

# Iterative and Adaptive Sampling with Spatial Attention for Black-Box Model Explanations

Bhavan Vasu\*  
Kitware Inc.  
Clifton Park, NY, USA 12065  
bhavan.vasu@kitware.com

Chengjiang Long\*  
Kitware Inc.  
Clifton Park, NY, USA 12065  
chengjiang.long@kitware.com

## Abstract

Deep neural networks have achieved great success in many real-world applications, yet it remains unclear and difficult to explain their decision-making process to an end-user. In this paper, we address the explainable AI problem for deep neural networks with our proposed framework, named IASSA, which generates an importance map indicating how salient each pixel is for the models prediction with an iterative and adaptive sampling module. We employ an affinity matrix calculated on multi-level deep learning features to explore long-range pixel-to-pixel correlation, which can shift the saliency values guided by our long-range and parameter-free spatial attention module. Extensive experiments on the MS-COCO dataset show that the proposed approach matches or exceeds the performance of state-of-the-art black-box explanation methods. Our source code is available at <https://github.com/vbhavank/IASSA-Saliency>.

## 1. Introduction

It is still unclear how a specific deep neural network works, how certain it is about the decision making, *etc.*, although deep networks have achieved remarkable success in multiple applications such as object recognition [42, 51, 9, 18, 16, 19, 17, 10, 38], object detection [21, 5, 30], image labeling [15, 8], media forensics [33, 20, 14], medical diagnosis [43, 44], and autonomous driving [23, 4, 22]. However, due to the importance of explanation in understanding and building trust in cognitive psychology and philosophy [12, 13, 28, 45, 31], it is very critical to make the deep neural networks more explainable and trustable, especially to ensure that the decision-making mechanism is transparent and easily interpretable. Therefore, the problem of Explainable AI, *i.e.*, providing explanations for an intelligent models decision, especially in explaining classification de-

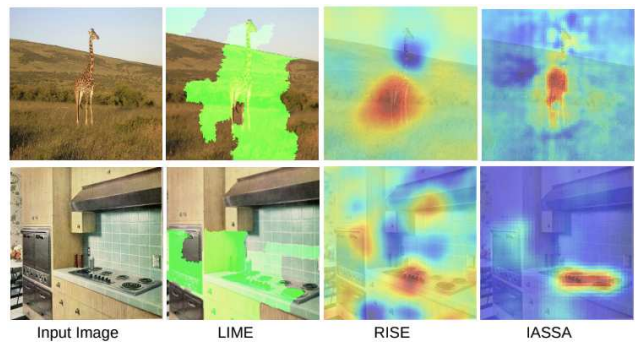


Figure 1. Visual comparison between the proposed method IASSA and two state-of-the-art black-box explanation algorithms, *i.e.*, LIME [26] and RISE [27], for the importance of producing explanations.

cisions made by deep neural networks on natural images, attracts much attention in artificial intelligence research [34].

Rather than explainable solutions [35, 50, 40, 29, 37] to certain white-box models via calculating importance based on the information like the network’s weights and gradients. We advocate a more general explainable approach to produce a saliency map for an arbitrary network that is treated as a black-box model, without requiring details about the architecture and implementation. Such a saliency map can show how important each image pixel is for the networks prediction.

Recently, multiple explainable approaches have been proposed for black-box models. LIME [26, 1] proposes to draw random samples around the instance for an explanation by fitting an approximate linear decision model. However, such a superpixel based saliency method may not group correct regions. RISE [27] explores the black-box model by sub-sampling the input image via random masks and generates the final importance map by a linear combination of the random binary masks. Although this is seemingly simple yet surprisingly powerful approach for black-box models, the results are still far from perfect, especially

\*Equal contributions. This work was supervised by Chengjiang Long.

in complex scenes.

In this paper, inspired by RISE [27], we propose a novel iterative and adaptive sampling with spatial attention (IASSA) for explaining black-box models. We do not access weights or gradients of the underlying model. We only sample the input image randomly using a sliding window during the initialization stage. And then an iterative and adaptive sampling module is designed to generate sampling masks for the next iteration, based on the adjusted attention map which is obtained with the saliency map at the current iteration and the long-range and parameter-free spatial attention. Such an iterative procedure continues until convergence *i.e.* no substantial change in the final saliency map. The visual comparison of explanations obtained with LIME and RISE is shown in Figure 1.

Regarding the long-range and parameter-free spatial attention module, we apply a pre-trained model trained on the large-scale ImageNet dataset to extract features for the input image. Note that we combine multi-level contextual features to better represent the image. Then we calculate an affinity matrix and apply a softmax function to get spatial attention. Since the affinity matrix covers pixel-to-pixel correlations irrespective of them being local neighbors or not, our attention covers long-range inter-dependencies. Also, no parameters are required to be learned in this procedure. Such a long-range and parameter-free spatial attention can guide the saliency values in the obtained saliency map towards highly correlated pixels. This can be very helpful to identify sampling regions for adaptive sampling in the next iteration.

Another contribution of our work is to evaluate for "goodness" of an explanation. Besides previously used metrics like *deletion*, *insertion* and "*Pointing Game*" [27], we also choose to use  $F-1$  and  $IoU$ . We evaluate the final saliency maps at both the image and pixel-level to highlight the success of our approach in maximizing information contained in each pixel. We argue that such a comprehensive evaluation should be more trustable when compared with the human-annotated importance. In our case, we assume ground truth masks as representative of human interpretation of an object, as they are human-annotated.

To sum up, the technical contributions are of three-folds: (1) we propose an iterative and adaptive sampling for generating accurate explanations, based on the adjusted saliency map generated by combining the saliency map obtained from the previous iteration and the long-range and parameter-free spatial attention map; (2) our long-range and parameter-free attention module that incorporates "objectness" and guides our adaptive sampler with the help of multi-level feature fusion; and (3) we further introduce an evaluation scheme that tries to estimate goodness of an explanation in a way that it is reliable and accurate.

We conduct extensive experiments on the popular and

vast MS-COCO dataset [11] and compare it with the state-of-the-art methods in the field. The experimental results demonstrate the efficacy of our proposed method.

## 2. Related work

The related work can be divided into two categories, *i.e.*, *white-box approaches* and *black-box approaches* for producing explanations.

**White-box approaches** rely on information such as the model parameters or gradients, as well as the intermediate feature maps. Zeiler *et. al.* [47] visualize the intermediate representation learned by CNNs using deconvolutional networks. Explanations are achieved in other methods [25, 36, 46] by synthesizing an input image that highly activates a neuron. Class activation maps (CAM) [52] achieve class-specific importance at each location in an image by computing a weighted sum of the activation values at each location across all channels using a Global Average Pooling layer (GAP). Such a method prevents us from using this approach to explain models lacking a native GAP layer without additional re-training. Later, CAM was extended to Grad-CAM [35] by weighting the feature activation values at every location with the average gradient of the class score (w.r.t. the feature activation values) for every feature map channel. In addition, Zhang *et. al.* [50] introduced a probabilistic winner-takes-all strategy to compute the relative importance of neurons towards model predictions. Fong *et. al.* [7] and Cao *et. al.* [2] learn a perturbation mask that maximally affects the models output by back-propagating the error signals through the model. However, all of the above methods assume that the internal parameters of the underlying model are accessible as a white-box. They achieve interpretability by incorporating changes to a white-box based model and are constrained to specific network architectures, limiting reproducibility on a new dataset.

**Black-box approaches** treat the learning models as purely black-box, without requiring access to any details of the architecture and the implementation. LIME [32] tries to fit an approximate linear decision model (LIME) in the vicinity of a particular input. For a sufficiently complex model, a linear approximation may not result in a faithful representation of the non-linear model. Even though LIME model produces good quality results on the MS-COCO dataset, due to its reliance on super-pixels, it is not the best at grouping object boundaries with activation. As an improvement over LIME, RISE model [27] was proposed to generate an importance map indicating how salient each pixel is for the black-box model's prediction. Such a method estimates importance empirically by probing the model with randomly masked versions of the input image and obtaining the corresponding output probability. Note that sampling methods to generate explanations have been

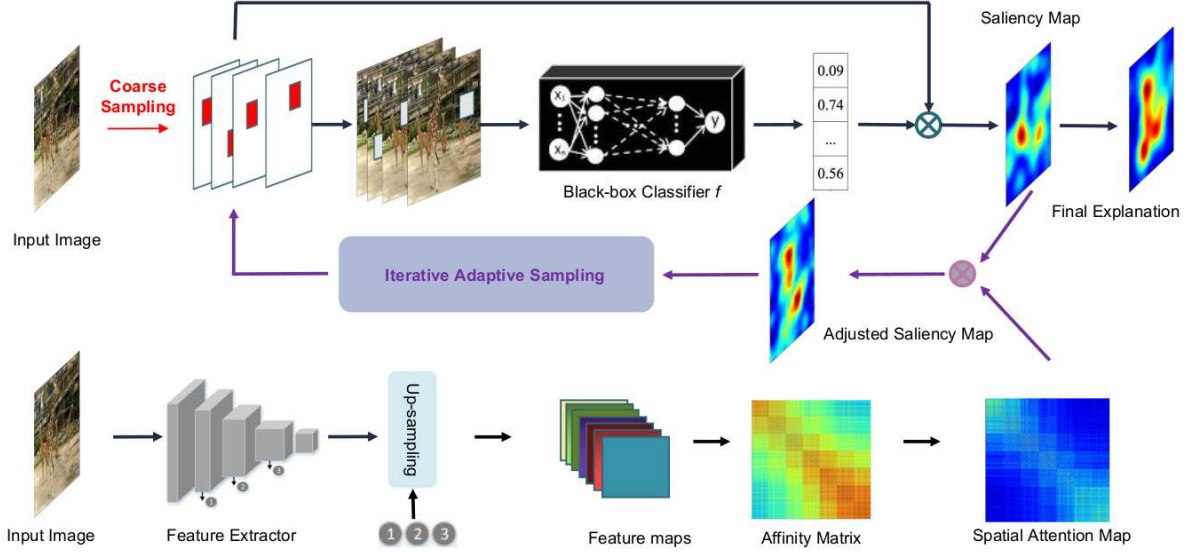


Figure 2. The framework of our unsupervised saliency map extraction method for an explanation of a black-box model. Given an input image, we perform a rough pass over the image to start the iterative process with large window size. The masked images are passed to the black box classifier that predicts logit scores for each sample, the predicted logit scores are used to weight image regions to produce a saliency map. Then an adjusted saliency map is generated by combining with the long-range and parameter-free spatial attention module to guide the iterative and adaptive sampling module to generate the new sampling masks, which leads to a new saliency map. Such iterative procedure continues until convergence. Note that the spatial attention module is built based on multi-level deep learning features via an affinity matrix.

explored in the past [27, 48, 6]. Even though they produce explanations for a wide variety of black-box model applications, their resolution is always limited by factors like sampling sensitivity and strength of classifier.

In this paper, unlike the existing methods, we explore a novel method to provide precise explanations for any application that uses a deep neural network for feature extraction, irrespective of the multi-level features. We leverage a long-range and parameter-free spatial attention to adjust the saliency map. We propose an iterative and adaptive sampling module with long-range and parameter-free attention to determine important regions in an image. The proposed system can also be adapted to perform co-saliency [6] by weighting the final saliency map using a standard feature comparison metric like Euclidean or Cosine distance. This makes our approach robust to the form of explanation desired and produces better quality saliency maps across different applications with little or no overhead in training.

### 3. Methodology

The proposed framework is illustrated in Figure 2. Given an input image, we perform a rough pass to initialize our approach. The sampled image regions are passed to the black box classifier that predicts logit scores for each sample, the predicted logit scores are used to weight image regions to produce an aggregated response map. The response map is combined with the attention map obtained from the long-

range and parameter-free spatial attention module (LRPF-SA) to get an adjusted saliency map. The attention module also guides the iterative and adaptive sampling to sample relevant regions in the next iteration. Such an iterative procedure continues until convergence. Note that the spatial attention module is built on multi-level deep learning features via an affinity matrix. In the following subsections, we further explain our approaches in detail.

#### 3.1. Iterative and Adaptive Sampling Module

We propose a novel iterative and adapting sampler that is guided by our LRPF-SA to automatically pick sampling regions of interest with an appropriate sampling factor rather than weighting all image regions equally. Sampling around important regions ensures faster convergence and finer saliency maps. The iterative quality of our approach also allows the users to control the quality of saliency maps, which is inversely proportional to the amount of time needed to generate them. We believe this is crucial in real-time applications where the same explanation generator system needs to be scaled according to user requirements with minimal changes.

Given an image  $I$ , a black-box model  $f$  produces a score vector of length  $c$ , where  $c$  is the number of classes the black-box model was trained for. We sample the input image  $I$ , using masks  $M: \Lambda \rightarrow \{0, 1\}$  be a sliding window of size  $W$  and stride  $S$ . Considering the masked version

$(I \odot M)$  of  $I$ , where  $\odot$  represents element-wise multiplication, we compute the confidence scores for all the masked images  $f(I \odot M)$ . We define the importance of a pixel  $\lambda \in \Lambda$  as the expected score over all possible masks  $M$  conditioned on the event that pixel  $\lambda$  is observed. In other words, when the scalar score  $f(I \odot M)$  is high for a chosen mask  $m \in M$ , we can infer that the pixels preserved by  $m$  are important. We define the importance of the pixel  $\lambda$  as the expected score over all possible masks conditional on the event that  $\lambda$  is observed, *i.e.*.

$$S(I, f, \lambda) = \sum_m f(I \odot M) P[M = m, M(\lambda) = 1], \quad (1)$$

where

$$P[M = m, M(\lambda) = 1] = \begin{cases} 0, & \text{if } m(\lambda) = 0 \\ P[M = m], & \text{if } m(\lambda) = 1 \end{cases} \quad (2)$$

With Equation 1 and 2, we arrive at

$$S(I, f, \lambda) = \frac{1}{P[M(\lambda) = 1]} \sum_m f(I \odot M) \cdot m(\lambda) \cdot P[M = m] \quad (3)$$

Considering that  $P[M(\lambda) = 1] = \mathbb{E}[M(\lambda)]$ , we rewrite Equation 3 in matrix notation as

$$S(I, f, \lambda) = \frac{1}{\mathbb{E}[M]} \sum_m f(I \odot M) \cdot m \cdot P[M = m] \quad (4)$$

Using Monte Carlo sampling, at the iteration 0, the final saliency map is computed as a weighted average of a collection of masks  $\mathbf{M}_k = \{M_1, \dots, M_N\}$  by the following approximation:

$$S(I, f, \lambda) \approx \frac{1}{\mathbb{E}[M] \cdot N} \sum_{i=1}^N f(I \odot M_i) \cdot M_i(\lambda). \quad (5)$$

When the black-box model  $f$  is associated with a class  $c$ , then we can obtain a saliency map corresponding to  $c$  according to Equation 4. Although most applications require only the top-1 saliency map, our approach can be used to obtain class specific salient structures.

The initial saliency map  $S_0$  is generated based on a sliding window  $\mathbf{M}_0$ . After the initialization, we take the long-range and parameter-free attention module  $A$  to adjust the saliency map from  $S_k$  to  $S'_k$  at the  $k$ -th iteration by the following rules

$$S'_k = \beta S_k + (\beta - 1) A \times S_k, \quad (6)$$

where  $\beta$  is a regularizer to control the amount of influence the attention network has towards generating the final explanation. The intuition behind using both saliency and attention maps is that, while the saliency map  $S_k$  is associated

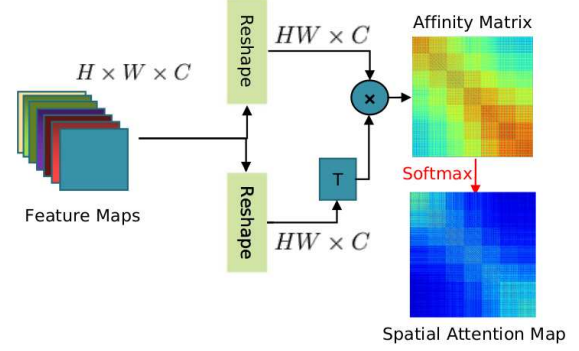


Figure 3. An illustration of our LRPF-SA module that produces attention maps used to guide the iterative and adaptive sampling module.

with the output of a back-box model, we provide a new insight with our proposed LRPF-SA (see next subsection) to apply some spatial constraints with respect to the extracted features. Therefore, by combining both forms of explanations we hope to converge on an aggregated saliency map that gives a complete picture of the image regions that interest the system while conforming to object boundaries.

Then we use  $S'_k$  to guide the adaptive sampling for the next iteration by

$$\mathbf{M}_{k+1} = \text{HAR}(S'_k), \quad (7)$$

where  $\text{HAR}(\cdot)$  denotes the highest activated region obtained by applying a threshold, resulting in a binary map that highlights all pixels containing the object of interest, *i.e.*,

$$\text{HAR}(S'_k) = S'_k > T_{\text{thresh}} \quad (8)$$

With the adaptive sampling masks  $\mathbf{M}_{k+1}$ , we are able to apply Equation 5 to obtain the saliency map  $S_{k+1}$  at the  $(k+1)$ -th iteration. And then  $S'_{k+1}$  is obtained using Equation 6 to get the adaptive sampling masks  $\mathbf{M}_{k+2}$  for generating the saliency map  $S_{k+2}$  at the  $(k+2)$ -th iteration. It is worth noting that the window size and stride can be gradually depreciated with respect to the number of iterations to increase the resolution of saliency maps until there is very little or no change in the quality of maps. The number of iterations can also be fixed based on application requirements, when the user is willing to sacrifice the quality of saliency maps for run-time.

### 3.2. Long-Range and Parameter-Free Spatial Attention

Obtaining an attention map from a deep learning model is a well-researched topic [39, 49]. Recent developments in minimizing attention generation overhead was proposed in [41]. Inspired by [41], we propose a novel long-range and parameter-free spatial attention (LRPF-SA) module. We



make use of a deep network for feature extraction that encompasses activations from different levels of the network. We believe by using activations from different levels of the network we provide a true explanation about how the image is perceived by the complete network, giving rise to hierarchical salient concepts in the attention map. The saliency maps are then used to choose from the hierarchical concepts that match with image boundaries, thus giving rise to accurate and reliable saliency maps.

In this paper, we use the pre-trained network learned on the ImageNet dataset. Note that in the case of a new domain, the network can be adapted into the target domain using methods proposed in [3]. Let  $\Phi(I)$  be a pre-trained deep network used to extract multi-level features that are combined by upsampling and fusion. Finally, we use a softmax operation over the resulting affinity matrix to obtain an attention map as showing in Figure 3.

Note that the Affinity matrix contains dependencies of every pixel with all other pixels. Let  $\Phi_1(I)$ ,  $\Phi_2(I)$ ,  $\Phi_3(I)$  and  $\Phi_4(I)$  be the features extracted from four different levels of the feature extractor. Since we use a  $\Phi_1(I)$  of  $H \times W \times C_1$  dimensions, where  $H$  and  $W$  are the height and width of the obtained feature maps and  $C_1$  is the number of channels. The feature maps  $\Phi_2(I)$ ,  $\Phi_3(I)$  and  $\Phi_4(I)$  are upsampled to  $H \times W$ , with channel numbers  $C_2$ ,  $C_3$ , and  $C_4$ . Upsampling the feature maps let us directly compute an aggregated response using the following Equation

$$\Phi(I) = \Phi_1(I) \oplus (\Phi_2(I)_{\uparrow} \oplus \Phi_3(I)_{\uparrow} \oplus \Phi_4(I)_{\uparrow}), \quad (9)$$

where the subscript  $\uparrow$  denotes an upsampling operation and  $\oplus$  is an concatenation operation, and the long-range and parameter-free spatial attention can be obtained by

$$A = \text{softmax}(\Phi'(I) \odot (\Phi'(I))^T), \quad (10)$$

where  $\Phi'$  is reshaped from  $\Phi$  of dimension  $H \times W \times C$  to  $HW \times C$ , and  $C = \sum_{i=1}^4 C_i$  is the number of channels in  $\Phi$ . Figure 3 shows an illustration of our LRPF-SA module that produces attention maps used to guide the iterative and adaptive sampling module. By using an attention mechanism we hope to gain information related to the "objectness", hidden among pixels in an image.

### 3.3. Iterative Saliency Convergence

We propose to find the best possible saliency map that captures the decision-making process of the underlying algorithm in an iterative manner. Generating high-quality explanations is a very time-consuming process and limits its usage in real-time applications that require generating precise maps on large datasets. By gradually converging on the optimal saliency map, we hope to let the user decide the rate

of convergence that fits their time budget, opening up possibilities to use explanations in a wide variety of applications.

## 4. Experiments

One would wonder if we should consider an explanation "good" if it represents the importance according to the black-box classifier or if it conforms to object boundaries, encouraging human trust in the explanation system. To verify the effectiveness of our proposed approach IASSA, we conduct experiments on the MS-COCO dataset [11] and evaluate explanations for their ability to best represent image regions that the underlying model relies on and also for their segmentation performance. By leveraging attention with model dependant saliency, the proposed approach achieves better performance when evaluated for insertion, deletion, intersection over union (IoU), F1-score, and a pointing game score [27]. We believe we can leverage the proposed explanation generation method to fine-tune models, especially deep learning classifiers in a closed loop using Attention Branch Networks [24].

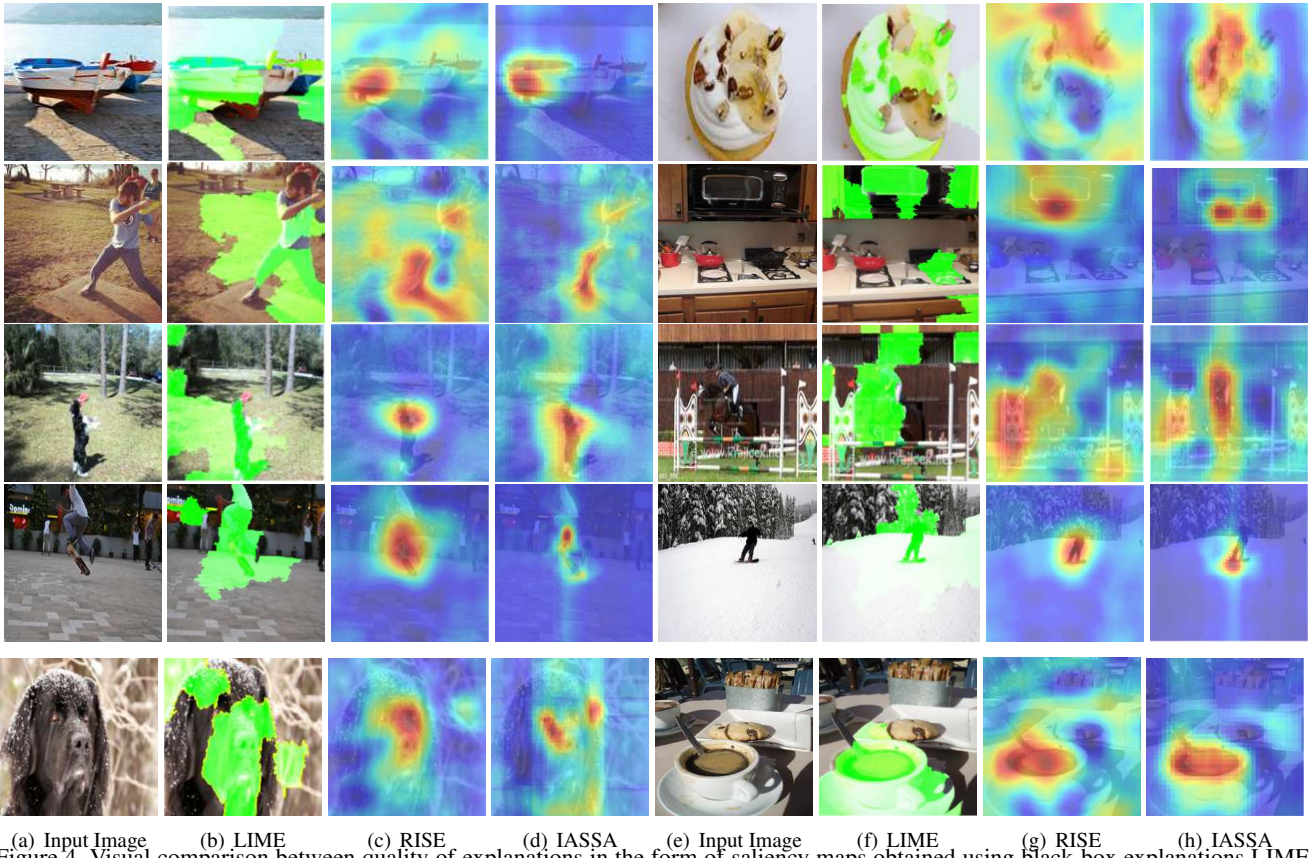
Note that in this paper, the input images are resized to  $224 \times 224$  to facilitate mask reuse and ease in feature extraction. The IASSA module is initialized with a window size of  $W$  of 45 and a stride  $S$  of 8 with step size 1.5 and 0.2 respectively. We use a  $\beta$  of 0.5 and a  $T_{thresh}$  of 0.3 to generate a new saliency map at the  $k$ -th iteration. We present results for 25 iterations.

### 4.1. Evaluation Metrics

Evaluating the quality of saliency maps can be subjective to the kind of explanation. We evaluate the quality of saliency maps using five different metrics: deletion, insertion, IoU, F1-score, along with a pointing game score [27].

In deletion, given a saliency map and input image  $I$  we gradually remove pixels based on their importance in the saliency map, meanwhile monitoring the Area Under the Curve (AUC). A sharp drop in activation as a function of the fraction of pixels removed can be used to quantify the quality of saliency maps. Analogously, in insertion, we reveal pixels gradually in the blurred image. The pixels can be removed or added in several ways like setting the pixels of interest to zero, image mean, gray value or blurring pixels. For deletion, we set pixels of interest to a constant grey value. But the same evaluation protocol cannot be used for insertion as the model would be biased towards grouping of pixels introduced on an empty canvas.

To prevent the introduction of bias towards pixels grouping shapes, for insertion we unblur regions of the image, under consideration. The IoU and F1-score are calculated by applying a threshold  $T_{thresh}$  of 0.3 on the aggregated saliency maps using Equation 8 and 9 obtained at the end of  $i$ -th iterations. We also use a pointing game that considers



(a) Input Image (b) LIME (c) RISE (d) IASSA (e) Input Image (f) LIME (g) RISE (h) IASSA  
Figure 4. Visual comparison between quality of explanations in the form of saliency maps obtained using black-box explanations LIME, RISE and our proposed IASSA on the MS-COCO dataset.

an explanation as a positive hit when the highest activated pixel lies on the object of interest. We average all performance metrics at both image and pixel-level by normalizing the performance by the number of pixels activated. The normalization for per-pixel performance lets us fairly evaluate explanations that might cover a region much larger than the object of interest but also include the object.

#### 4.2. Effectiveness of Iterative Adaptive Sampling Module with LRPF-SA

We consider explanation generation as an optimization problem, assuming there exists an optimal explanation that encapsulates both model dependence and human interpretable cues in an image. Converging on this optimal explanation is conditioned upon parameters such as the number of iterations  $k$ , regularizer  $\beta$ , and threshold  $T_{thresh}$  (where  $\beta$  and  $T_{thresh}$  decide the convergence rate). We fix the value for  $\beta$  and  $T_{thresh}$ , and evaluate the impact of  $k$ .

A qualitative analysis of the proposed explanation system’s ability to converge on an optimal explanation can be visualized in Figure 5. The obtained explanations contain well-defined image boundaries at iteration 10 and slowly converges to its peak performance at iteration 15. Figure 5 shows the improvement in the quality of explanations with the increasing number of iterations. Figures 6 and 7 show

the quantitative performance both at an image and pixel-level with increasing number of iterations. As we can observe, at the image-level, the proposed IASSA seems to reach its peak performance at iteration 15 and deteriorate post-peak due to oversampling. Whereas, when evaluated at the pixel level, the proposed method IASSA’s performance increases across all metrics but deletion, suggesting the reduction in the influence of model-dependent saliency.

#### 4.3. Comparison with State-of-the-art approaches

Figure 4 shows results comparing explanations obtained using the proposed method with LIME and RISE. The saliency maps obtained by our IASSA highlight regions of interest more accurately when compared to other state-of-the-art approaches. For example, the success of our approach can be qualitatively visualized in the test image for class “snowboard” in Figure 4 (row 4, column 5), while there exists some ambiguity if the person in the input image contributes to classification while using either LIME or RISE. The IASSA clearly points at the snowboard to classify the image.

We also summarize the results quantitatively in Table 1. From Table 1, we can observe that our proposed method IASSA outperforms the two black-box model explanation approaches with the added flexibility of generating expla-



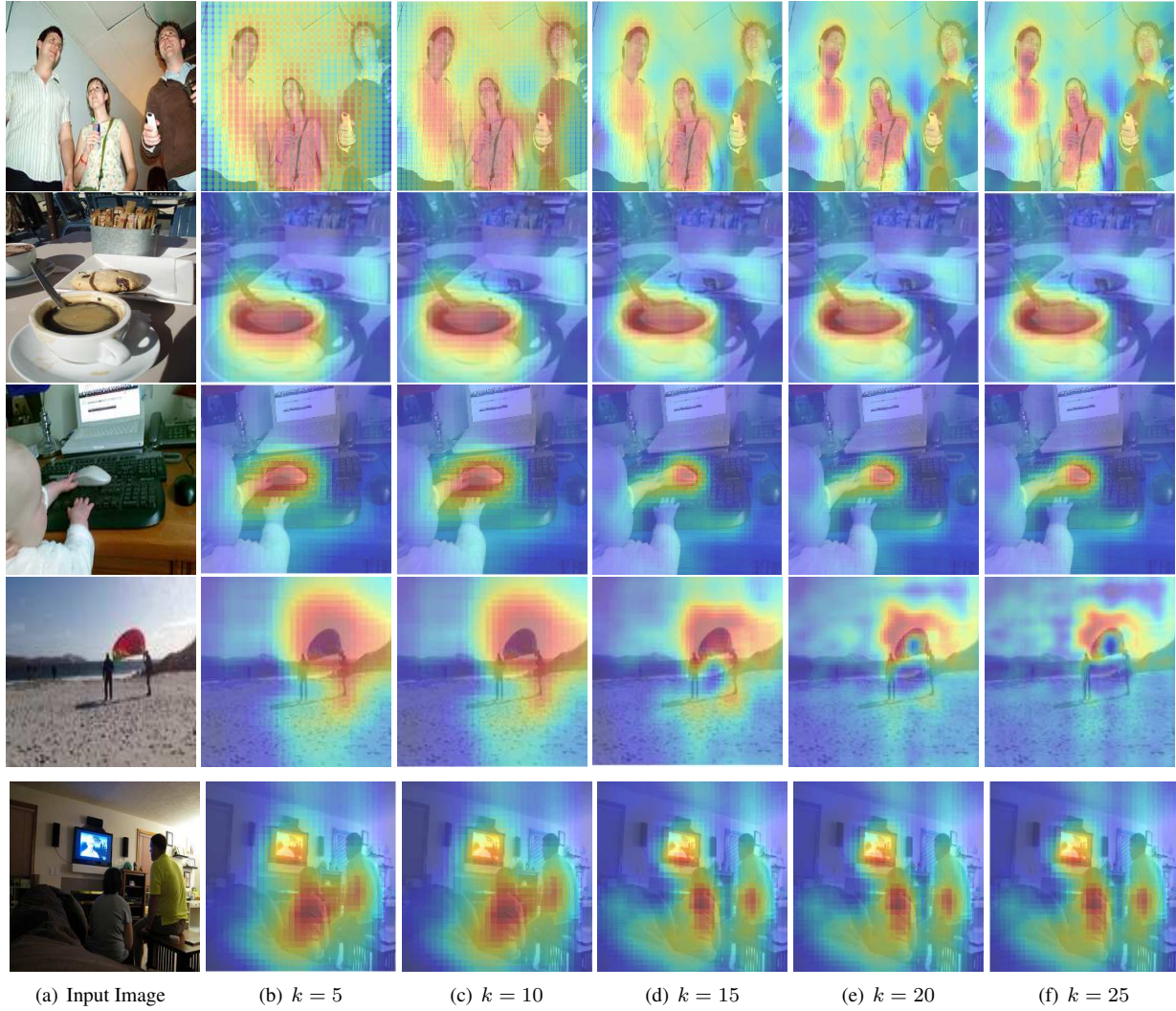


Figure 5. Visualization of our IASSA’s saliency maps with increasing iteration number  $k$  on the MS-COCO dataset.

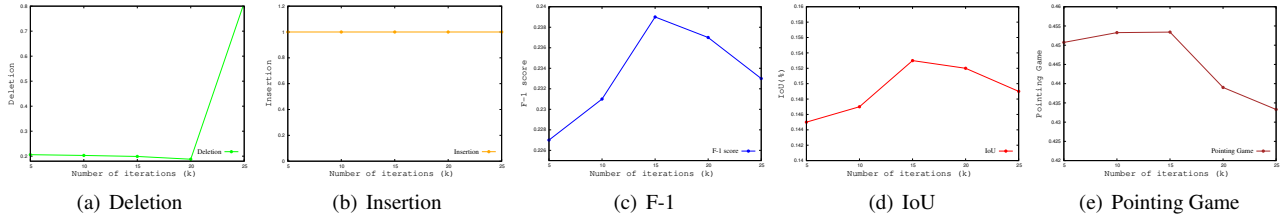


Figure 6. Performance of our IASSA at the image level with increasing number of iterations on the MS-COCO dataset.

nations in an iterative manner enabling its application in speed-critical explanation systems. When averaged at an image level, LIME is severely affected, especially in pointing game due to instances when the pixels with the highest activation were not aligned with the ground truth mask. The proposed model not only outperforms other explanation mechanisms when evaluated for “goodness” for the underlying model but also maintains human trust in explanation.

Even though RISE obtains deletion metrics close to the proposed system, IASSA gives the best of both worlds by explaining the model underneath and encapsulating object-

ness information at the same time. While our IASSA performs close to the best when evaluated at the image level, the true merit of our approach can only be appreciated at the pixel level. In an ideal explanation, we would expect all the contributing regions to carry the highest activation possible. Black box explanation approaches are prone to error in interpretation of an explanation due to extraneous image regions that affect human trust in explanation. Normalizing saliency maps with the number of pixels carrying the top 30% of the activations resolve this issue, resulting in a fair evaluation. The iterative aspect of our IASSA makes it a

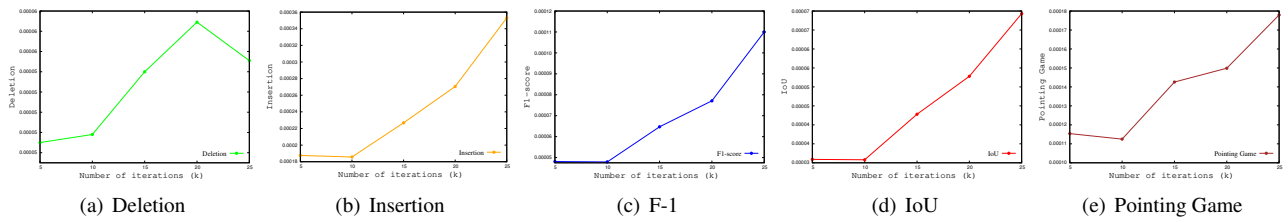


Figure 7. Performance of our IASSA at the pixel level with increasing number of iterations on the MS-COCO dataset.

Table 1. Comparative evaluation in terms of deletion (lower is better) and insertion (higher is better), F-1 (higher is better), IoU (higher is better), and Pointing Game (higher is better) scores at both image and pixel levels on the MS-COCO dataset.

	Method	Deletion ↓	Insertion ↑	F-1 ↑	IoU ↑	Pointing Game ↑
Image-level	LIME	0.900967	0.99	0.15390	0.09745	0.16461
	RISE	<b>0.1847</b>	<b>1.0</b>	0.13837	0.13653	0.25
	IASSA	0.18803	<b>1.0</b>	<b>0.23658</b>	<b>0.15153</b>	<b>0.4216</b>
Pixel-level	LIME	10.8526e-05	10.96158e-05	1.71177e-05	1.08447e-05	0.43671e-05
	RISE	5.5423e-05	28.8669e-05	4.26672e-05	2.69240e-05	8.95937e-05
	IASSA	<b>5.50534e-05</b>	<b>35.33639e-05</b>	<b>10.5960e-05</b>	<b>6.9282e-05</b>	<b>17.79331e-05</b>

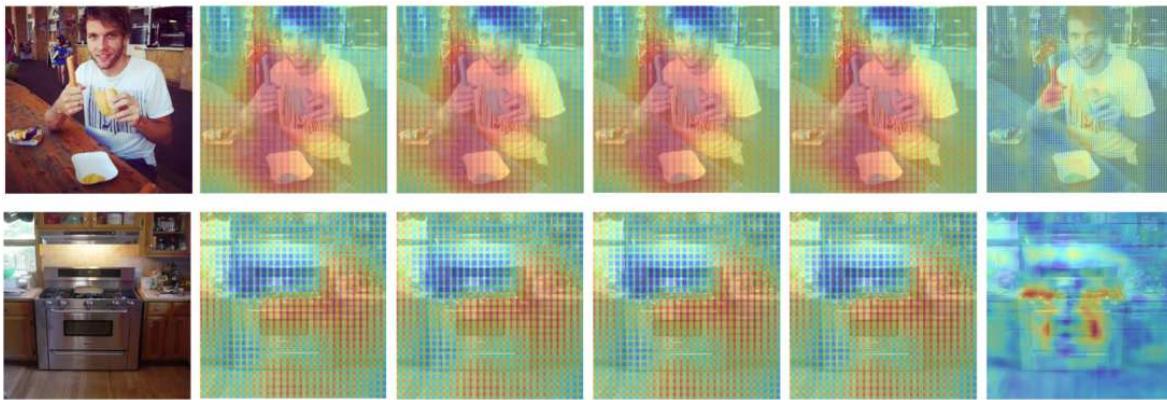


Figure 8. Explanations affected by sampling artifact that results in what can be an accurate explanation for class handbag and oven, but causes ambiguity due artifacts in the form of lines.

perfect match for applications that require the system to be scaled with minimal overhead.

#### 4.4. Discussion

Fine-tuning hyper-parameters such as  $W$  and  $S$ ,  $\beta$  and  $T_{thresh}$  plays a crucial role in determining performance. Hyperparameters help users control the quality of explanations and the algorithms convergence rate. Even though setting hyperparameters requires some knowledge about the underlying algorithm, we limit the range of values between a standard range of (0.0, 1.0) as opposed to it being arbitrary. The proposed system can result in explanations containing sampling artifacts due to a mismatch between window size  $W$  of stride  $S$ . To prevent this, we plan to look into other sampling methods that are both faster and can get a consensus on a larger image region at a time. Some examples of sampling artifacts are shown in Figure 8. Ultimately, the proposed system takes an average of approximately 800 milliseconds per iteration to compute explanation on an image of size  $224 \times 224$  using a ResNet-50 architecture in

batches of 256. Since a majority of the run-time is spent in loading the deep learning feature extractor, we advice using large batch sizes to minimize model load time.

#### 5. Conclusion

In this paper, we propose a novel iterative and adaptive sampling with a parameter-free long-range spatial attention for generating explanations for black-box models. The proposed approach assists in bridging the gap between model dependant explanation and human trustable explanation by laying the path for future research into methodologies to define “goodness” of an explanation. We prove the above claim by evaluating our approach using a plethora of metrics like deletion, insertion, IoU, F-1 score, and pointing game, at both the image and pixel levels. We believe the explanations obtained using the proposed approach could not only be used for the human to reason model decision but also contains generalized class specific information that could be fed back into the model to form a closed loop.



## References

- [1] S. A. Bargal, A. Zunino, V. Petsiuk, J. Zhang, K. Saenko, V. Murino, and S. Sclaroff. Guided zoom: Questioning network evidence for fine-grained classification. *ArXiv*, abs/1812.02626, 2018. [1](#)
- [2] C. Cao, X. Liu, Y. Yang, Y. Yu, J. Wang, Z. Wang, Y. Huang, L. Wang, C. Huang, W. Xu, et al. Look and think twice: Capturing top-down visual attention with feedback convolutional neural networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2956–2964, 2015. [2](#)
- [3] M. Caron, P. Bojanowski, A. Joulin, and M. Douze. Deep clustering for unsupervised learning of visual features. In *European Conference on Computer Vision (ECCV)*, 2018. [5](#)
- [4] J. Choi, D. Chun, H. Kim, and H.-J. Lee. Gaussian yolov3: An accurate and fast object detector using localization uncertainty for autonomous driving. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. [1](#)
- [5] B. Ding, C. Long, L. Zhang, and C. Xiao. Argan: Attentive recurrent generative adversarial network for shadow detection and removal. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 10213–10222, 2019. [1](#)
- [6] B. Dong, R. Collins, and A. Hoogs. Explainability for content-based image retrieval. In *The IEEE Conference on Computer Vision and Pattern Recognition Workshop (CVPRW)*, June 2019. [3](#)
- [7] R. C. Fong and A. Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 3429–3437, 2017. [2](#)
- [8] T. Hu, C. Long, L. Zhang, and C. Xiao. Vital: A visual interpretation on text with adversarial learning for image labeling. *arXiv preprint arXiv:1907.11811*, 2019. [1](#)
- [9] G. Hua, C. Long, M. Yang, and Y. Gao. Collaborative active learning of a kernel machine ensemble for recognition. In *IEEE International Conference on Computer Vision (ICCV)*, pages 1209–1216. IEEE, 2013. [1](#)
- [10] G. Hua, C. Long, M. Yang, and Y. Gao. Collaborative active visual recognition from crowds: A distributed ensemble approach. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 40(3):582–594, 2018. [1](#)
- [11] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision (ECCV)*, pages 740–755. Springer, 2014. [2](#), [5](#)
- [12] T. Lombrozo. The structure and function of explanations. *Trends in Cognitive Sciences*, 10(10):464–470, 2006. [1](#)
- [13] T. Lombrozo. The instrumental value of explanations. *Philosophy Compass*, 6(8):539–551, 2011. [1](#)
- [14] C. Long, A. Basharat, , and A. Hoogs. A coarse-to-fine deep convolutional neural network framework for frame duplication detection and localization in forged videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1–10. IEEE, 2019. [1](#)
- [15] C. Long, R. Collins, E. Swears, and A. Hoogs. Deep neural networks in fully connected crf for image labeling with social network metadata. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1607–1615. IEEE, 2019. [1](#)
- [16] C. Long and G. Hua. Multi-class multi-annotator active learning with robust gaussian process for visual recognition. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2839–2847, 2015. [1](#)
- [17] C. Long and G. Hua. Correlational gaussian processes for cross-domain visual recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 118–126, 2017. [1](#)
- [18] C. Long, G. Hua, and A. Kapoor. Active visual recognition with expertise estimation in crowdsourcing. In *IEEE International Conference on Computer Vision (ICCV)*, pages 3000–3007. IEEE, 2013. [1](#)
- [19] C. Long, G. Hua, and A. Kapoor. A joint gaussian process model for active visual recognition with expertise estimation in crowdsourcing. *International Journal of Computer Vision (IJCV)*, 116(2):136–160, 2016. [1](#)
- [20] C. Long, E. Smith, A. Basharat, and A. Hoogs. A c3d-based convolutional neural network for frame dropping detection in a single video shot. In *IEEE International Conference on Computer Vision and Pattern Recognition Workshop (CVPRW) on Media Forensics*. IEEE, 2017. [1](#)
- [21] C. Long, X. Wang, G. Hua, M. Yang, and Y. Lin. Accurate object detection with location relaxation and regionlets re-localization. In *The 12th Asian Conf. on Computer Vision (ACCV)*, pages 3000–3016. IEEE, 2014. [1](#)
- [22] X. Ma, Z. Wang, H. Li, P. Zhang, W. Ouyang, and X. Fan. Accurate monocular 3d object detection via color-embedded 3d reconstruction for autonomous driving. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. [1](#)
- [23] A. I. Maqueda, A. Loquercio, G. Gallego, N. Garca, and D. Scaramuzza. Event-based vision meets deep learning on steering prediction for self-driving cars. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [1](#)
- [24] M. Mitsuhara, H. Fukui, Y. Sakashita, T. Ogata, T. Hirakawa, T. Yamashita, and H. Fujiyoshi. Embedding human knowledge in deep neural network via attention map. *arXiv preprint arXiv:1905.03540*, 2019. [5](#)
- [25] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 3387–3395, 2016. [2](#)
- [26] T. L. Pedersen and M. Benesty. lime: Local interpretable model-agnostic explanations. *R Package version 0.4*, 1, 2018. [1](#)
- [27] V. Petsiuk, A. Das, and K. Saenko. Rise: Randomized input sampling for explanation of black-box models. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2018. [1](#), [2](#), [3](#), [5](#)
- [28] B. A. Plummer, M. I. Vasileva, V. Petsiuk, K. Saenko, and D. Forsyth. Why do these match? explaining the behavior of

- image similarity models. *arXiv preprint arXiv:1905.10797*, 2019. 1
- [29] Z. Qi, S. Khorram, and F. Li. Visualizing deep networks by optimizing with integrated gradients. *arXiv preprint arXiv:1905.00954*, 2019. 1
- [30] F. U. Rahman, B. Vasu, and A. Savakis. Resilience and self-healing of deep convolutional object detectors. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 3443–3447. IEEE, 2018. 1
- [31] M. T. Ribeiro, S. Singh, and C. Guestrin. "why should i trust you?": Explaining the predictions of any classifier. In *HLT-NAACL Demos*, 2016. 1
- [32] M. T. Ribeiro, S. Singh, and C. Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining (SIGKDD)*, pages 1135–1144. ACM, 2016. 2
- [33] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Niessner. Faceforensics++: Learning to detect manipulated facial images. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 1
- [34] W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K.-R. Müller. Explainable ai: Interpreting, explaining and visualizing deep learning. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, 2019. 1
- [35] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision (CVPR)*, pages 618–626, 2017. 1, 2
- [36] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013. 2
- [37] B. Vasu, F. U. Rahman, and A. Savakis. Aerial-cam: Salient structures and textures in network class activation maps of aerial imagery. In *2018 IEEE 13th Image, Video, and Multidimensional Signal Processing Workshop (IVMSP)*, pages 1–5. IEEE, 2018. 1
- [38] B. Vasu and A. Savakis. Visualizing the resilience of deep convolutional network interpretations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 107–110, 2019. 1
- [39] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 5998–6008, 2017. 4
- [40] J. Wagner, J. M. Kohler, T. Gindele, L. Hetzel, J. T. Wiedemer, and S. Behnke. Interpretable and fine-grained visual explanations for convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9097–9107, 2019. 1
- [41] H. Wang, Y. Fan, Z. Wang, L. Jiao, and B. Schiele. Parameter-free spatial attention network for person re-identification. *arXiv preprint arXiv:1811.12150*, 2018. 4
- [42] J. Wei, C. Long, H. Zou, and C. Xiao. Shadow inpainting and removal using generative adversarial networks with slice convolutions. *Computer Graphics Forum (CGF)*, 38(7):381–392, 2019. 1
- [43] B. Wu, X. Sun, L. Hu, and Y. Wang. Learning with unsure data for medical image diagnosis. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 1
- [44] X. Xing, Q. Li, H. Wei, M. Zhang, Y. Zhan, X. S. Zhou, Z. Xue, and F. Shi. Dynamic spectral graph convolution networks with assistant task training for early mci diagnosis. In *Medical Image Computing and Computer Assisted Intervention (MICCAI)*, pages 639–646, 2019. 1
- [45] C.-K. Yeh, C.-Y. Hsieh, A. S. Suggala, D. W. Inouye, and P. D. Ravikumar. On the (in)fidelity and sensitivity for explanations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 1
- [46] J. Yosinski, J. Clune, A. Nguyen, T. Fuchs, and H. Lipson. Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*, 2015. 2
- [47] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision (ECCV)*, pages 818–833. Springer, 2014. 2
- [48] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision (ECCV)*, pages 818–833. Springer, 2014. 3
- [49] H. Zhang, I. J. Goodfellow, D. N. Metaxas, and A. Odena. Self-attention generative adversarial networks. *arXiv:1805.08318*, 2018. 4
- [50] J. Zhang, S. A. Bargal, Z. Lin, J. Brandt, X. Shen, and S. Sclaroff. Top-down neural attention by excitation backprop. *International Journal of Computer Vision (IJCV)*, 126(10):1084–1102, 2018. 1, 2
- [51] L. Zhang, C. Long, X. Zhang, and C. Xiao. Ris-gan: Explore residual and illumination with generative adversarial networks for shadow removal. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2020. 1
- [52] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2921–2929, 2016. 2