

MACE: Model Agnostic Concept Extractor for Explaining Image Classification Networks

Ashish Kumar ^{*}, Karan Sehgal [†], Prerna Garg [‡], Vidhya Kamakshi [§] and Narayanan C Krishnan [¶]

IIT Ropar

Abstract

Deep convolutional networks have been quite successful at various image classification tasks. The current methods to explain the predictions of a pre-trained model rely on gradient information, often resulting in saliency maps that focus on the foreground object as a whole. However, humans typically reason by dissecting an image and pointing out the presence of smaller concepts. The final output is often an aggregation of the presence or absence of these smaller concepts. In this work, we propose MACE: a Model Agnostic Concept Extractor, which can explain the working of a convolutional network through smaller concepts. The MACE framework dissects the feature maps generated by a convolution network for an image to extract concept based prototypical explanations. Further, it estimates the relevance of the extracted concepts to the pre-trained model's predictions, a critical aspect required for explaining the individual class predictions, missing in existing approaches. We validate our framework using VGG16 and ResNet50 CNN architectures, and on datasets like Animals With Attributes 2 (AWA2) and Places365. Our experiments demonstrate that the concepts extracted by the MACE framework increase the human interpretability of the explanations, and are faithful to the underlying pre-trained black-box model.

1 Introduction

The state of the art convolutional networks has been quite successful at various computer vision tasks. However, the improved performance has come at the cost of reduced human understanding of the model. It is crucial to bring transparency in these networks to help understand their decisions. Recently, there has been a lot of work on designing models that explain the behavior of a pre-trained convolutional network. These approaches generate saliency maps for explaining the model's behavior or explain the prediction based on training instances. Approaches such as GradCAM Selvaraju et al. [2017] and its variants, use gradient information to generate saliency maps as illustrated in Figure 1. Other approaches like Excitation Backpropagation Zhang et al. [2016] use top-down neural attention to generate attention maps. These explanations provide a high-level insight into the working of the model. However, we observe that invariably almost the entire object is highlighted in all the saliency maps and their explanations are almost identical for the predicted class and any other class for a given image. This reduces the interpretability of their explanations. While the explanations can detect the foreground object, they are unable to accurately highlight regions in the object that contributed towards the prediction.

^{*}2016csb1033@iitrpr.ac.in

[†]2016csb1080@iitrpr.ac.in

[‡]2016csb1050@iitrpr.ac.in

[§]2017csz0005@iitrpr.ac.in

[¶]ckn@iitrpr.ac.in

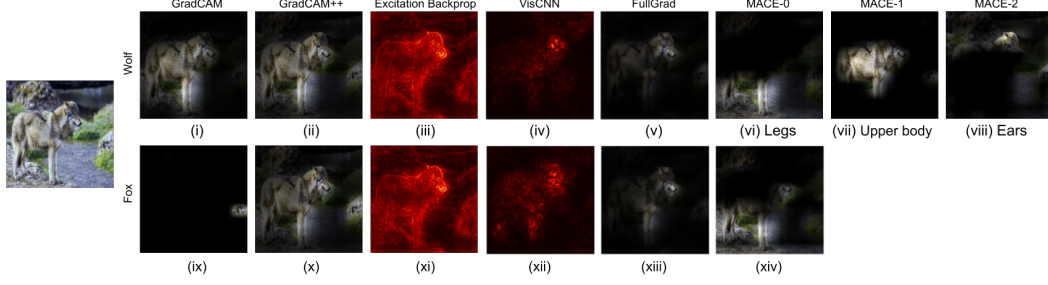


Figure 1: [Best viewed in color] Explanations generated by different approaches for the predicted class and another class for an image of wolf.

Another style of reasoning for a prediction is by dissecting an image and referring to the presence of smaller concepts. For example, if the image is that of a lion, the smaller concepts could be legs, ears, face, skin texture, etc. The aggregation of these smaller concepts is used to explain the final output. The fundamental objective of our framework: Model-Agnostic Concept Extractor (MACE) is to mimic this style of reasoning; which is to explain the model’s behavior in terms of smaller concepts in the image. Our approach learns multiple concepts for a class without any supervision or additional information. These concepts are visualized through multiple high quality localized saliency maps, rather than a single region. Figure 2 shows the visualizations of some of the concepts learned by our approach.

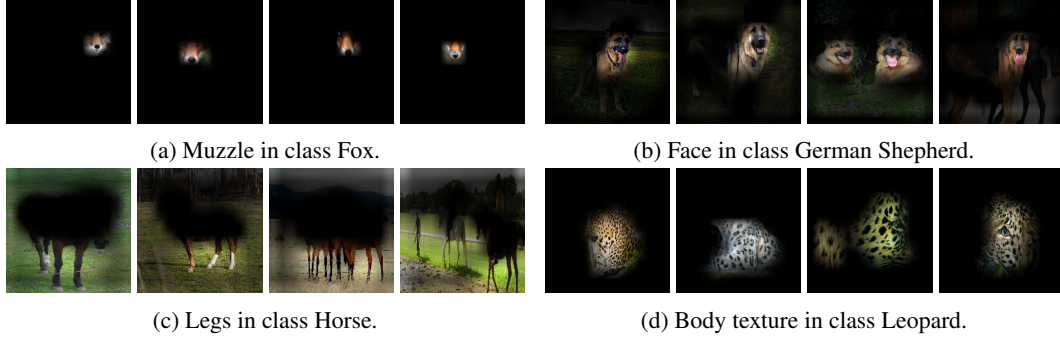


Figure 2: [Best viewed in color] Visualization of concepts across multiple images of the same class.

Further, the proposed framework also estimates the relevance of each concept with respect to the output of the model. This is particularly helpful when the explanation for any output of the model is desired. In Figure 1 the first row corresponds to the explanations from different approaches for the predicted class and the second row presents the explanations for the class with second-highest prediction probability. It can be noticed that there is no significant difference in the explanations generated by current approaches for the two outputs of the model. However, our approach can generate class-specific concepts as well as their relevance for a more comprehensive explanation. For an incorrect class, most of the concepts have negative relevance thereby suppressing the prediction probability. As illustrated in Figure 1 while the predicted class is the wolf, the second-highest probability is for class fox. The MACE framework extracts only a single concept (that of legs) with positive relevance for the class fox.

We also ensure that the MACE framework generates explanations that satisfy a few desirable properties. *Stability*: The learned concepts should be consistent i.e. the localized saliency maps generated for a particular concept should be similar across all the images of the same class, thus achieving higher interpretability. This can be observed in Figure 2 where the visualizations for a single concept obtained from different images appear very similar. *Faithfulness*: The explanations generated by the framework should accurately represent the working of the pre-trained model. This will improve the trustworthiness of the explanations. Finally, *Robustness*: The explanations should be robust to perturbations such as noise, translation, and rotation in the input. Explanations of perturbed inputs should not change significantly if there is no significant change in the corresponding outputs.

2 Related Work

The last few years have seen significant work on developing post-hoc explanation models Selvaraju et al. [2017], Chattopadhyay et al. [2018], Zhou et al. [2016], Ribeiro et al. [2016, 2018], Lundberg and Lee [2017]. Some of the popular approaches construct saliency maps that highlight important regions in the image. The Class Activation Map (CAM) Zhou et al. [2016] is one of the earliest approaches to construct a saliency map specific to a classification model. The gradient-class activation maps (Grad-CAM) Selvaraju et al. [2017] is a generalization of CAM that uses gradients flowing into the final convolutional layer to localize salient image regions. GradCAM is among the most popular approach for generating explanations for an image classification network. However, it has multiple drawbacks including poor localization in the presence of multiple instances of the same object in an image and poor distinction of the class-specific saliency maps generated for different classes, given an input Wang et al. [2019b]. GradCAM++ Chattopadhyay et al. [2018], a variant of GradCAM, overcomes the limitation of generating accurate saliency maps for multiple instances of an object in an image. However, GradCAM++ requires the estimation of higher-order derivatives, the cost of which grows with increasing network complexity.

Vis-CNN Simonyan et al. [2013] also uses gradients to explain the contribution of each input pixel towards the output. But pixel-wise quantification appears like a foreground detector looking at most of the results presented in the paper and also during our experiments. Full Grad Srinivas and Fleuret [2019] combines the idea of using gradients flowing back to intermediary layers as well as gradients flowing back till input to generate a saliency map attributing importance of each pixel towards a prediction. On the other hand, Wang et al Wang et al. [2019b] suggest a possible compromise in the faithfulness of explanations when gradients are used to generate the explanations. There have been other variants of CAM like Ablation-CAM Desai and Ramaswamy [2020] and Score-CAM Wang et al. [2019a] that avoid the need for gradients for computing the saliency maps. Our approach also estimates saliency maps for explaining the pre-trained model without utilizing the gradient information. Further, we automatically synthesize multiple maps each characterizing a different concept present in the image.

Fong et al, Fong and Vedaldi [2017] propose a perturbation based approach for estimating the salient regions in the image, that was further extended to also take into account when the region was preserved Fong et al. [2019]. These complementary approaches aim to modify the input image successively and map the region integral for the final prediction. There are many challenges with the feasibility of the approach given the high dimensional nature of the images and voluminous amount of possible perturbations. Excitation backpropagation Zhang et al. [2016] proposes a probabilistic attention mechanism to explain the working of black boxes in a fine-grained manner. But as recent works like Mohankumar et al. [2020], Xu et al. [2020], Grimsley et al. [2020] state, attention need not be human interpretable .

Another category of explainable models is the ante-hoc models that are explainable by design Li et al. [2018], Chen et al. [2019], Hase et al. [2019], Chen et al. [2020]. Due to explainability being a parameter considered in the design, the explanations are faithful to the underlying model. However, the need to retrain the model from scratch to incorporate the explainability aspect is an obstacle to explain a pre-trained model that has already been deployed.

Our proposed approach is a posthoc explainability technique like Zhou et al. [2016], Selvaraju et al. [2017], Chattopadhyay et al. [2018] that leverages some aspects from ante-hoc explainability techniques like Chen et al. [2019], Li et al. [2018] to generate fine-grained and faithful explanations. Like GradCAM, our approach generates human interpretable explanations in the form of a saliency map. However, unlike GradCAM, we do not utilize the gradient backflow information for generating the saliency maps. Our work is closely related to the exciting paradigm of ‘this looks like that’ proposed by Chen et al Chen et al. [2019]. Specifically, the MACE framework generates multiple saliency maps containing different parts/concepts for explaining the output of a single image. However, unlike Chen et al. [2019], our approach does not require modifying and retraining the black box architecture nor is there a need to slide across the activation maps or tune additional parameters to visualize a concept.

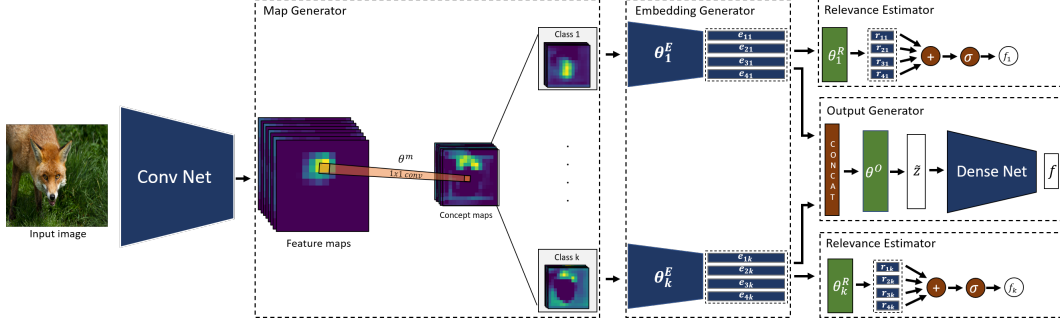


Figure 3: [Best viewed in color] Architecture of the MACE Framework.

3 Methodology

The MACE framework is envisioned as a modular network that can be attached to any convolution layer of a pre-trained network for probing its functionality. However, for describing the framework, we assume that the MACE network is inserted in between the final convolutional layer and the first dense layer of a pre-trained network. This lateral connection taps into the spatial distribution of concepts (as gathered from the downstream convolutional layers) and the discriminative features for classification (from the upstream dense layers). Given, a query image and a class label, the MACE framework extracts the class-specific low-level concepts along with their relevance towards the output of the model.

The process of learning and extracting the concepts and their relevance is divided across four modules namely; the map generator, embedding generator, relevance estimator, and output generator. The MACE framework is illustrated in Figure 3. The map generator learns an attended spatial map representing the spread of a concept using the output of the last convolutional layer of the pre-trained model. The embedding generator transforms the manifestation of a concept into a latent representation that is invariant to spatial information. The relevance estimator determines the importance of a class-specific concept towards the output of the model. Finally, the output generator completes the loop by linking the extracted concepts back to the first dense layer of the pre-trained model.

Let f be a pre-trained model for a K -way image classification task; trained using a dataset \mathcal{D} . The output of the last convolutional layer in f for a given input is denoted as \mathbf{x} . Thus $\mathbf{x} \in \mathbb{R}^{H \times W \times D}$, where H and W refer to the size of the feature map and D refers to the number of filters in the last convolution operation. Let $\mathbf{z} \in \mathbb{R}^L$ represent the output of the dense layer following the convolutional layer. The MACE framework is laterally connected in between the transformation of \mathbf{x} to \mathbf{z} .

For the sake of simplicity, we assume that every class can be explained by the same number of concepts. (In practice as discussed in the supplementary material section S2, pruning may result in a varying number of concepts per class). A concept has two aspects to it. The first is its actual manifestation in an image, termed as the concept map, for example, the region corresponding to the ear or legs. The second aspect is the latent representation, termed as the concept embedding, that is invariant to the manifestation. Irrespective of the location and orientation of the concept ear in an image, the encoding of the concept remains the same. $\mathbf{c}_{jk} \in \mathbb{R}^{H \times W}$ denotes the j^{th} concept map for the k^{th} class and $\mathbf{e}_{jk} \in \mathbb{R}^Q$ denotes the corresponding concept embedding.

3.1 Map Generator

The map generator takes as input the convolutional feature map \mathbf{x} and estimates a concept map \mathbf{c}_{jk} . The concept map encodes two properties - the activation pattern of a concept and the salient region in the image that expresses the activation pattern. The activation pattern is used by the embedding generator to extract a distinct and invariant encoding for the concept. The salient region in the image is used for visualizing the concept post-training, facilitating the human interpretability of the concept.

Due to the nature of convolution operations, the activations of a concept are distributed across the different channels of the feature map \mathbf{x} . The MACE framework combines the channels in the feature map to produce a single activation pattern for a concept. This is achieved as a weighted average of

the channels in the feature map. Specifically, we employ 1D convolutions to obtain the concept map. Let θ_{jk}^M denote the weights of the 1D convolution filter for generating the j^{th} concept map of the k^{th} class, then $\mathbf{c}_{jk} = \text{ReLU}(\mathbf{x} * \theta_{jk}^M)$. The ReLU operation is performed because, the locations in the activation map corresponding to large positive values may denote the presence of the concept Selvaraju et al. [2017], Chattopadhyay et al. [2018]. We have to learn $K \times C$ 1D convolution filters, assuming C concepts per class. The concept maps \mathbf{c}_{jk} are of the same size as the input feature map \mathbf{x} .

Each of these concept maps represents an activation pattern. The region in the concept map exhibiting the highest activation is most likely to have expressed the concept. This region in the original image contributing to this activation can be visualized by resizing the concept map to the original image dimension. We use this resized concept map to visualize the concept. The MACE framework does not impose any constraints on the size of the concept. Thus allowing it to learn small concepts such as ears to large concepts like body and background.

3.2 Embedding Generator

The next module - the embedding generator takes as input the concept maps that are stacked class-wise. The map generator only uses the convolution layer activation maps to generate the concept map and thus retains the spatial information about the location of the concept in the image. Consider two images one in which the object lion is present on the right half of the image, and a different orientation of the object lion is present in the left half of the second image. The concept maps of these two images will have high activations in different regions even though the underlying concept is the same. Thus, we need to abstract out the concept from its actual manifestation. We refer to this encoding as the concept embedding. We would like the embeddings for a concept extracted from different images belonging to the same class exhibiting the concept to form a tight cluster.

The abstraction of the concept through the embedding also enables a comparison of different concepts. We model the embedding generator as a multi-layer dense network. Having a sufficiently large network enables the extraction of different concepts. However, the parameters of the network are shared across the concepts of a class to reduce the overall learning complexity. Formally, let the embedding generator network for class k , be parameterized by θ_k^E . A stack of the concept maps $\{\mathbf{c}_{jk}\}_{j=1}^C$ extracted from an input image are fed as inputs to embedding generator resulting in the set of embeddings $\{\mathbf{e}_{jk}\}_{j=1}^C$. Restricting to only a positive subspace (as was the case with the concept map generation) is not necessary for learning the concept embedding. We use tanh activations enabling the use of the full space to learn well-separated embeddings.

Triplet loss is used to learn the invariant concept embeddings. We normalize the embeddings before applying the triplet loss. We use the training images belonging to class k from a mini-batch to learn the C concept embeddings for the class. The embeddings generated for concept j across the training images of a particular class are chosen as anchor positives. The embeddings for the other concepts (except j) for a specific training image are the anchor negatives. Similar to Schroff et al. [2015], we use all anchor-positive pairs and select semi-hard negatives for anchor-negative pairs. The margin α is set to 1 to make the embeddings orthogonal to each other. Let $\mathbf{e}_{jk}^a(i)$ represent an anchor for the i^{th} image in the mini-batch, the anchor positives are represented as $\mathbf{e}_{jk}^p(i)$ and the anchor negatives as $\mathbf{e}_{jk}^n(i)$. Then the triplet loss for class k using a mini-batch of B images is defined as

$$\mathcal{L}_k^E = \sum_{i=1}^B \sum_{j=1}^C \left[\|\mathbf{e}_{jk}^a(i) - \mathbf{e}_{jk}^p(i)\|_2^2 - \|\mathbf{e}_{jk}^a(i) - \mathbf{e}_{jk}^n(i)\|_2^2 + \alpha \right]_+ \quad (1)$$

3.3 Relevance Estimator

A key aspect of the MACE framework is the estimation of the relevance of the concept towards the model prediction. This score, termed as concept relevance, quantifies the contribution of a concept towards the output of the model for a given class. The concept relevance would enable a comparison of different concepts for a given class and provide better insight into the relationship between the concepts and the model's output.

Given a training image, the concept embeddings of class k are concatenated and passed through a sigmoid activated dense layer (parameterized by θ_k^R) with a single output to learn the probability

of classification for class k as estimated by the pre-trained model. θ_k^R can be divided into chunks representing the connection weights for the individual concepts, i.e. $\theta_k^R = [\theta_{1k}^R, \dots, \theta_{Ck}^R]$. Then, the concept relevance for concept j with respect to the output for class k is denoted as $r_{jk} = \theta_{jk}^{R^T} \mathbf{e}_{jk}$. Note that the concept relevance $r_{jk} \in (-\infty, \infty)$. As the argument to the sigmoid activation of the dense layer is $\sum_{j=1}^C r_{jk}$, the impact of the concept relevance scores on the prediction probabilities can be easily understood. Positive (negative) concept relevance increases (decreases) the prediction probability, thus emphasizing (reducing) the importance of the concept for a particular prediction.

The parameters of the dense layer, θ_k^R , are learned by minimizing the cross-entropy between the output of the sigmoid and the pre-trained model’s prediction probability for class k . Specifically, for a mini-batch of training instances, the loss for learning the parameters is defined as

$$\mathcal{L}_k^R = - \sum_{i=1}^B f_k(i) \log \left(\sigma \left(\sum_{j=1}^C r_{jk}(i) \right) \right) \quad (2)$$

where $f_k(i)$ and $r_{jk}(i)$ are the pre-trained models output for class k and the concept relevance for concept j with respect to class k respectively for the training instance i .

3.4 Output Generator

To increase the faithfulness of the explanations we loop back the concept embeddings into the pre-trained model. The concatenated concept embeddings are passed through a dense layer (parameterized by θ^O) to approximate the output, \mathbf{z} , of the first dense layer of the pre-trained model. An L_2 loss between the approximation, $\tilde{\mathbf{z}}$, and \mathbf{z} , defined as $\mathcal{L}^D = \|\mathbf{z} - \tilde{\mathbf{z}}\|_2^2$ is used to learn the approximation. If the approximation is accurate we should obtain the output $f(\mathbf{z})$ (\mathbf{x} is replaced by \mathbf{z} , a slight abuse of the notation), when $\tilde{\mathbf{z}}$ is passed through the dense layers of the pre-trained model. This requirement is enforced by minimizing the divergence between the pre-trained model’s outputs for $\tilde{\mathbf{z}}$ and \mathbf{z} defined as $\mathcal{L}^O = KL(f(\tilde{\mathbf{z}}) \| f(\mathbf{z}))$. The L_2 loss is required to faithfully mimic the behavior of the pre-trained model. On the other hand, leaving out the divergence loss can still lead to a good approximation $\tilde{\mathbf{z}}$, but $f(\tilde{\mathbf{z}})$ can be inconsistent with $f(\mathbf{z})$.

Overall, the MACE framework minimizes $\mathcal{L} = \sum_{k=1}^K (\mathcal{L}_k^E + \mathcal{L}_k^R) + \mathcal{L}^D + \mathcal{L}^O$ with respect to the parameters θ^M , θ^E , θ^R , and θ^O . Adam optimizer Kingma and Ba [2014] is used to minimize the loss function.

4 Experiments

We validate the MACE framework on VGG Simonyan and Zisserman [2014], and ResNet He et al. [2015] architectures trained on the Imagenet dataset Russakovsky et al. [2014]. The networks are fine-tuned on AWA2 dataset AWA, Xian et al. [2018] and the Places365 dataset Zhou et al. [2017]. We discuss the results for the AWA2 dataset using the VGG model in the paper and present the results for the Places365 dataset and ResNet architecture in the supplementary material (sections S8 and S7 respectively). The AWA2 dataset consists of 37322 images of 50 animal classes. We select a subset of 10 classes with approximately 400 images per class for our experiments. The VGG network that is fine-tuned with this dataset is our pre-trained model. For every class, we aim to learn 10 concepts ($C = 10$) with the concept embedding dimension to be 32 ($Q = 32$). Examples of the concepts extracted for various classes are presented in Figure 2. The architectural details of the embedding generator, along with the choices for the hyper-parameters are presented in the supplementary material section S1. Additional examples of the visualization of the concepts for the AWA2 dataset are presented in sections S3 and S4 of the supplementary material. While we describe the results from three salient experiments here, investigations into the stability and robustness of the explanations as well as ablation studies are discussed in sections S5, S6, and S9 of the supplementary material respectively.

4.1 Faithfulness of the Explanations

Faithfulness, the ability to accurately explain the working of a pre-trained model, is one of the most important traits of a post-hoc explanation model. As suggested in prior literature, we measure

faithfulness using the drop in the classification probability upon masking. We generate a binary mask by thresholding each of the concept maps. As the concepts are trained to be different from each other, masking out just one concept might not give a significant drop in the probability of the predicted class, as the model can still predict with the help of other concepts. For example, as seen in Figure 4a masking the ear concept does not alter the classification probabilities significantly as the model compensates for it through the other concepts. Thus, we take the union of all binary masks. The region in the generated mask is removed from the image and the resulting drop in the probability of the predicted class is measured. We also compare the drop in the prediction probability when the mask is synthesized using other saliency-map generating approaches as well as a random process (activation maps combined randomly). The results of this experiment are presented in Figure 4b. The highest drop in the probability is observed for the mask synthesized by the MACE framework for every threshold.

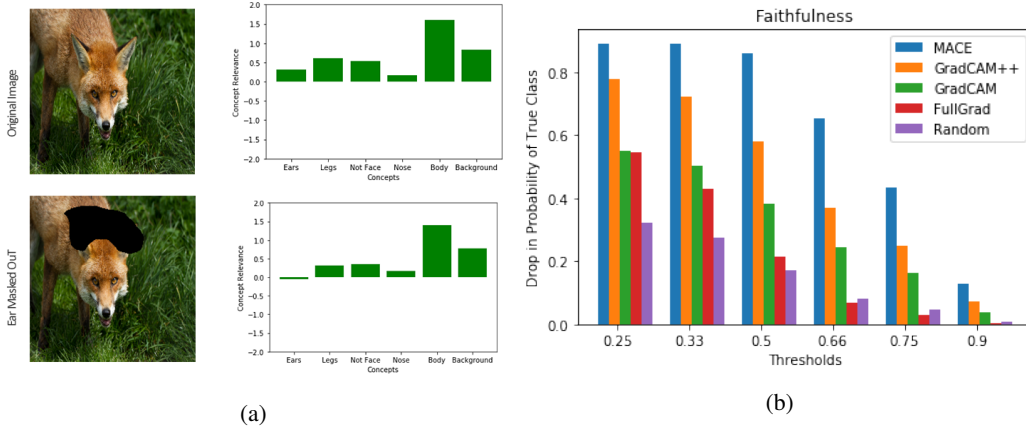


Figure 4: [Best viewed in color] Faithfulness of the explanations extracted by MACE (a) No change in the prediction distribution after masking a concept due to compensation by the model using other concepts, (b) Drop in the probability of classification across different threshold values for generating the binary mask.

4.2 Explaining all the outputs for a test image

An image classification network outputs a probability distribution over the class labels. An important aspect of any explanation model is to explain all the non-zero probabilities in the output of the classifier. In our approach, this is explained by looking at the concepts generated for every class, and the corresponding concept relevances. We specifically look at the concepts that have positive concept relevance for any class. Figure 5a shows the results for a few images. For example, the first two rows present the explanations generated by different models for why the classifier has non-zero probabilities for the German Shepherd and Wolf classes when the correctly predicted class is Fox. The original image in the first column of 5a is a fox image, the visualization of explanation given by MACE (second column) for German Shepherd class and wolf class are meaningful and highlight smaller regions in the fox such as ears and legs, while other approaches give the same explanation for both the classes. We can also observe from the figure that MACE consistently gives different meaningful explanations for different classes, while other approaches give the same or no explanation.

We further analyze the behavior of the concept relevance scores for different outputs of the pre-trained model. We first compute the average concept relevance for a concept embedding e_{jk} using all the images belonging to class k . Over the test set, we compute the percentage of concepts with positive concept relevance less than the average score using only the images for which class k is the second-highest ranked class according to the pre-trained model’s predictions. We compute the average percentage across all classes for a particular rank. The result is summarized in Figure 5b. We can observe from the figure that the percentage of concepts with positive relevance score less than the average increases as the prediction ranking decreases. This is intuitive as concept relevance for the concept of an incorrect class should decrease as we go down the class ranking.

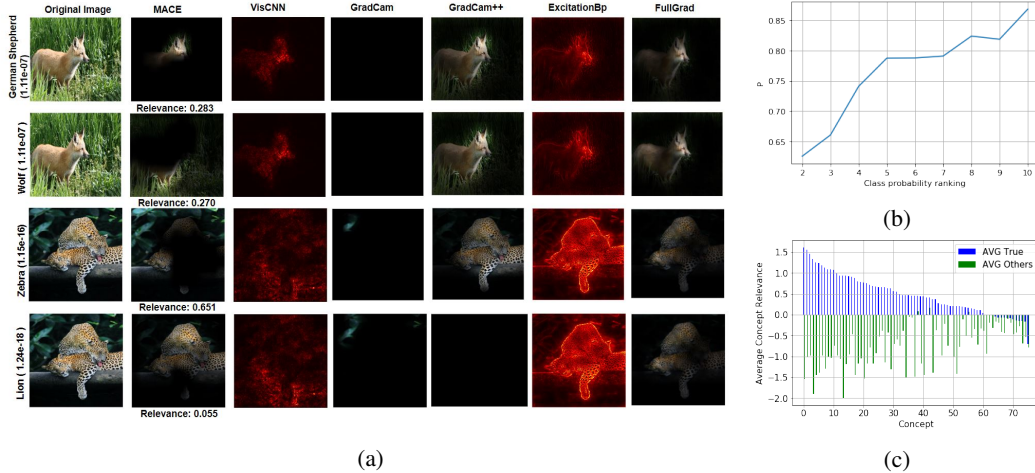


Figure 5: [Best viewed in color] (a) contains explanations of different approaches for incorrect class. Each row corresponds to the explanation for some output probability for a wrong class. The wrong class is written on the left-most side in each row. The image shown for MACE approach is the visualization of the concept with the highest concept relevance for that class, and the concept relevance is written below the visualization; (b) Percentage of concepts with positive concept relevance less than the average score using only the images for different classes (See text for more details); (c) Average relevances of all the concepts for the true class and for all the classes except the true class.

Table 1: Preferences of the Human Subjects for Explaining the Classifier’s Predictions.

Approach	Vote Percentage
MACE (ours)	48.29%
Excitation Backpropagation Zhang et al. [2016]	40.73%
GradCAM++ Chattopadhyay et al. [2018]	9.51%
GradCAM Selvaraju et al. [2017]	1.21%
VisCNN Simonyan et al. [2013]	0.24%

Figure 5c shows the average concept relevances for all the concepts where the average is split between relevances of concepts belonging to the true class of a test image (AVG True) and relevances of concepts belonging to any other class except the true class (AVG Others). It can be observed that AVG Others (green bars) is significantly lower than AVG True (blue bars) for all the concepts. This further strengthens the credibility of concept relevances as concept relevance should be higher only if the concept actually belongs to the true class of the image than any other class.

4.3 Human Evaluation of Explanations

We evaluate the quality of our explanations generated for the VGG16 network fine-tuned on 10 classes of the AWA2 dataset. We present to every user the saliency maps generated by MACE, GradCAM Selvaraju et al. [2017], GradCAM++ Chattopadhyay et al. [2018], Excitation Backpropagation Zhang et al. [2016] and VisCNN Simonyan et al. [2013] for a random set of 10 images. In addition, we randomly select 2 images from this set for repetition and create a final questionnaire of 12 examples and the corresponding explanations. The subjects were asked to choose the approach whose explanation helped them better understand the classifier’s prediction. The responses in which the answers to the repeated questions mismatched were removed to maintain the consistency in the responses. We received 41 consistent responses, resulting in a total of 410 votes. The results have been summarised in Table 1. We observe that explanations from MACE are the most preferred choice amongst the participants. Further, approaches like MACE and Excitation Backpropagation are preferred over gradient-based approaches like VisCNN, GradCAM and GradCAM++.

5 Summary

In this work, we define a new form of reasoning for explaining the outputs of an image classification network using multiple concepts/parts of an object. We present the MACE framework, for extracting and visualizing these multiple concepts. We also propose a mechanism for estimating the relevance of a concept towards the output of the model. We perform extensive experiments on the MACE framework using VGG16 and ResNet50 architectures for animal and places classification tasks. Our results confirm the faithfulness of the explanations as well as their human interpretability.

References

- Animals with attributes 2 (awa2) dataset. URL <https://cvml.ist.ac.at/AwA2/>.
- Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 839–847. IEEE, 2018.
- Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K Su. This looks like that: deep learning for interpretable image recognition. In *Advances in Neural Information Processing Systems*, pages 8928–8939, 2019.
- Zhi Chen, Yijie Bei, and Cynthia Rudin. Concept whitening for interpretable image recognition. *arXiv preprint arXiv:2002.01650*, 2020.
- Saurabh Desai and Harish Guruprasad Ramaswamy. Ablation-cam: Visual explanations for deep convolutional network via gradient-free localization. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 983–991, 2020.
- Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Understanding deep networks via extremal perturbations and smooth masks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2950–2958, 2019.
- Ruth C Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3429–3437, 2017.
- Christopher Grimsley, Elijah Mayfield, and Julia RS Bursten. Why attention is not explanation: Surgical intervention and causal reasoning about neural models. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 1780–1790, 2020.
- Peter Hase, Chaofan Chen, Oscar Li, and Cynthia Rudin. Interpretable image recognition with hierarchical prototypes. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 7, pages 32–40, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. URL <http://arxiv.org/abs/1512.03385>.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Oscar Li, Hao Liu, Chaofan Chen, and Cynthia Rudin. Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in neural information processing systems*, pages 4765–4774, 2017.
- Akash Kumar Mohankumar, Preksha Nema, Sharan Narasimhan, Mitesh M Khapra, Balaji Vasan Srinivasan, and Balaraman Ravindran. Towards transparent and explainable attention models. *arXiv preprint arXiv:2004.14243*, 2020.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144, 2016.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014. URL <http://arxiv.org/abs/1409.0575>.

- Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2014.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- Suraj Srinivas and François Fleuret. Full-gradient representation for neural network visualization. In *Advances in Neural Information Processing Systems*, pages 4126–4135, 2019.
- Haofan Wang, Mengnan Du, Fan Yang, and Zijian Zhang. Score-cam: Improved visual explanations via score-weighted class activation mapping. *arXiv preprint arXiv:1910.01279*, 2019a.
- Yulong Wang, Hang Su, Bo Zhang, and Xiaolin Hu. Learning reliable visual saliency for model explanations. *IEEE Transactions on Multimedia*, 2019b.
- Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2251–2265, 2018.
- W. Xu, J. Wang, Y. Wang, G. Xu, L. Daoyu, W. Dai, and Y. Wu. Where is the model looking at –concentrate and explain the network attention. *IEEE Journal of Selected Topics in Signal Processing*, pages 1–1, 2020.
- Jianming Zhang, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. Top-down neural attention by excitation backprop. *CoRR*, abs/1608.00507, 2016. URL <http://arxiv.org/abs/1608.00507>.
- Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016.
- Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.

MACE: Model Agnostic Concept Extractor for Explaining Image Classification Networks (Supplementary Material)

Ashish Kumar ^{*}, Karan Sehgal [†], Prerna Garg [‡], Vidhya Kamakshi [§] and Narayanan C Krishnan [¶]

IIT Ropar

S1 Experimental Details

Our implementation can be accessed at <https://github.com/mace19/MACE>. In this section, we discuss our experimental details for various architectures and datasets. For each dataset, we select 10 classes and for every class, we aim to learn 10 concepts with concept embedding dimension equal to 32.

S1.1 VGG16 on AWA2

We train our MACE framework on VGG16 Simonyan and Zisserman [2014] architecture using a subset of AWA2 dataset AWA, Xian et al. [2018], having 10 classes with approximately 400 images per class. The learning rate is set to 10^{-4} and the model is trained for 64 epochs using an ADAM optimizer Kingma and Ba [2014]. The accuracy of the model was 92.7%

S1.2 VGG16 on Places365

We train our MACE framework on VGG16 Simonyan and Zisserman [2014] architecture using a subset of Places365 Zhou et al. [2017], having 10 classes with approximately 5000 images per class. The learning rate is set to $5 * 10^{-4}$ and the model is trained for 32 epochs using an ADAM optimizer Kingma and Ba [2014]. The accuracy of the model was 93.1%

S1.3 ResNet-50 on AWA2

We train our MACE framework on ResNet-50 He et al. [2015] architecture using a subset of AWA2 dataset AWA, Xian et al. [2018], having 10 classes with approximately 400 images per class. The learning rate is set to 10^{-3} and the model is trained for 128 epochs using an ADAM optimizer Kingma and Ba [2014]. The accuracy of the model was 96.5%

S2 Pruning

The set of concepts generated by following the proposed methodology generate some concepts that are not meaningful. We prune these concepts so that only the meaningful concepts remain. We use the following strategy for pruning the concepts:

^{*}2016csb1033@iitrpr.ac.in

[†]2016csb1080@iitrpr.ac.in

[‡]2016csb1050@iitrpr.ac.in

[§]2017csz0005@iitrpr.ac.in

[¶]ckn@iitrpr.ac.in

- We fetch the top T images in terms of concept relevance, and if more than S of those images don't belong to the same class as the concept, then we prune the concept. In our experiments, we set $T = 10$ and $S = 5$.
- If a large part of test dataset has positive concept relevance for a concept, we prune that concept. In our experiments with AWA2 dataset, we pruned the concept if it had positive concept relevance for more than 50% of the test images.
- If the concept masks in almost the entire image, we prune the concept. In our case, we pruned the concept if on an average it masked in 95% of the image.
- If a concept has negative concept relevance for most of the images that belong to the same class as the concept, then we prune the concept. In our experiments we pruned the concept if less than 5% of the images belonging to the same class had positive concept relevance.

After pruning the concepts, we fine-tune the MACE model. We use the fine-tuned model to perform all the experiments. Visualizations of some of the pruned concepts are shown in Figure SF1

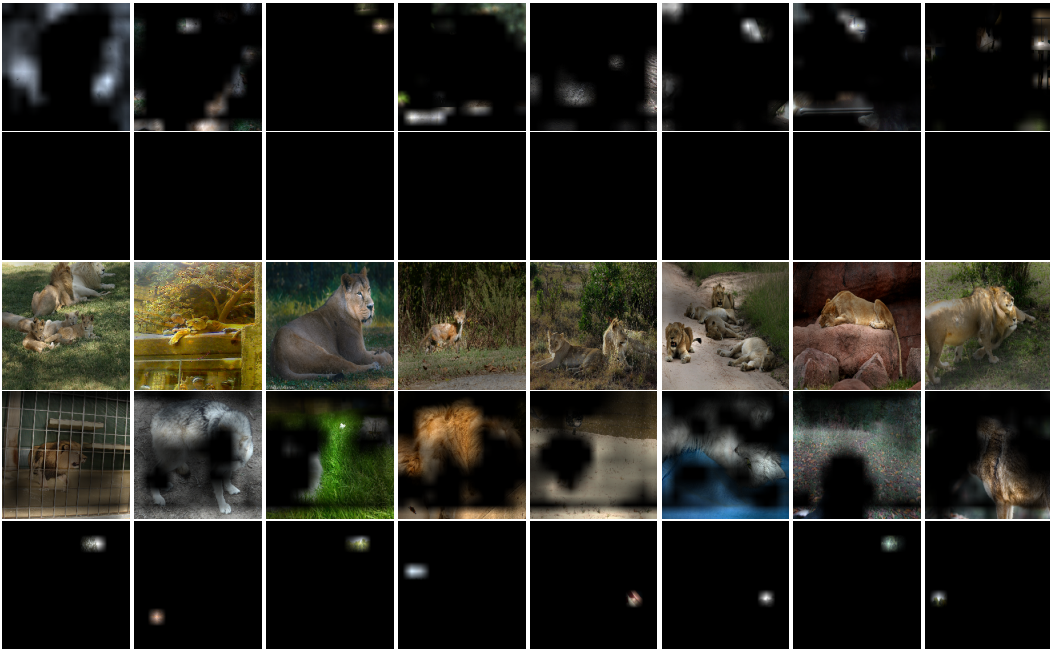


Figure SF1: [Best viewed in color] Some of the pruned concepts. Each row contains prototypical images of a pruned concept.

S3 Concept visualization and relevance for true class prediction

Figure SF2 shows the visualizations of the learned concepts and their corresponding relevance values for a few images. We observe that for most images the background concept had very high relevance thus suggesting that it plays an important contribution in the model's prediction. We also observe some repetitive pattern in the visualizations of the learned concepts. Such repetitions have been removed from this figure.



Figure SF2: [Best viewed in color] Concepts with relevances corresponding to the predicted class

S4 Concept visualization and relevance for other class prediction

One of the major contributions of our method is that we can provide rich explanations for why a non-zero prediction probability is assigned to the class other than that of the predicted class. Figures SF3, SF4, SF5, SF6, SF7, SF8, SF9, SF10 show few such explanations.

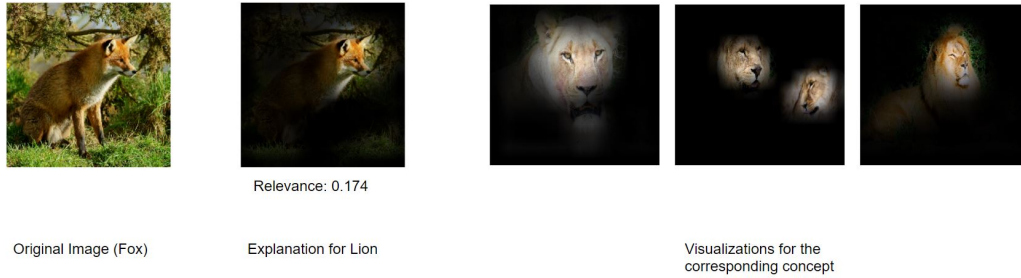


Figure SF3: [Best viewed in color] The model is looking at the face region in the fox image in this concept. The visualizations for lion for this concept is also face.

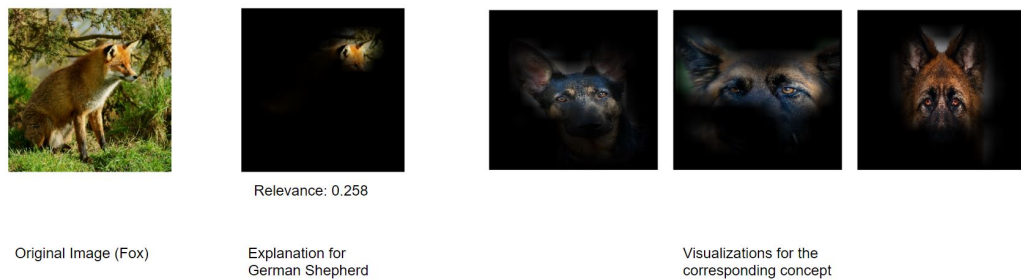


Figure SF4: [Best viewed in color] The model is looking at the ear region in the fox image in this concept. The visualizations for dog for this concept is the ear and eye region. Since eyes are not clearly visible in the fox image, the model is looking at just the ears in this concept.

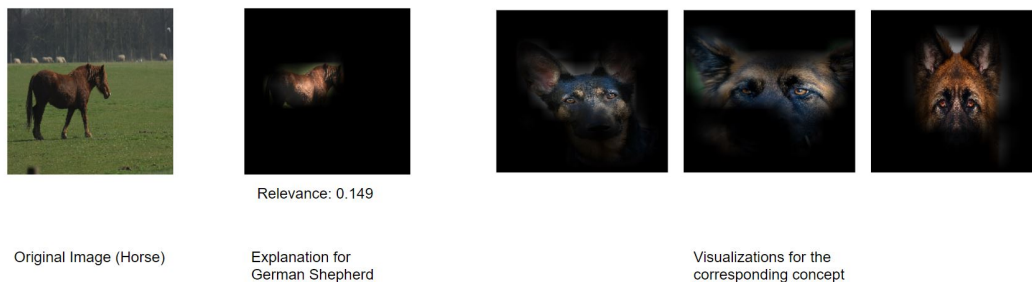


Figure SF5: [Best viewed in color] The bulges on back and front side of horse make the region highlighted in the explanation for german shepherd look like ears (based on color). Since the concept itself is also of ear and eye region, this concept contributes towards german shepherd probability in this image.

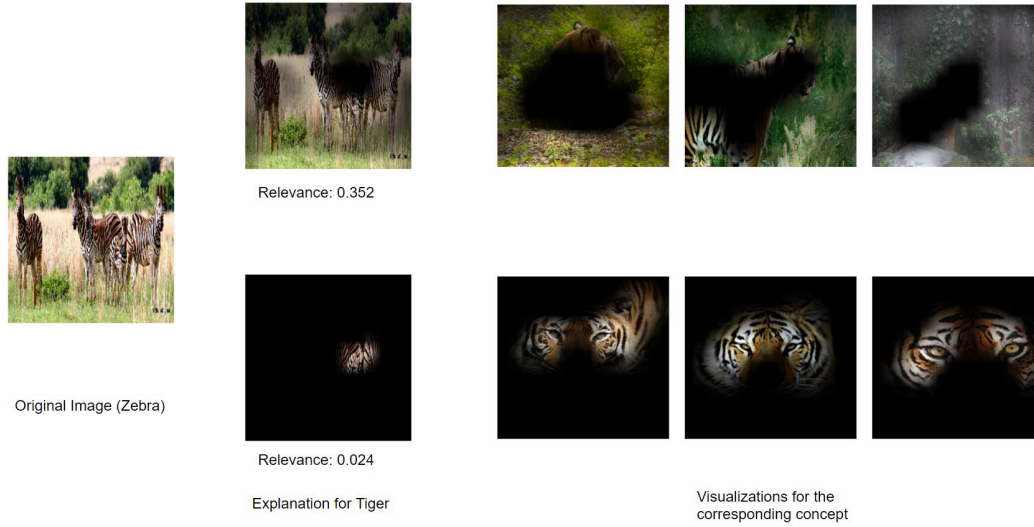


Figure SF6: [Best viewed in color] In the concept in the first row, the model is looking at the background. Background of Zebra image and that of the images for the actual concept is quite similar. In the concept in second row, the model is looking at the stripes region in the zebra image, and the concept is also the stripes on tiger.

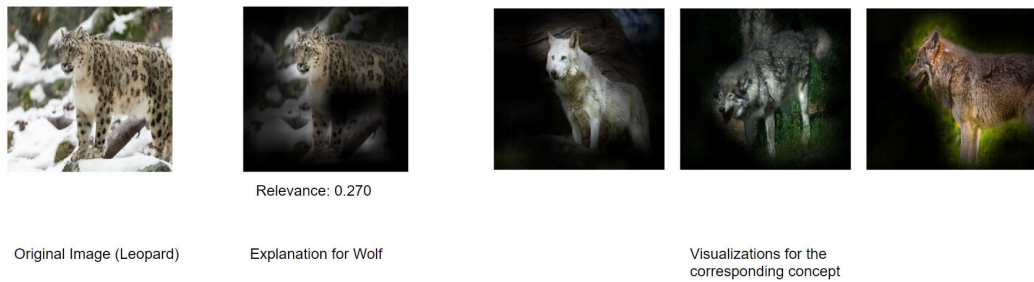


Figure SF7: [Best viewed in color] The model is looking at the body of leopard in the explanation for wolf class. The concept itself also represents the body of wolf.

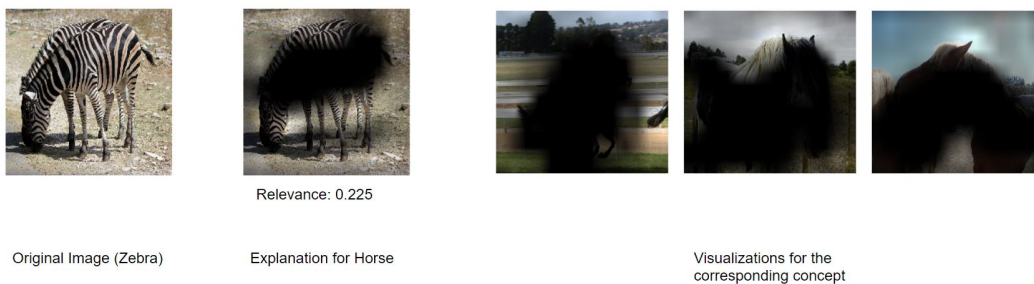


Figure SF8: [Best viewed in color] The model is looking at the background. Background of Zebra image and that of the images for the actual concept is quite similar



Figure SF9: [Best viewed in color] The model is looking at the background. Background of Tiger image and that of the images for the actual concept is quite similar

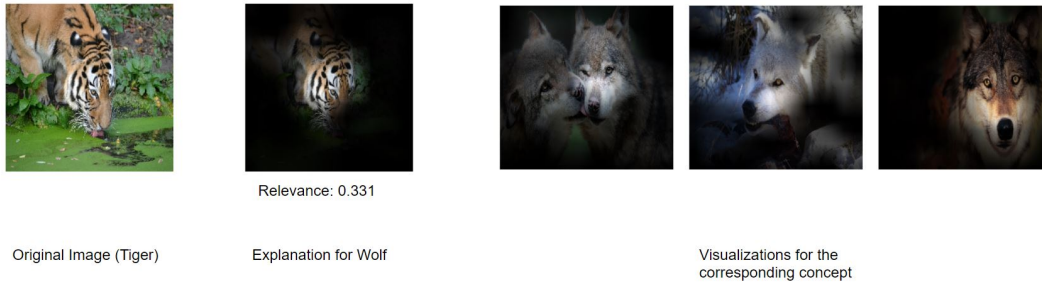


Figure SF10: [Best viewed in color] The model is looking at the face of tiger in the explanation for wolf class. The concept itself also represents the face region of wolf

S5 Stability

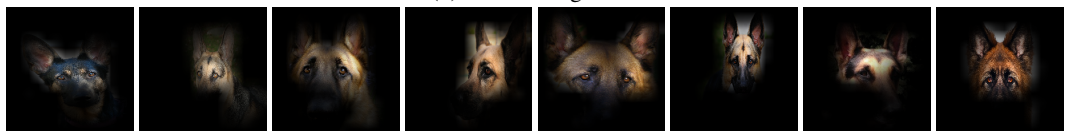
Stability of an interpretation refers to visual comparison of the learned concepts for images of the same class with varying orientation, size or position. The explanations of such images should be similar as well i.e., the learned concepts should be consistent. Given the saliency maps of a few such images, we should be able to understand the underlying concept easily. Figures SF11 and SF12 show that our learned concepts have very high stability. We also demonstrate that the concept embeddings for a particular concept for images with varying orientation, size or position are very close to each other. We take a small set of 10 images from the fox class and randomly choose 5 concepts of this class. For each concept we calculate the pairwise Euclidean distance between the concept embedding of the images. The results are shown in Figure SF13 in the form of a box diagonal matrix. Let c_i denote the concept and f_i denote the image then the axis of the matrix are defined as $[(c_1, f_1), \dots, (c_1, f_{10}), (c_2, f_1), \dots, (c_2, f_{10}), \dots, (c_5, f_1), \dots, (c_5, f_{10})]$. As it can be seen that the distances along the diagonal is lesser compared to the non-diagonal distances. This shows that the embeddings of similar (dissimilar) concepts are closer (farther).



(a) Fox - Ear



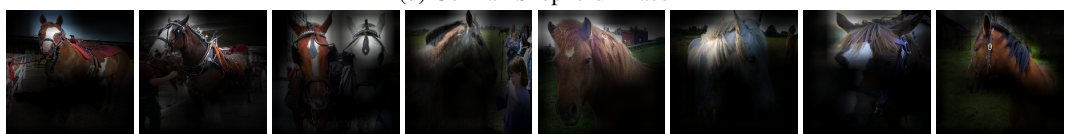
(b) Fox - Background



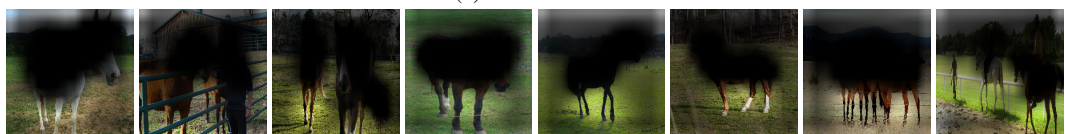
(c) German Shepherd - Eyes



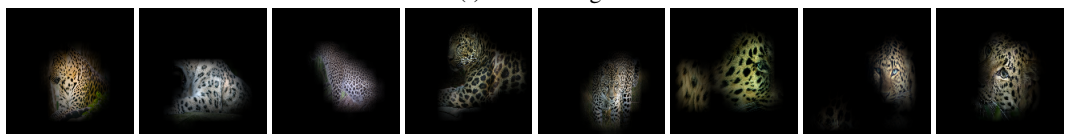
(d) German Shepherd - Face



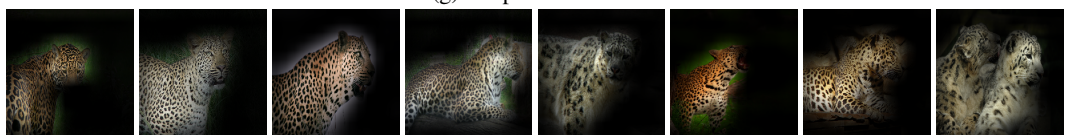
(e) Horse - Crest



(f) Horse - Legs



(g) Leopard - Texture



(h) Leopard - Body



(i) Lion - Body and Background



(j) Lion - Face

Figure SF11: [Best viewed in color] Visualization of concepts across multiple images of the same class

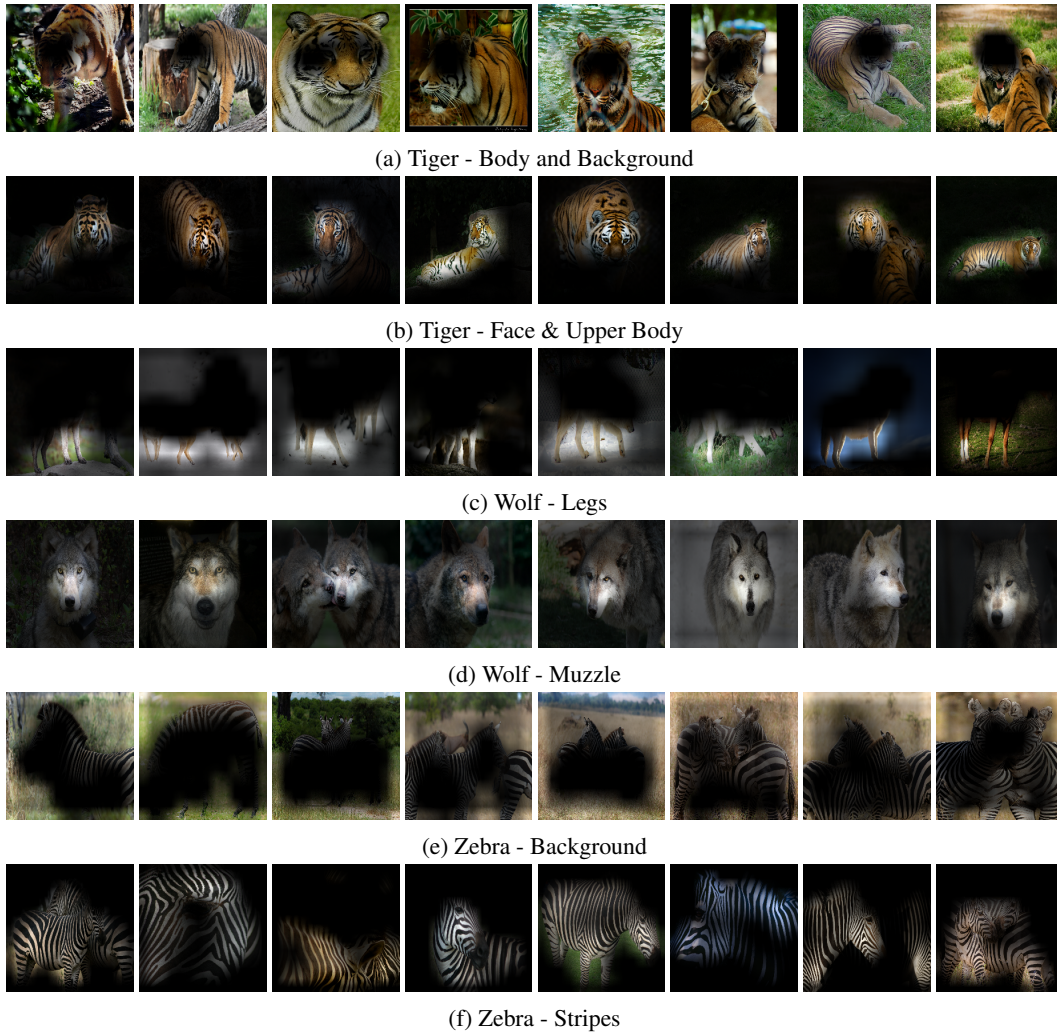


Figure SF12: [Best viewed in color] Visualization of concepts across multiple images of the same class

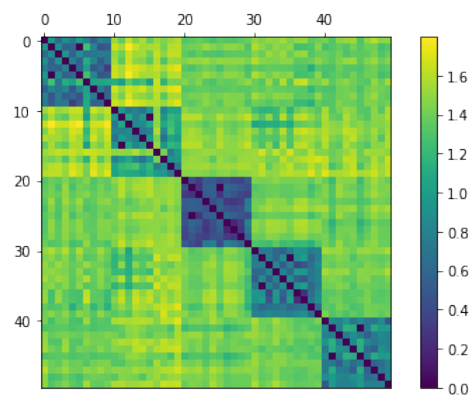


Figure SF13: [Best viewed in color] Pairwise Euclidean distance between concept embeddings

S6 Comparison of robustness

Robustness is an important desideratum of an explanation. The explanation generated by any interpretability method should be robust to local perturbations in the input image. Figure SF14 shows that this is not the case for popular interpretability methods; even adding minimal noise to the input introduces visible changes in the explanations. We formally quantify this parameter of evaluation as the intersection over union between the explanations of the image and its perturbed variation:

$$R(x_i) = IoU(f_{expl}(x_i), f_{expl}(\tilde{x}_i)) \quad (1)$$

Here x_i is an input image and \tilde{x}_i is the perturbed image. f_{expl} refers to the explanation function whose output is the saliency map for the interpretation method. The saliency maps are thresholded to create a binary map before calculating the IoU. We perform robustness experiment on various parameters including variations in brightness, contrast, random noise and rotation. For each parameter we define a range of values that describes the intensity of perturbation, eg. standard deviation of the noise distribution to be added in an image, delta value to increase the brightness/contrast in the image, angle of rotation etc. We define a wide range of threshold values in the range $[0.3 - 0.7]$ to generate the binary maps and compare their robustness for three approaches including MACE, GradCAM Selvaraju et al. [2017] and GradCAM++ Chattopadhyay et al. [2018].

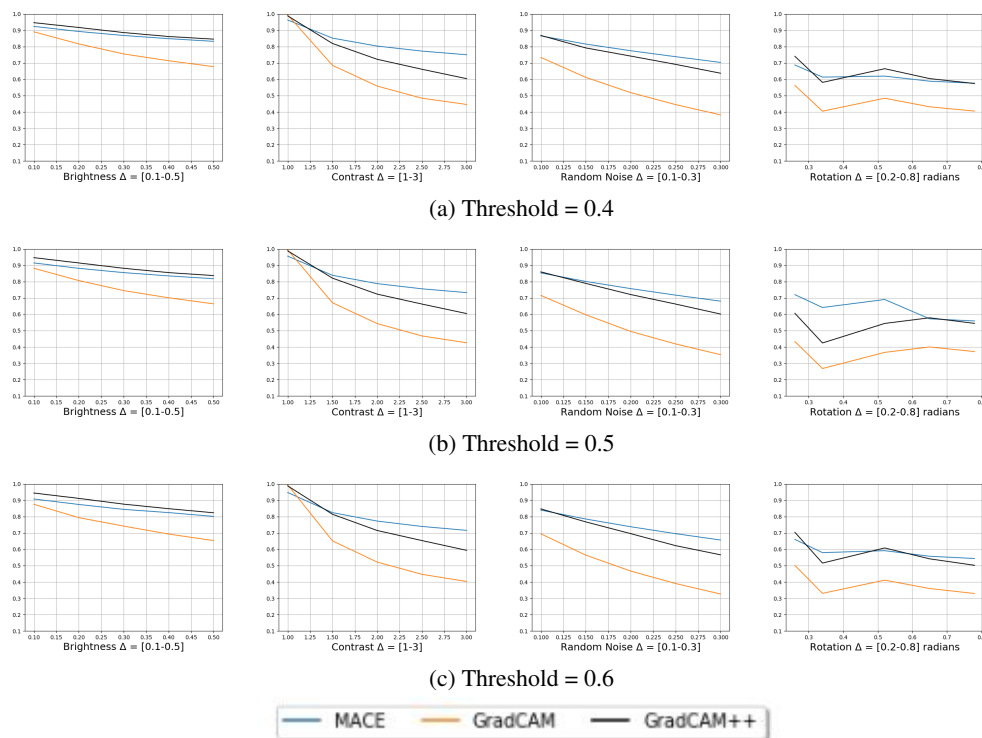


Figure SF14: [Best viewed in color] Intersection Over Union for thresholded saliency maps of original and perturbed image

S7 Places Experiments

The visualizations for the VGG model, trained on Places365 are shown in Figure SF15. We observed that most of the concepts were pruned and we had just 2-3 useful concepts per class. One of the reasons could be that the classes in the dataset were different from each other and didn't require a large set of concepts to make the prediction.

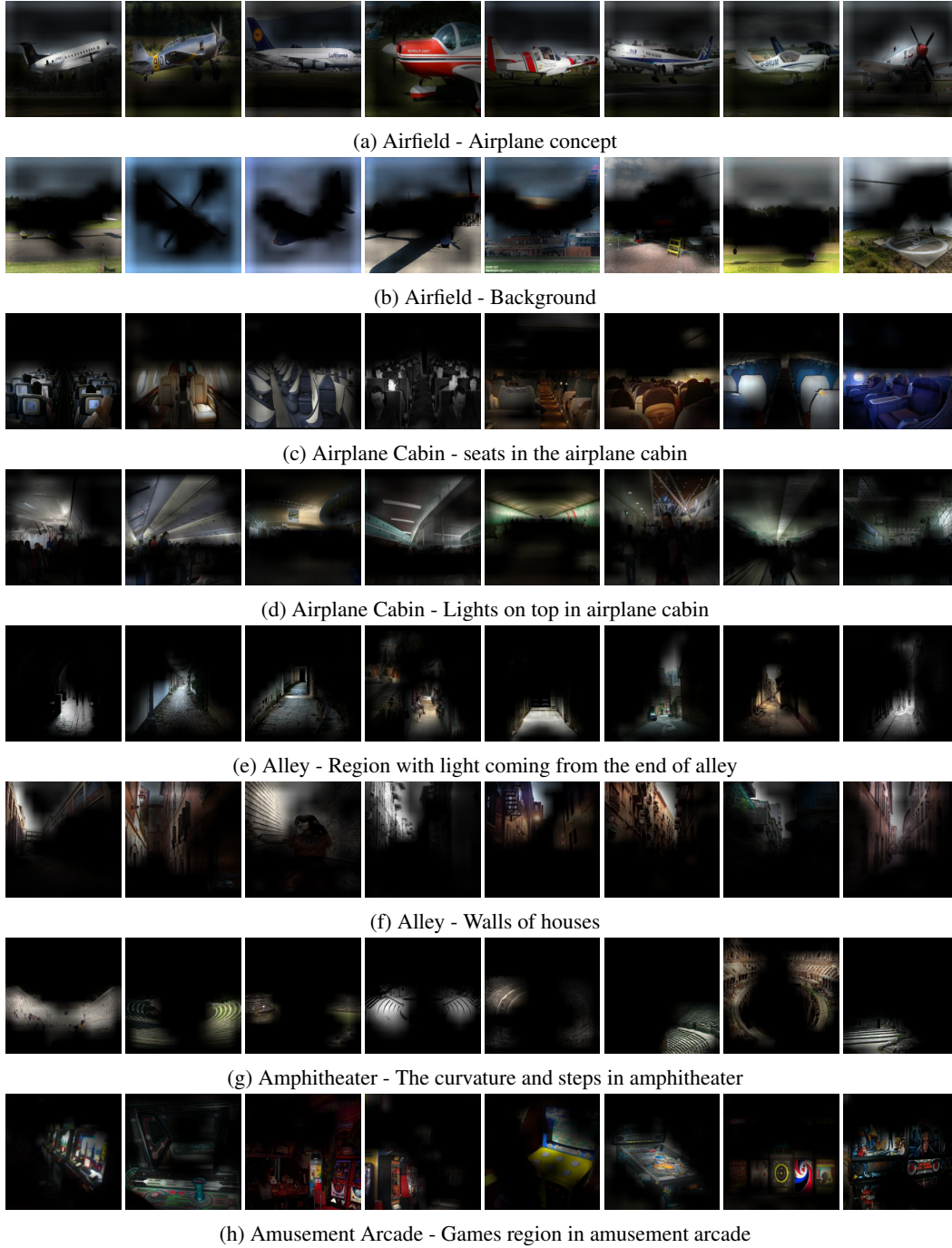
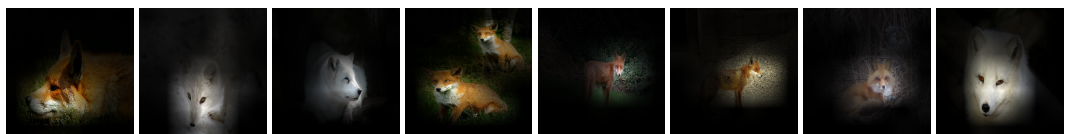


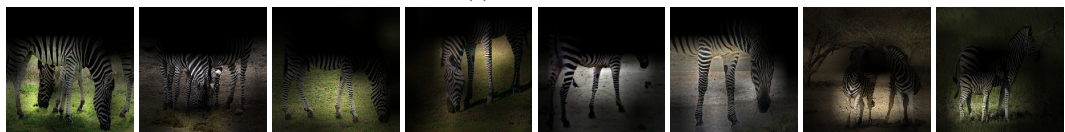
Figure SF15: [Best viewed in color] Concepts for Places365 Dataset.

S8 ResNet Experiments

The visualizations for the ResNet model, trained on AWA2 are shown in Figure SF16. We observed more head and body concepts in the case of the ResNet model, as compared to the VGG model. A possible reason for this could be that ResNet is a much deeper network than the VGG, and hence the features present in the activation maps of the last convolution layer majorly contain high-level body concepts.



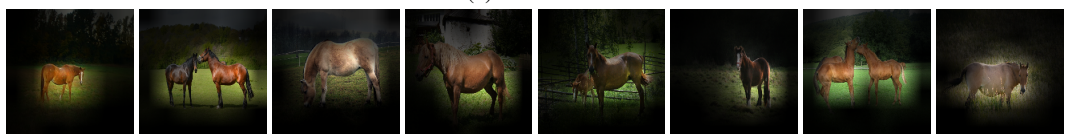
(a) Fox - Head



(b) Zebra-Legs



(c) Lion-Head



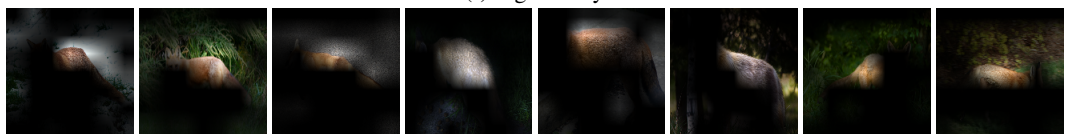
(d) Horse-Body



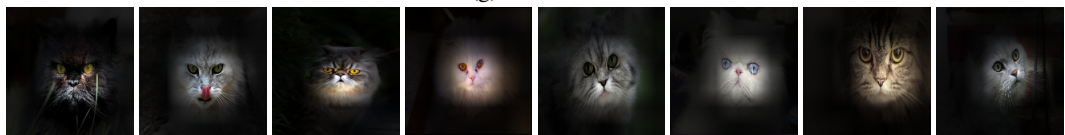
(e) Wolf-Head



(f) Tiger-Body



(g) Fox-Back



(h) Cat-Face

Figure SF16: ResNet Outputs

S9 Ablation Study

In order to justify the need for both the \mathcal{L}^O and \mathcal{L}^D losses, we compared the faithfulness of the approach with and without these losses. According to our hypothesis, we require the \mathcal{L}^O loss to help us recreate the output of the first dense layer of a classifier, so as to keep our concept embeddings faithful to the classifier. We require \mathcal{L}^D Loss to avoid possibility of inconsistency with the final output of the classifier.

We train two models, one without using the \mathcal{L}^O loss, while other without the \mathcal{L}^D loss. We train both the models for 50 epochs, setting the learning rate to be 10^{-4} .

In Figure SF17, we see that there is a significant difference in drop in the true class probability when either of the two losses are reduced, showing both these losses are necessary for the faithfulness of our MACE unit.

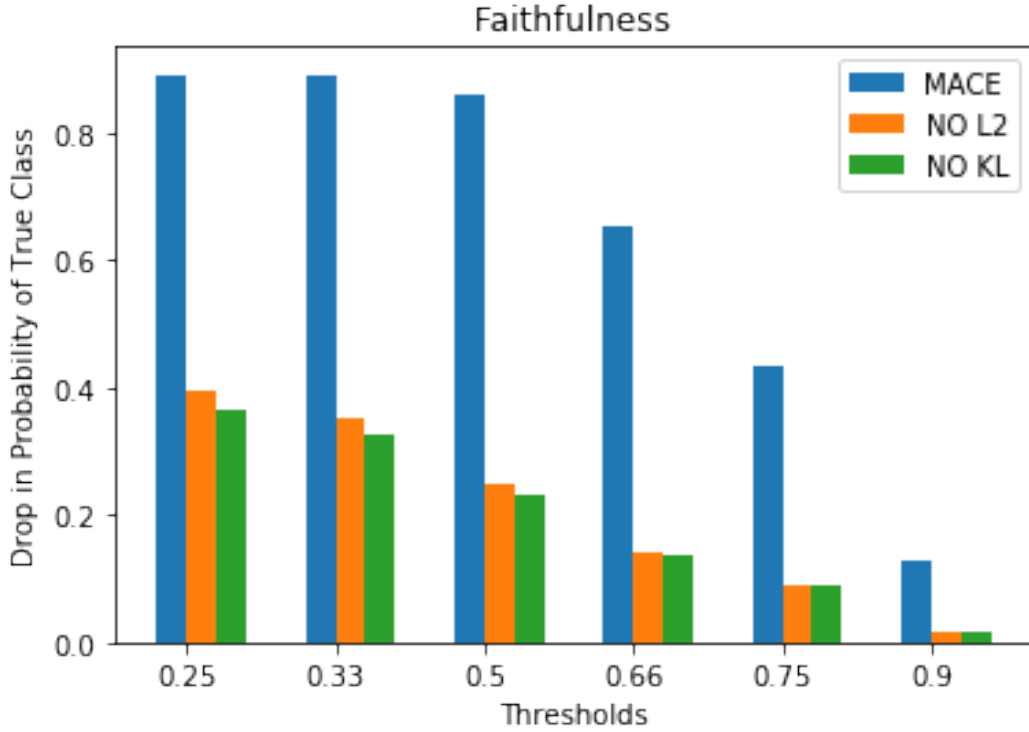


Figure SF17: [Best viewed in color] Effect of \mathcal{L}^O and \mathcal{L}^D losses on model faithfulness

References

- Animals with attributes 2 (awa2) dataset. URL <https://cvml.ist.ac.at/AwA2/>.
- Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 839–847. IEEE, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. URL <http://arxiv.org/abs/1512.03385>.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2014.
- Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2251–2265, 2018.
- Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.