResNet). We note that we use relatively simple, class-averaged representations for captioning which may contribute to the lower absolute performance.

**Table 1.** Generalized Zero-Shot Learning performances on the proposed splits [31] for AWA2 and CUB. We report class-normalized accuracies on seen and unseen classes and harmonic mean. <sup>1</sup>reproduced from [31]. <sup>2</sup>based on code provided by the authors by tuning hyper-parameters on the test-set to convey an upper-bound of performance.

Method		AWA2 [32]			CUB [30]			
		Acc <sub>U</sub>	Acc <sub>S</sub>	H	Acc <sub>U</sub>	Acc <sub>S</sub>	H	
ResNet101 [13]	Fixed	ALE [2] <sup>1</sup>	14.0	81.8	23.9	23.7	<b>62.8</b>	<b>34.4</b>
		ALE [2] <sup>2</sup>	20.9	<b>88.8</b>	33.8	<b>24.7</b>	62.3	<b>34.4</b>
		Deep Embed. [35] <sup>2</sup>	<b>28.5</b>	82.3	<b>42.3</b>	22.3	45.1	29.9
		NIWT-Attributes	21.6	37.8	27.5	10.2	57.7	17.3
	FT	ALE [2] <sup>2</sup>	22.7	<b>75.1</b>	34.9	24.1	<b>60.8</b>	<b>34.5</b>
		Deep Embed. [35] <sup>2</sup>	21.5	59.6	31.6	<b>24.7</b>	57.4	<b>34.5</b>
		NIWT-Attributes	<b>42.3</b>	38.8	<b>40.5</b>	20.7	41.8	27.7
		NIWT-Caption		N/A		22.1	25.7	23.8
VGG16 [28]	Fixed	ALE [2] <sup>2</sup>	17.9	<b>84.3</b>	29.5	22.2	<b>54.8</b>	<b>31.6</b>
		Deep Embed. [35] <sup>2</sup>	<b>28.8</b>	81.7	<b>42.6</b>	<b>24.1</b>	45.2	31.5
		NIWT-Attributes	43.8	30.7	36.1	17.0	54.6	26.7
	FT	ALE [2] <sup>2</sup>	16.9	<b>91.5</b>	28.5	25.3	<b>62.6</b>	36.0
		Deep Embed. [35] <sup>2</sup>	26.6	83.3	38.2	27.0	49.7	35.0
		NIWT-Attributes	<b>35.3</b>	75.5	<b>48.1</b>	<b>31.5</b>	44.9	<b>37.0</b>
		NIWT-Caption		N/A		15.9	46.5	23.6

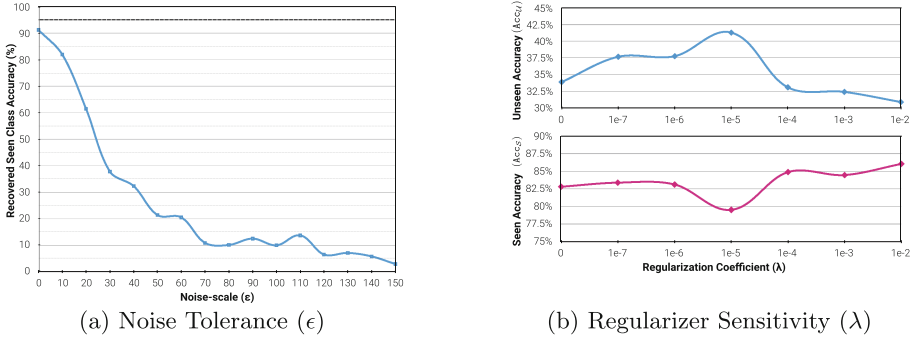
## 5 Analysis

To better understand the different stages of NIWT, we perform a series of experiments to analyze and isolate individual components in our approach.

### 5.1 Effect of Regularization Coefficient $\lambda$

One key component to our importance to weight optimization is the regularizer which enforces that learned unseen weights be close to the mean seen weight – avoiding arbitrary scaling of the learned weights and the bias this could introduce. To explore the effect of the regularizer, we vary the coefficient  $\lambda$  from 0 to  $1e^{-2}$ . Figure 3b shows the final seen and unseen class-normalized accuracy for the AWA2 dataset at convergence for different  $\lambda$ 's.

Without regularization ( $\lambda = 0$ ) the unseen weights tend to be a bit too small and achieve an unseen accuracy of only 33.9% on AWA2. As  $\lambda$  is increased the unseen accuracy grows until peaking at  $\lambda = 1e^{-5}$  with an unseen accuracy of 41.3% – an improvement of over 8% from the unregularized version! Of course, this improvement comes with a trade-off in seen accuracy of about 3% over the same interval. As  $\lambda$  grows larger  $>1e^{-4}$ , the regularization constraint becomes too strong and NIWT has trouble learning anything for the unseen classes.



**Fig. 3.** Analysis of the importance vector to weight optimization for VGG-16 trained on AWA2 (a). We find that ground-truth weights can be recovered for a pre-trained network even in the face of high magnitude noise. (b) We also show the importance of the regularization term to final model performance.

## 5.2 Noise Tolerance in Neuron Importance to Weight Optimization

One important component of NIWT is the ability to ground concepts learnt by a convolutional network in some referable domain. Due to the inherent noise involved in this mapping  $W_{\mathcal{K} \rightarrow A}$ , the classifiers obtained for unseen classes in the expanded network  $\text{NET}_{\text{SUL}}$  are not perfect. In order to judge the capacity of the optimization procedure, we experiment with a toy setting where we initialize an unseen classifier head with the same dimensionality as the seen classes and try to explicitly recover the seen class weights with supervision only from the *oracle*  $\mathbf{a}_c$  obtained from the seen classifier head for the seen classes. To simulate for the error involved in estimating  $\mathbf{a}_c$ , we add increasing levels of zero-centered gaussian noise and study recovery performance in terms of accuracy of the recovered classifier head on the seen-test split. That is, the supervision from importance vectors is constructed as follows:

$$\tilde{\mathbf{a}}_c = \mathbf{a}_c + \epsilon \overline{\|\mathbf{a}_c\|_1} \mathcal{N}(0, I) \quad (4)$$

We operate at different values of  $\epsilon$ , characterizing different levels of corruption of the supervision from  $\mathbf{a}_c$  and observe recovery performance in terms of accuracy of the recovered classifier head. Figure 3a shows the effect of noise on the ability to recover seen classifier weights (fc7) for a VGG-16 network trained on 40 seen classes of AWA2 dataset with the same objective as the one used for unseen classes.

In the absence of noise over  $\mathbf{a}_c$  supervision, we find that we are exactly able to recover the seen class weights and are able to preserve the pre-trained accuracy on seen classes. Even with a noise-level of  $\epsilon = 10$  (or adding noise with a magnitude 10x the average norm of  $\mathbf{a}_c$ ), we observe only minor reduction in the accuracy of the recovered seen class weights. As expected, this downward trend continues as we increase the noise-level until we reach almost chance-level performance on the recovered classifier head. This experiment shows that the importance vector to weights optimization is quite robust even to fairly extreme noise.

### 5.3 Network Depth of Importance Extraction

In this section, we explore the sensitivity of NIWT with respect to the layer from which we extract importance vectors in the convolutional network. As an experiment (in addition to Table 1) we evaluate NIWT on AWA2 with importance vectors extracted at different convolutional layers of VGG-16. We observe that out of those we experimented with `conv5_3` performs the best with  $H = 48.1$  followed by `conv4_3` ( $H = 39.3$ ), `conv3_3` ( $H = 35.5$ ), `conv2_2` ( $H = 23.8$ ) and `conv2_2` ( $H = 20.8$ ). We also experimented with the fully-connected layers `fc6` and `fc7` resulting in values of  $H$  being 40.2 and 1 respectively.

Note that performing NIWT on importance vectors extracted from the penultimate layer `fc7` is equivalent to learning the unseen head classifier weights directly from the domain space representation ( $\mathbf{k}_c$ ). Consistent with our hypothesis, this performs very poorly across all the metrics with almost no learning involved for the unseen classes at all. Though we note that this may be due to the restricted capacity of the linear transformation  $W_{\mathcal{K} \rightarrow A}$  involved in the process.

### 5.4 Importance to Weight Input Images

We evaluate performance with differing input images during weight optimization (random noise images, ImageNet images, and seen class images). We show results of each in Table 2. As expected, performance improves as input images more closely resemble the unseen classes; however, we note that learning occurs even with random noise images.

## 6 Explaining NIWT

In this section we demonstrate how we can use NIWT to provide visual and textual explanations for the decisions made by the newly learned classifiers on the unseen classes. In addition to the visual explanations provided by Grad-CAM [26], we utilize a mapping (similar to the one in Sect. 3.3) learned in the inverse direction –  $W_{a \rightarrow \mathcal{K}}$ , i.e., neuron-importance(s)  $\mathbf{a}_c$  to domain knowledge

$\mathcal{K}$  to ground the predictions made in the textual domain used as external knowledge. Since this mapping explicitly grounds the *important* neurons in the interpretable domain, we automatically obtain neuron names.

**Visual Explanations.** Since NIWT learns the classifier associated with the unseen classes as an extension to the existing deep network for the seen classes, it preserves the end-to-end differentiable pipeline for the novel classes as well. This allows us to directly use any of the existing deep network interpretability

**Table 2.** Results by sampling images from different sets for NIWT-Attributes on VGG-CUB.

Sampling mode	$\text{Acc}_U$	$\text{Acc}_S$	$H$
Random normal	23.9	41.0	30.2
ImageNet	31.5	44.9	37.0
Seen-classes	36.4	40.0	38.1

mechanisms to visually explain the decisions made at inference. We use Grad-CAM [26] on instances of unseen classes to visualize the support for decisions (see Fig. 4) made by the network with NIWT learnt classification weights.

**Evaluating Visual Explanations.** Quantitatively, we evaluate the generated maps for both seen and unseen classes by the mean fraction of the Grad-CAM activation present inside the bounding box annotations associated with the present objects. On seen classes, we found this number to be  $0.80 \pm 0.008$  versus  $0.79 \pm 0.005$  for the unseen classes on CUB – indicating that the unseen classifier learnt via NIWT is indeed capable of focusing on relevant regions in the input image while making a prediction.

**Textual Explanations.** In Sect. 3.3, we learned a mapping  $W_{\mathcal{K} \rightarrow a}$  to embed the external domain knowledge (attributes or captions) into the neurons of a specific layer of the network. Similarly, by learning an inverse mapping from the neuron importances to the attributes (or captions), we can ground the former associated with a prediction in a human-interpretable domain. We utilize such an inverse mapping to obtain scores in the attribute-space (given  $\mathbf{a}_c$ ) and retrieve the top-k attributes as explanations. A high scoring  $\mathbf{k}_c$  retrieved via  $W_{a \rightarrow \mathcal{K}}$  from a certain  $\mathbf{a}_c$  emphasizes the relevance of that attribute for the corresponding class  $c$ . This helps us ground the class-score decisions made by the learnt unseen classifier head in the attribute space, thus, providing an explanation for the decision.

**Evaluating Textual Explanations.** We evaluate the fidelity of such textual explanations by the percentage of associated ground truth attributes captured in the top-k generated explanations on a per instance level – 83.9% on CUB using a VGG-16 network. Qualitative results in Fig. 4 show visual and textual explanation(s) demonstrating the discriminative attributes learned by the model for predicting the given target category.

**Neuron Names and Focus.** Neuron names are referable groundings of concepts captured by a deep convolutional network. Unlike previous approaches, we obtain neuron names in an automatic fashion (without the use of any extra annotations) by feeding a one-hot encoded vector corresponding to a neuron being activated to  $W_{a \rightarrow \mathcal{K}}$  and performing a similar process of top-1 retrieval (as the textual explanations) to obtain the corresponding ‘*neuron name*’.

Figure 4 provides qualitative examples for named neurons and their activation maps. The green block shows instances where the unseen class images were correctly classified by  $\text{NET}_{\text{SUN}}$ . Conversely, those in red correspond to errors. The columns correspond to the class-labels, images, Grad-CAM visualizations for the class, textual explanations in the attribute space and top-3 neuron names responsible for the target class and their corresponding activation maps. For instance, notice that in the second row, for the image – correctly classified as a yellow-headed blackbird – the visualizations for the class focus specifically at the union of attributes that comprise the class. In addition, the textual explanations also filter out these attributes based on the neuron-importance scores - *has throat*

GT Class	Original Image	Visual Explanations	Text Explanations	Important neurons with corresponding activation maps		
Yellow-headed blackbird			<i>has_eye_color = black,</i> <i>has_underparts_color = white,</i> <i>has_belly_color = white,</i> <i>has_breast_color = white,</i> <i>has_breast_pattern = solid</i>	neuron_id = 145 <i>has_eye_color = black</i>	neuron_id = 299 <i>has_crown_color = yellow</i>	neuron_id = 20 <i>has_wing_color = black</i>
Yellow-headed blackbird			<i>has_eye_color = black,</i> <i>has_throat_color = yellow,</i> <i>has_wing_color = black,</i> <i>has_upperparts_color = black,</i> <i>has_bill_color = black</i>	neuron_id = 145 <i>has_eye_color = black</i>	neuron_id = 126 <i>has_throat_color = yellow</i>	neuron_id = 20 <i>has_wing_color = black</i>
Groove-billed Ani			<i>has_throat_color = black,</i> <i>has_primary_color = black,</i> <i>has_nape_color = black,</i> <i>has_forehead_color = black,</i> <i>has_crown_color = black</i>	neuron_id = 131 <i>has_throat_color = black</i>	neuron_id = 259 <i>has_throat_color = black</i>	neuron_id = 193 <i>has_nape_color = black</i>
Groove-billed Ani			<i>has_throat_color = black,</i> <i>has_primary_color = black,</i> <i>has_nape_color = black,</i> <i>has_forehead_color = black</i>	neuron_id = 131 <i>has_throat_color = black</i>	neuron_id = 116 <i>has_breast_color = black</i>	neuron_id = 50 <i>has_underparts_color = black</i>
Yellow-headed blackbird			<i>has_eye_color = black,</i> <i>has_throat_color = yellow,</i> <i>has_wing_color = black,</i> <i>has_breast_color = yellow,</i> <i>has_bill_color = black</i>	neuron_id = 145 <i>has_eye_color = black</i>	neuron_id = 126 <i>has_throat_color = yellow</i>	neuron_id = 20 <i>has_wing_color = black</i>
Northern Fulmer			<i>has_forehead_color = white,</i> <i>has_crown_color = white,</i> <i>has_throat_color = white,</i> <i>has_bill_shape = hooked_seabird,</i> <i>has_nape_color = white</i>	neuron_id = 305 <i>has_crown_color = white</i>	neuron_id = 132 <i>has_throat_color = white</i>	neuron_id = 4 <i>has_bill_shape = hooked_seabird</i>
GT Class: Yellow bellied Flycatcher			<i>has_eye_color = black,</i> <i>has_bill_length = shorter_than_head,</i> <i>has_shape = perching_like,</i> <i>has_underparts_color = yellow,</i> <i>has_primary_color = yellow</i>	neuron_id = 126 <i>has_throat_color = yellow</i>	neuron_id = 45 <i>has_underparts_color = yellow</i>	neuron_id = 111 <i>has_breast_color = yellow</i>
Predicted Class: Yellow throated Vireo			<i>has_throat_color = yellow,</i> <i>has_underparts_color = yellow,</i> <i>has_breast_color = yellow,</i> <i>has_primary_color = yellow,</i> <i>has_belly_color = yellow</i>	neuron_id = 145 <i>has_eye_color = black</i>	neuron_id = 151 <i>has_bill_length shorter than head</i>	neuron_id = 235 <i>has_shape perching like</i>

**Fig. 4.** Success and failure cases for unseen classes using explanations for NIWT: Success cases: (a) the ground truth class and image, (b) Grad-CAM visual explanations for the GT category, (c) textual explanations obtained using the inverse mapping from  $\mathbf{a}_c$  to domain knowledge, (d) most important neurons for this decision, their names and activation maps. The last 2 rows show failure cases, where the model predicted a wrong category. We show Grad-CAM maps and textual explanations for both the ground truth and predicted category. By looking at the explanations for the failure cases we can see that the model’s mistakes are not completely unreasonable.

*color yellow, has wing color black*, etc. In addition, when we focus on the individual neurons with relatively higher importance we see that individual neurons focus on the visual regions characterized by their assigned ‘names’. This shows that our neuron names are indeed representative of the concepts learned by the network and are well grounded in the image.

Consider the misclassified examples (rows 7 and 8). Looking at the regions in the image corresponding to the intersection of the attributes in the textual explanations for the ground truth as well as the predicted class, we can see that the network was unable to focus on the primary discriminative attributes. Similarly, the neuron names and corresponding activations have a mismatch with the predicted class with the activation maps focusing on a ‘yellowish’ area rather than a visual region corresponding to a fine-grained attribute.

## 7 Conclusion

To summarize, we propose an approach we refer to as Neuron Importance-aware Weight Transfer (NIWT), that learns to map domain knowledge about novel classes directly to classifier weights by grounding it into the importance of network neurons. Our weight optimization approach on this grounding results in classifiers for unseen classes which outperform existing approaches on the generalized zero-shot learning benchmark. We further demonstrate that this grounding between language and neurons can also be learned in reverse, linking neurons to human interpretable semantic concepts, providing visual and textual explanations.

**Acknowledgements.** We thank Yash Goyal and Nirbhay Modhe for help with figures; Peter Vajda and Manohar Paluri for helpful discussions. This work was supported in part by NSF, AFRL, DARPA, Siemens, Google, Amazon, ONR YIPs and ONR Grants N00014-16-1-{2713,2793}. 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.

## References

1. Akata, Z., Perronnin, F., Harchaoui, Z., Schmid, C.: Label-embedding for attribute-based classification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 819–826 (2013)
2. Akata, Z., Perronnin, F., Harchaoui, Z., Schmid, C.: Label-embedding for image classification. *IEEE Trans. Pattern Anal. Mach. Intell.* **38**(7), 1425–1438 (2016)
3. Akata, Z., Reed, S., Walter, D., Lee, H., Schiele, B.: Evaluation of output embeddings for fine-grained image classification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2927–2936 (2015)
4. Bau, D., Zhou, B., Khosla, A., Oliva, A., Torralba, A.: Network dissection: quantifying interpretability of deep visual representations. In: *Computer Vision and Pattern Recognition* (2017)
5. Changpinyo, S., Chao, W.L., Gong, B., Sha, F.: Synthesized classifiers for zero-shot learning. In: *CVPR* (2016)



6. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: a large-scale hierarchical image database. In: CVPR (2009)
7. Elhoseiny, M., Elgammal, A., Saleh, B.: Write a classifier: predicting visual classifiers from unstructured text. *IEEE Trans. Pattern Anal.* **PP**(99), 1 (2017)
8. Elhoseiny, M., Saleh, B., Elgammal, A.: Write a classifier: zero-shot learning using purely textual descriptions. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 2584–2591 (2013)
9. Elhoseiny, M., Zhu, Y., Zhang, H., Elgammal, A.: Link the head to the “beak”: zero shot learning from noisy text description at part precision. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017
10. Farhadi, A., Endres, I., Hoiem, D., Forsyth, D.: Describing objects by their attributes. In: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2009, pp. 1778–1785. IEEE (2009)
11. Frome, A., et al.: DeViSE: a deep visual-semantic embedding model. In: Advances in Neural Information Processing Systems (2013)
12. Goyal, Y., Mohapatra, A., Parikh, D., Batra, D.: Interpreting visual question answering models. CoRR abs/1608.08974 (2016). <http://arxiv.org/abs/1608.08974>
13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
14. Kodirov, E., Xiang, T., Gong, S.: Semantic autoencoder for zero-shot learning. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017
15. Konam, S., Quah, I., Rosenthal, S., Veloso, M.: Understanding convolutional networks with apple : automatic patch pattern labeling for explanation. In: First AAAI/ACM Conference on AI, Ethics, and Society (2018)
16. Lampert, C.H., Nickisch, H., Harmeling, S.: Attribute-based classification for zero-shot visual object categorization. *IEEE Trans. Pattern Anal. Mach. Intell.* **36**(3), 453–465 (2014)
17. Larochelle, H., Erhan, D., Bengio, Y.: Zero-data learning of new tasks. In: AAAI, vol. 1, p. 3 (2008)
18. Lei Ba, J., Swersky, K., Fidler, S., et al.: Predicting deep zero-shot convolutional neural networks using textual descriptions. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 4247–4255 (2015)
19. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: NIPS (2013)
20. Nair, V., Hinton, G.E.: Rectified linear units improve restricted boltzmann machines. In: ICML, pp. 807–814 (2010)
21. Norouzi, M., et al.: Zero-shot learning by convex combination of semantic embeddings. In: ICLR (2014)
22. Novak, R., Bahri, Y., Abolafia, D.A., Pennington, J., Sohl-Dickstein, J.: Sensitivity and generalization in neural networks: an empirical study. *arXiv preprint arXiv:1802.08760* (2018)
23. Qiao, R., Liu, L., Shen, C., Hengel, A.V.D.: Less is more: zero-shot learning from online textual documents with noise suppression. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
24. Reed, S., Akata, Z., Lee, H., Schiele, B.: Learning deep representations of fine-grained visual descriptions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 49–58 (2016)
25. Romera-Paredes, B., Torr, P.: An embarrassingly simple approach to zero-shot learning. In: Proceedings of The 32nd International Conference on Machine Learning, pp. 2152–2161 (2015)



26. Selvaraju, R.R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., Batra, D.: Grad-CAM: why did you say that? visual explanations from deep networks via gradient-based localization. In: ICCV (2017)
27. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR (2015)
28. Socher, R., Ganjoo, M., Manning, C.D., Ng, A.: Zero-shot learning through cross-modal transfer. In: Advances in Neural Information Processing Systems (2013)
29. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The Caltech-UCSD Birds-200-2011 Dataset. Technical report CNS-TR-2011-001. California Institute of Technology (2011)
30. Xian, Y., Akata, Z., Sharma, G., Nguyen, Q., Hein, M., Schiele, B.: Latent embeddings for zero-shot classification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2016)
31. Xian, Y., Schiele, B., Akata, Z.: Zero-shot learning - the good, the bad and the ugly. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017
32. Yu, R., et al.: NISP: pruning networks using neuron importance score propagation. In: CVPR (2018)
33. Zeiler, M.D., Fergus, R.: Visualizing and understanding convolutional networks. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) ECCV 2014. LNCS, vol. 8689, pp. 818–833. Springer, Cham (2014). [https://doi.org/10.1007/978-3-319-10590-1\\_53](https://doi.org/10.1007/978-3-319-10590-1_53)
34. Zhang, L., Xiang, T., Gong, S.: Learning a deep embedding model for zero-shot learning. In: CVPR (2017)
35. Zhang, Z., Saligrama, V.: Zero-shot learning via semantic similarity embedding. In: Proceedings of the IEEE International Conference on Computer Vision (2015)
36. Zhou, B., Khosla, A., Lapedriza, À., Oliva, A., Torralba, A.: Object detectors emerge in deep scene CNNs. CoRR abs/1412.6856 (2014)