

# ADA-SISE: ADAPTIVE SEMANTIC INPUT SAMPLING FOR EFFICIENT EXPLANATION OF CONVOLUTIONAL NEURAL NETWORKS

*Mahesh Sudhakar<sup>\*</sup>, Sam Sattarzadeh<sup>\*</sup>, Konstantinos N. Plataniotis<sup>\*</sup>,  
 Jongseong Jang<sup>†</sup>, Yeonjeong Jeong<sup>†</sup>, Hyunwoo Kim<sup>†</sup>*

<sup>\*</sup>Department of Electrical & Computer Engineering, University of Toronto

<sup>†</sup>Fundamental Research Lab, LG AI Research

## ABSTRACT

Explainable AI (XAI) is an active research area to interpret a neural network’s decision by ensuring transparency and trust in the task-specified learned models. Recently, perturbation-based model analysis has shown better interpretation, but backpropagation techniques are still prevailing because of their computational efficiency. In this work, we combine both approaches as a hybrid visual explanation algorithm and propose an efficient interpretation method for convolutional neural networks. Our method adaptively selects the most critical features that mainly contribute towards a prediction to probe the model by finding the activated features. Experimental results show that the proposed method can reduce the execution time up to 30% while enhancing competitive interpretability without compromising the quality of explanation generated.

**Index Terms**— CNNs, Deep Learning, Explainable AI, Interpretable ML, Neural Network Interpretability.

## 1. INTRODUCTION

Over the recent past years, access to a lot of digital data, the advances in computing facilities, and the facile access to many readily available pre-trained models have fueled the growth in deep learning. Although such models produce high accuracy in object recognition, the interpretability [1] of their decisions is also essential to convince the stakeholders or locate any potential bias in the underlying data. With AI currently being employed in various fields such as in healthcare, consumer retail, and banking, it is high time to develop “Responsible AI” [2] for society. To ensure the uniformity of the training data’s distribution, lately, there is an increase in modern open-source toolkits [3, 4] that acts as a common framework to evaluate a model’s fairness.

Explainable AI (XAI) has recently been offering many algorithms to interpret a model’s behavior. Based on their usage at the training process’s timeline, XAI approaches can be broadly classified into *ad-hoc* and *post-hoc* methods. In terms of their explanation ability to interpret a single instance or the whole decision process, XAI can be classified into *local* and

*global*. They can also be categorized into *model-agnostic* and *model-specific* methods, based on the requirement to specify the model’s architecture.

In this work, we study such a recent *post-hoc*, *local*, and *model-specific* XAI algorithm - Semantic Input Sampling for Explanation (SISE) [5] developed for image classification tasks. Building on this method, we propose a way to improve its run-time while retaining its overall performance without compromising the visual explanation’s quality. Our approach introduces a novel way to adaptively select the most important feature information to be considered for the subsequent steps of the algorithm’s operation. This modification acts as a smart filtering procedure that mutates the existing method into an automated, unified solution by eliminating the need for an end-user to tune the hyper-parameters. To demonstrate this claim, we evaluate our approach with the original algorithm’s performance in terms of the visual explanation quality, overall benchmark analysis, and execution time.

## 2. BACKGROUND

### 2.1. Existing methods

The prior works on *post-hoc* visual XAI can be divided into three main groups: ‘backpropagation-based’, ‘perturbation-based’, and ‘CAM-based’ methods. The backpropagation-based methods mainly operate by backpropagating the signals from the output neuron of the model to the inputs [6, 7] or the hidden nodes of the model [8], in order to calculate gradient [7] of relevance [9] terms. Perturbation-based approaches rely on feed-forwarding the model with perturbed copies of the input image. They interpret the model’s behavior using techniques such as probing the model with random masks [10] or optimizing a perturbation mask for each input [11]. Moreover, CAM-based methods are built based on the Class Activation Mapping (CAM) method [12] and are used specifically for CNNs by taking advantage of the phenomenon of this type of networks in weak object localization, as stated in [13]. Most of these methods are developed by backpropagation techniques [14, 15] or perturbation techniques [16].

## 2.2. Semantic Input Sampling for Explanation

SISE is a recent explanation method that spans among all three mentioned visual XAI methods, although it is generally classified as a perturbation-based algorithm. SISE employs the feature information underlying the model’s various depths to generate a saliency-based high-resolution visual explanation map. For a given trained classification model  $\delta : I \rightarrow \mathbb{R}^C$  with  $N$  convolutional blocks that outputs a confidence score over  $C$  classes for each input image  $I$ , SISE generates a 2-dimensional explanation map  $Y_{I,\delta(\lambda)}$  for  $\lambda$  in the domain of feature maps  $\Lambda$ , through its four-phased architecture.

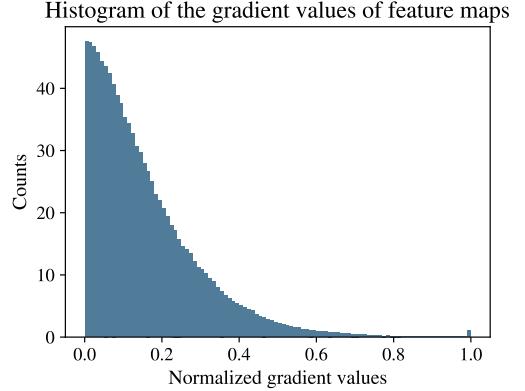
In the first phase (*Feature map Extraction*), pooling layers  $p_l$  of the model for  $l \in \{1, \dots, N\}$  are targeted, and their corresponding feature maps  $F_k^{[p]}$  for  $k \in \{1, \dots, M^p\}$  are collected. As this operation is independent of the classifier part, there would be a lot of irrelevant feature information about the background or other object classes (if present), in addition to the class of interest. The excess information is filtered out in the second phase (*Feature map Selection*) based on their backpropagation scores. Here, the feature maps with positive gradients towards a particular class are selected and post-processed to be converted into attribution masks  $A_k^{[p]}$ , via bilinear interpolation followed by normalization in the range  $[0, 1]$ .

The generated attribution masks are then scored by weighing based on their classification scores in the third phase (*Attribution mask Scoring*) and later combined to form a layer visualization map  $V_{I,\delta(\lambda)}^{[p]}$ . These preceding steps are repeated for all pooling (down-sampling) layers  $p_l$  of the network and then passed to the final phase (*Fusion*) of the algorithm, where they are fused in a cascading manner under a series of operations including addition, multiplication, and adaptive binarization, to reach the final explanation map.

## 3. PROBLEM STATEMENT

The gradient-based feature map selection policy in SISE is aimed to distinguish the feature maps containing essential attributions for the model’s prediction (‘positive-gradient’) against the ones representing outliers or background information. That was achieved using a threshold parameter  $\mu$  that was set to 0 by default to discard the ones with ‘negative-gradient’ scores.

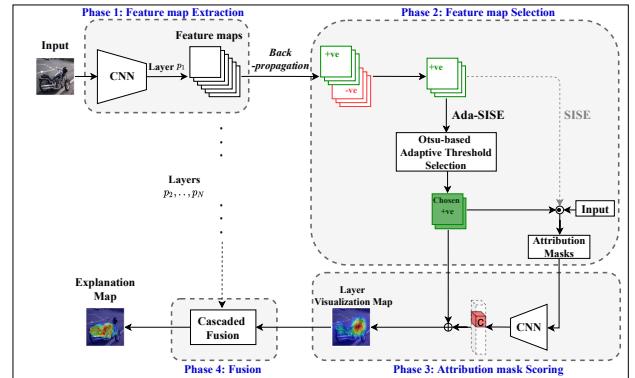
However, most of the elected activation maps with positive gradients are relatively ineffective in the prediction procedure, thereby increasing SISE’s computational overhead unnecessarily. We identify that the average gradient distribution of the positive-gradient feature maps follows a pattern, as in Fig. 1, where several trivial features are represented with low gradient values. Thus, only a fraction of the most critical feature maps is passed to the third phase of SISE. Hence, we focus on developing an adaptive selection policy for the parameter  $\mu$  of SISE to estimate the least number of required



**Fig. 1.** Histogram of the gradient values recorded from the feature maps in the last convolutional layer of a ResNet-50.

features to generate an explanation map without any notable compromise (and even in some cases, a slight enhancement) in terms of visual quality.

## 4. ADAPTIVE MASK SELECTION



**Fig. 2.** Architecture of the proposed Ada-SISE XAI method.

To tune the strictness of feature map selection adaptively for each of the layers, we employ an Otsu-based framework [17]. For a selected layer  $p$ , we reach the set of feature maps  $F_k^{[p]}$  and their corresponding gradient values  $\sigma_k^{[p]}$ , and determine its maximum as  $\rho^{[p]}$ . Denoting the normalized gradient values for the feature maps as  $v_k^{[p]} = \frac{\sigma_k^{[p]}}{\rho^{[p]}}$ , we define the set of positive-gradient feature maps as:

$$\Upsilon^{[p]} \equiv \Upsilon^{[p]+} = \{v_k^{[p]} > 0 | k \in \{1, \dots, M^{[p]}\}\} \quad (1)$$

where  $M^{[p]}$  is the number of feature maps extracted from layer  $p$ . Otsu’s method is applied to the set of positive-gradient feature maps to implement an updated threshold on them, based on the histogram of their average gradient scores. Assuming  $\Upsilon^{[p]}(i) \forall i \in \{1, \dots, |\Upsilon^{[p]}|\}$  to be the  $i$ -th

---

**Algorithm 1:** Ada-SISE: Adaptive Semantic Input Sampling for Explanation

---

**Input :** An input image  $I$  and a trained model  $\delta$ .  
 $\eta \leftarrow$  post-processing function.  
 $\zeta \leftarrow$  heatmap fusion function.

**for**  $n \leftarrow 1, \dots, N$  **do**

- Select the pooling layer  $p$  and collect feature maps  $F_k^{[p]} \forall k \in \{1, \dots, M^p\}$ ;
- Let the domain of the feature maps be  $\Lambda^{[p]}$ ;
- $\sigma_k^{[p]} = \sum_{\lambda^{[p]} \in \Lambda^{[p]}} \frac{\partial \delta(I)}{\partial F_k^{[p]}(\lambda^{[p]})} \text{ } \& \text{ } \rho^{[p]} = \max(\sigma_k^{[p]})$  ;
- $A_k^{[p]} \leftarrow []$  ;  $\Upsilon^{[p]} \leftarrow \{v_k^{[p]} > 0 | k \in \{1, \dots, M^{[p]}\}\}$ ;
- $\mu^{[p]} \leftarrow \Upsilon^{[p]}(\text{argmax}_{j \in \{1, \dots, |\Upsilon^{[p]}|\}} (\tau^{[p]}(j)))$ ;
- foreach**  $k \leftarrow \{1, \dots, m^p\}$  **do**

  - if**  $\frac{\sigma_k^{[p]}}{\rho^{[p]}} > \mu^{[p]}$  **then**
  - $A_k^{[p]} \leftarrow A_k^{[p]} \cup \eta(F_k^{[p]})$ ;
  - else**
  - $A_k^{[p]} \leftarrow A_k^{[p]}$ ;
  - end**

- end**
- $V_{I, \delta(\lambda)}^{[p]} = \mathbb{E}_{A^{[p]}} [\delta(I \odot m) \cdot C_m(\lambda)]$ ;

**end**

SISE explanation:  $Y_{I, \delta(\lambda)} = \zeta(V_{I, \delta(\lambda)}^{[p]})$

**Output:** A 2D explanation map  $Y_{I, \delta(\lambda)}$ .

---

value in  $\Upsilon^{[p]}$ , we can formulate the mean value of the masks with less/more gradient values than  $\Upsilon^{[p]}(i)$  respectively, as follows:

$$\omega_L^{[p]}(i) = \frac{\sum_{j=1}^i (\Upsilon^{[p]}(j))}{i} \times |\Upsilon^{[p]}| \quad (2)$$

$$\omega_H^{[p]}(i) = \frac{\sum_{j=i}^{|\Upsilon^{[p]}|} (\Upsilon^{[p]}(j))}{|\Upsilon^{[p]}| - i} \times |\Upsilon^{[p]}| \quad (3)$$

If we set  $\mu = \Upsilon^{[p]}(i)$  to divide the set of positive-gradient feature maps into two low and high subsets, the inter-class variance of these sets are calculated as follows:

$$\tau^{[p]}(i) = \omega_L^{[p]}(i) \times \omega_H^{[p]}(i) \times \left[ \frac{|\Upsilon^{[p]}| - i}{|\Upsilon^{[p]}|} - \frac{i}{|\Upsilon^{[p]}|} \right]^2 \quad (4)$$

which can be simplified as:

$$\tau^{[p]}(i) = \omega_L^{[p]}(i) \times \omega_H^{[p]}(i) \times \left[ \frac{|\Upsilon^{[p]}| - 2i}{|\Upsilon^{[p]}|} \right]^2 \quad (5)$$

According to [17], minimizing the intra-class variance for both classes simultaneously is equivalent to maximizing the inter-class variance in equation (5). Hence, we can identify the most deterministic feature maps in each layer by applying

a threshold which maximizes the inter-class variance accordingly:

$$\mu^{[p]} = \Upsilon^{[p]} \left( \underset{j \in \{1, \dots, |\Upsilon^{[p]}|\}}{\text{argmax}} (\tau^{[p]}(j)) \right) \quad (6)$$

The argmax operation in equation (6) is achieved by a simple search method. If the number of feature maps derived from a layer is not noticeably large, and if some of these feature maps are discarded as negative-gradient activation maps, a simple search method would not add any significant additional complexity to SISE framework. We term our method Ada-SISE and show its architecture in Fig. 2 and its methodology in Algorithm 1.

## 5. RESULTS

To compare Ada-SISE’s performance abreast with SISE, experiments were performed on the test set of the Pascal VOC 2007 dataset [18]. Two pre-trained models, a shallow VGG16 (with a test accuracy of 87.18%) and a residual ResNet-50 network (with 87.96% test accuracy), are directly loaded from the TorchRay library [10] to replicate the original experimentation setup. As it was reported in [5] that SISE meets or outperforms most of the state-of-the-art XAI methods like Grad-CAM [14], RISE [10] and Score-CAM [16], we restrict our comparisons only with Extremal Perturbation [11] (as it is one of the sophisticated perturbation-based methods) and SISE.

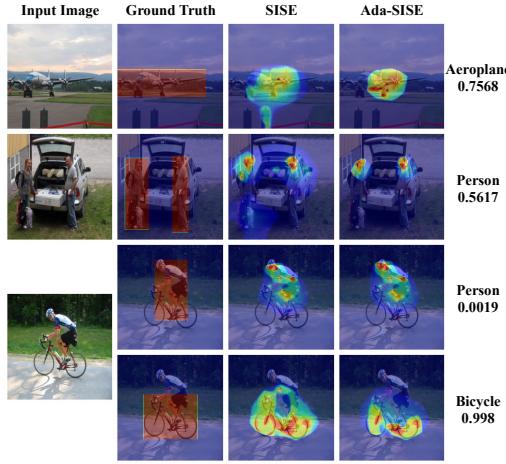
### 5.1. Benchmark Analysis

	Metric	Extremal Perturbation	SISE	Ada-SISE
VGG16	<b>EBPG</b>	<b>61.19</b>	60.54	60.79
	<b>Bbox</b>	51.2	55.68	<b>55.73</b>
	<b>Drop%</b>	43.9	<b>38.40</b>	38.87
	<b>Increase%</b>	32.65	37.96	<b>38.25</b>
	<b>Run-time (s)</b>	87.42	5.96	<b>4.23</b>
ResNet-50	<b>EBPG</b>	63.24	66.08	<b>66.4</b>
	<b>Bbox</b>	52.34	61.59	<b>61.77</b>
	<b>Drop%</b>	39.38	<b>30.92</b>	<b>30.92</b>
	<b>Increase%</b>	34.27	40.22	<b>40.75</b>
	<b>Run-time (s)</b>	78.37	9.21	<b>6.29</b>

**Table 1.** Results of benchmark evaluation of Ada-SISE on pre-trained models on the PASCAL VOC 2007 [18] dataset.

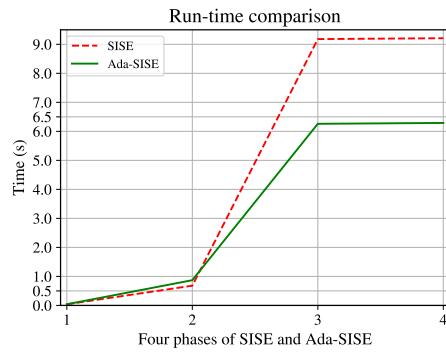
Table 1 shows the benchmark evaluation of Ada-SISE concerning various metrics and their execution time. As the depicted results are achieved through the same experimental setup as SISE paper, the readers can refer to [5] to infer further head-to-head comparison of Ada-SISE with other state-of-the-art methods. Energy-Based Pointing Game (EBPG) [16]

and Bbox [19] use the ground-truth annotations available to determine the precision of an XAI algorithm. Concurrently, Drop and Increase rates [20] measure the contribution of pixels captured in the explanation map towards the model’s predictive accuracy. Ada-SISE outperforms SISE in almost all of the metrics<sup>1</sup> while executing about 30% faster. Fig. 3 compares the explanation maps qualitatively on a ResNet-50 model and shows the ground-truth class and their annotations along with the model’s corresponding confidence score for each image.



**Fig. 3.** Comparison of Ada-SISE with SISE [5] on a ResNet-50 model with images from Pascal VOC 2007 dataset [18] demonstrating their class-discriminative explanation ability.

## 5.2. Run-time Analysis



**Fig. 4.** Comparison of the average run-times of Ada-SISE with SISE on a sample of images with a ResNet-50.

The bottleneck in SISE’s run-time is its third phase, where too many positive gradient feature maps are feed-forwarded

<sup>1</sup>For each metric in Table 1, the best is shown in bold. Besides Drop% and run-time (in seconds), the higher is better for all other metrics.

to compute their weights for scoring. As Ada-SISE chooses only a fraction of them that it considers crucial, our algorithm’s run-time is reduced significantly in the scoring phase.

Fig. 4 shows the comparison of run-times, where it can be noted that Ada-SISE executes under 6.3 seconds while SISE takes about 9.21 seconds. The small rise in the execution time at the second phase of Ada-SISE is the effect of our proposed adaptive thresholding procedure. The reported numbers are the average of experimentation performed over 100 random images from the Pascal VOC dataset on an NVIDIA Tesla T4 GPU with 16 GB of RAM.

## 5.3. Discussion

	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
<b>No. of feature maps available</b>	64	256	512	1024	2048
<b>SISE</b>	31	130	262	515	1008
<b>Ada-SISE</b>	26	114	179	420	551

**Table 2.** The number of feature maps chosen by Ada-SISE (on average) over the five pooling layers of a ResNet-50 compared with that of SISE and the corresponding number of available maps.

The number of feature maps selected for each pooling layer  $p_l$  of the network was recorded over a data sample of 500 images from the Pascal dataset, averaged, and reported in Table 2. As the deeper layers contribute more feature maps, it can be noticed that Ada-SISE chooses only a fraction of them, justifying the run-time reduction after the second phase. This validates our claim that by neglecting comparatively lower gradient values, dominant feature maps that contribute more towards a prediction can be extracted without compromising the explanation quality. Although an ablation study could be performed to identify a suitable value for  $\mu$  by fine-tuning SISE through extensive experiments, this solution would be profoundly dependent on the training data and would be brittle when expanded to new unseen data. Therefore, Ada-SISE generalizes SISE to be scaled for any application.

## 6. CONCLUSION

In this work, we propose Ada-SISE as an improvement to the recent SISE method that makes it a fully automated XAI algorithm. We also report a reduction in run-time and an overall improvement in the benchmark analysis without losing its visual explanation quality. Such identified important features would be adopted in future works to analyze a model’s behavior by studying its effect on the model’s prediction when replaced with noises or other classes’ attributions.

## 7. REFERENCES

- [1] Fenglei Fan, Jinjun Xiong, and Ge Wang, “On interpretability of artificial neural networks,” *arXiv preprint arXiv:2001.02522*, 2020.
- [2] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al., “Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai,” *Information Fusion*, vol. 58, pp. 82–115, 2020.
- [3] Julius A Adebayo et al., *FairML: ToolBox for diagnosing bias in predictive modeling*, Ph.D. thesis, Massachusetts Institute of Technology, 2016.
- [4] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, A Mojsilović, et al., “Ai fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias,” *IBM Journal of Research and Development*, vol. 63, no. 4/5, pp. 4–1, 2019.
- [5] Sam Sattarzadeh, Mahesh Sudhakar, Anthony Lem, Shervin Mehryar, KN Plataniotis, Jongseong Jang, Hyunwoo Kim, Yeonjeong Jeong, Sangmin Lee, and Kyunghoon Bae, “Explaining convolutional neural networks through attribution-based input sampling and block-wise feature aggregation,” *arXiv preprint arXiv:2010.00672*, 2020.
- [6] Mukund Sundararajan, Ankur Taly, and Qiqi Yan, “Axiomatic attribution for deep networks,” *arXiv preprint arXiv:1703.01365*, 2017.
- [7] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, “Deep inside convolutional networks: Visualising image classification models and saliency maps,” *arXiv preprint arXiv:1312.6034*, 2013.
- [8] Suraj Srinivas and François Fleuret, “Full-gradient representation for neural network visualization,” in *Advances in Neural Information Processing Systems*, 2019, pp. 4126–4135.
- [9] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek, “On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation,” *PloS one*, vol. 10, no. 7, pp. e0130140, 2015.
- [10] Vitali Petsiuk, Abir Das, and Kate Saenko, “Rise: Randomized input sampling for explanation of black-box models,” *arXiv preprint arXiv:1806.07421*, 2018.
- [11] Ruth Fong, Mandala Patrick, and Andrea Vedaldi, “Understanding deep networks via extremal perturbations and smooth masks,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2950–2958.
- [12] 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*, 2016, pp. 2921–2929.
- [13] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba, “Object detectors emerge in deep scene cnns,” *arXiv preprint arXiv:1412.6856*, 2014.
- [14] 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*, 2017, pp. 618–626.
- [15] Aditya Chattpadhyay, 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)*. IEEE, 2018, pp. 839–847.
- [16] Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and Xia Hu, “Score-cam: Score-weighted visual explanations for convolutional neural networks,” 2020.
- [17] N. Otsu, “A threshold selection method from gray-level histograms,” *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 9, no. 1, pp. 62–66, 1979.
- [18] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results,” 2007.
- [19] Karl Schulz, Leon Sixt, Federico Tombari, and Tim Landgraf, “Restricting the flow: Information bottlenecks for attribution,” *arXiv preprint arXiv:2001.00396*, 2020.
- [20] Harish Guruprasad Ramaswamy et al., “Ablation-cam: Visual explanations for deep convolutional network via gradient-free localization,” in *The IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 983–991.