

SUNY: A Visual Interpretation Framework for Convolutional Neural Networks from a Necessary and Sufficient Perspective

Xiwei Xuan^{1,*} Ziquan Deng¹ Hsuan-Tien Lin² Zhaodan Kong¹ Kwan-Liu Ma¹

¹University of California, Davis

{xwxuan, ziqdeng}@ucdavis.edu

²National Taiwan University

htlin@csie.ntu.edu.tw

*Corresponding Author

{zdkong, klma}@ucdavis.edu

Abstract

In spite of the ongoing evolution of deep learning, Convolutional Neural Networks (CNNs) remain the de facto choice for numerous vision applications. To foster trust, researchers have proposed various methods for visually interpreting CNNs via heatmaps, which highlight the input regions important to a specific model decision. However, in terms of the underlying design logic, existing approaches often concentrate on model parameters, overlooking the fundamental “why” question integral to human cognition. Thus they fail to embrace the two critical and complementary sides in reasoning: necessity and sufficiency. To address these issues, we introduce SUNY, a framework designed to rationalize the explanations toward better human understanding from both necessary and sufficient perspectives in a bi-directional manner. Extensive evaluations justify that SUNY not only yields more informative and convincing explanations from both angles, but also achieves performances competitive to other approaches across different CNN architectures over different datasets.

1. Introduction

Despite the unprecedented strides in deep learning, the interpretation of Convolutional Neural Networks (CNNs) continues to be an essential field of study, attributed to their pervasive application [35, 40], proven robustness [2, 20, 31], and inherently opaque nature [10]. This paper addresses the eXplainable Artificial Intelligence (XAI) [11] problem corresponding to CNN for natural image classification, i.e., reasoning why a classifier makes particular decisions. Specifically, we study visual explanation techniques that present heatmaps highlighting image portions associated with a model’s class prediction. A series of gradient-weighted CAMs [4, 17, 23] in the CAM [42] family are widely-adopted in applications. However, gradients’ saturation and vanishing issues can lead to noise explanations for such CAMs [8]. To bypass the shortcomings of gradi-

ents, Score-CAM [29] and Group-CAM [39] weight feature maps by contribution scores, referring to the corresponding input features’ importance to the model output. Reflecting on these methods, the design of Score-CAM and Group-CAM, which measures model’s prediction (outcome) by retaining specific input features (cause), aligns with the concept of causal *sufficiency* (S). Conversely, the design principle behind perturbation-based techniques [19, 21, 38], measuring model’s prediction (outcome) when changing input features (cause), aligns with the idea of causal *necessity* (N). In general, N involves changing hypothetical causes and measuring the resultant differences in outcomes, while S investigates whether preserving specific causes can maintain the outcome by quantifying outcome stability. Many studies underscore the cruciality of both N and S as “*two desirable and essential perspectives for a successful explanation*”, as they resonate with the bi-directional counterfactual thinking intrinsic to human cognition [9, 14, 32, 34].

Given the significance of these two facets, we propose an explanation framework called **S**UFFiciency and **N**ecessit**Y** explanation (**SUNY**), which interprets the CNN classifier by regarding the input features as the hypothesized causes and quantifying each cause’s importance towards the class prediction from angles of both N and S . Based on the qualitative evaluation, including the semantic evaluation and the sanity check, we demonstrate that **SUNY** provides a more faithful and interpretable visual explanation for CNN models. Comprehensive experiment evaluations on benchmark datasets, including ILSVRC2012 [22] and CUB-200-2011 [33], validate that **SUNY** outperforms other popular heatmap-based visual explanation methods.

2. The Proposed Approach

Our method aims to (1) measure the importance of each individual cause in a group of coordinating causes ($G1$), and (2) quantifying their actual impact on the outcome while considering both N and S ($G2$).

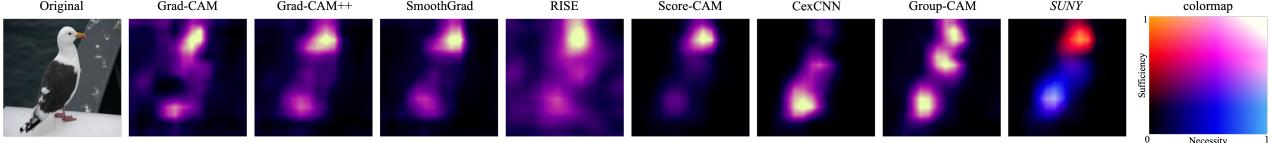


Figure 1. Visual comparison of heatmap explanations provided by different methods. All heatmaps in this paper follow the colormap on the right (X-axis: *Necessity*; Y-axis: *Sufficiency*; $y = x$: *Importance* (applicable to methods with one-dimensional information, “importance”).

2.1. Bi-directional Importance Quantification

The Shapley value’s framework for assessing marginal contributions across various coalitions lays the groundwork for achieving $G1$. Additionally, we further define the necessity and sufficiency value functions to accomplish $G2$. However, previous SHAP [15] image analyses segment input images into equally-sized patches, restricting finer distinctions. Our model-integrated method regards a single feature map (or a set of feature maps) as a cause f_i (or a set of causes F_*), thereby providing more granular explanations. We utilize a general formulation of Shapley values in the following definition.

For a set of causes (i.e., a set of feature maps) to be analyzed as F_* , we define \mathbf{N} value function to measure the degree of the outcome change when removing F_* :

$$E_N(F_*) = [p_c(I) - p_c(\text{do}(F \setminus F_*))]/p_c(I), \quad (1)$$

where $\text{do}(F \setminus F_*)$ represents the intervention of removing F_* . And $p_c(\cdot)$ refers to the model’s prediction probability w.r.t. a target class c , where $p_c(I)$ is its original value without any intervention; $p_c(\text{do}(F \setminus F_*))$ represents the value after the removal intervention. Similarly, \mathbf{S} value function is defined as:

$$E_S(F_*) = p_c(\text{do}(F_*))/p_c(I), \quad (2)$$

where $\text{do}(F_*)$ represents the intervention of only keeping F_* , i.e., removing $\{F \setminus F_*\}$.

Different from covering all elements, we tend to focus on more necessary (sufficient) ones. $\forall f_i \in F$, where f_i is a single cause, we set $F_* = \{f_i\}$ to calculate $E_N(f_i)$ ($E_S(f_i)$) and construct a set $F_N \subseteq F$ ($F_S \subseteq F$) by combining the relatively more necessary (*sufficiency*) f_i . Then to analyze a single cause $f_n \in F_N$, we calculate the N Shapley value as:

$$R_N(f_n) = \sum_{F' \subseteq \{F_N \setminus f_n\}} \frac{|F'|!(|F_N| - |F'| - 1)!}{|F_N|!} \times [E_N(F' \cup f_n) - E_N(F')] \quad (3)$$

Similarly, we can calculate S Shapley value for $f_s \in F_S$ as:

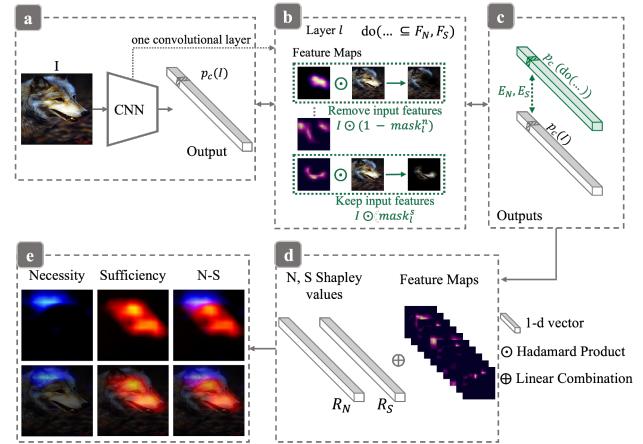


Figure 2. Overview of SUNY framework. Phase **a** is a forward pass of input image I through a CNN model. In Phase **b**, we obtain feature maps of a specified layer, and intervene on model filters or the corresponding input features. We get new prediction probabilities after the intervention and calculate E_N , E_S in Phase **c**, which are fed back to Phase **b** to construct hypothesized cause sets F_N and F_S . Through intervening on coalitions in F_N and F_S (Phase **b**), we can obtain their E_N , E_S (Phase **c**) and R_N and R_S (Phase **d**). The saliency maps are generated by linearly combining feature maps and examples of SUNY results are shown in Phase **e**.

$$R_S(f_s) = \sum_{F' \subseteq \{F_S \setminus f_s\}} \frac{|F'|!(|F_S| - |F'| - 1)!}{|F_S|!} \times [E_S(F' \cup f_s) - E_S(F')] \quad (4)$$

In the implementation, we reduce the amount of computation by estimating Eqns.(3), (4) using Shapley sampling values method [27]. Additionally, for $f'_n \in \{F \setminus F_N\}$ and $f'_s \in \{F \setminus F_S\}$, we set $R_N(f'_n) = 0$ and $R_S(f'_s) = 0$.

2.2. SUNY Implementation

Fig. 2 presents the SUNY framework applied to explain a CNN classifier w.r.t. the predicted class of an input image (a). Aligned with existing approaches ([4, 23, 25, 29, 39]), we consider feature maps of a convolutional layer as feature extractors to support our intervention. Specifically, each feature map can be upsampled into the image size as

a weighted mask, including values in the range of [0, 1]. Utilizing them, *SUNY* can remove/keep corresponding input regions for intervention, providing explanations for any convolutional layer. We first obtain its feature maps by forwarding an image into the model (Phase **b**). Next, we calculate E_N and E_S (refer to Eqns. (1)(2)) by intervening on every single cause, which is conducted by masking out specific input features ($I \odot (1 - mask_l^n)$; $I \odot mask_l^s$) (Phase **b** **c**). We then construct F_N and F_S by combining single causes with higher E_N and E_S , respectively, and repeat the aforementioned intervention operations on F_N and F_S (Phase **b** **c**). Finally, we calculate importance scores R_N and R_S based on Eqns. (3)(4)) and generate final visualizations (Phase **d** **e**).

Note that the general definition in Sec.2.1 is not limited to feature maps as causes. In Appendix, we detail the analysis when considering model filters as alternative causes.

3. Experiments

This section compares *SUNY*'s effectiveness with established visual explanation benchmarks.

3.1. Experimental Setup

Baseline Methods. To evaluate the effectiveness of our proposed method in pinpointing the crucial region for models' decisions, we carefully chose relevant methods for a comparative analysis with *SUNY*. The criteria for selection were twofold. First, these methods must possess class-discriminative capabilities, meaning they should provide unique explanations for different specified classes. Second, they should generate heatmaps that localize input regions relevant to the model's decision. Our selection encompasses seven methods that satisfy these requirements, covering a broad spectrum of approaches predominant in CNN visual explanations. These include gradient-based methods (Grad-CAM [23], Grad-CAM++ [4], and SmoothGrad [25]), a score-based method (Score-CAM [29]), causality-driven methods (CexCNN [6], and Group-CAM [39]), and a perturbation-based method (RISE [19]). Notably, we exclude pixel-space gradient visualizations such as Guided Backpropagation [26] or LRP [16] from our comparison due to their lack of class-discriminative ability, ensuring our baseline selection is aligned with our objectives.

Datasets and CNN Models. The experiments involve two datasets, ILSVRC2012 (ILSVRC) [22] with 1000 classes and 50k images and CUB-200-2011 (CUB) [33] with 200 classes and 5794 images. We use all explanation methods to explain three CNN models with different architectures, including VGG16 [24], Inception-v3 [28], and ResNet50 [13].

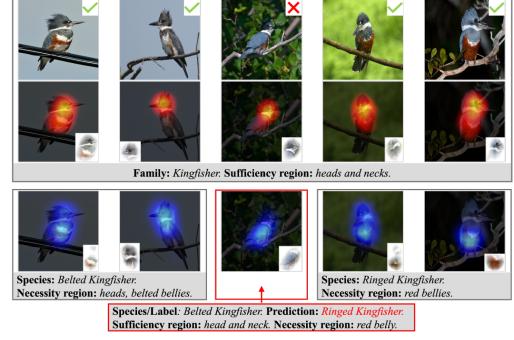


Figure 3. *SUNY* for a VGG16 trained on CUB. The first row displays bird images from four species across two families, with correct and incorrect predictions marked by \checkmark and \times . Misclassifications occur within the same family. The second and third rows show *sufficiency* and *necessity* heatmaps, with a small image in each heatmap's bottom corner highlighting the image portion.

3.2. Qualitative Evaluation

In Fig. 1, we visually compare *SUNY* with other explanation approaches and observe two advantages: (1) Saliency maps provided by *SUNY* contain fewer noises. (2) *SUNY* uniquely provides both *necessary* and *sufficient* information to support interpretation. For example, in Fig. 1, *SUNY* tells that the bottom wing is *necessary* and the head is *sufficient* for Gull prediction.

SUNY explanations for failure cases. Fig. 3 presents CUB images of two bird species (i.e., belted kingfisher and ringed kingfisher) of one family (i.e., kingfisher). The *sufficiency* heatmaps in the second row reveal family-specific features: all kingfisher display characteristic heads and necks. This explains why the model correctly identifies the bird family for every image. The *necessity* heatmaps in the third row provide additional insights to distinguish between different species within the same family: the belted bellies are highlighted in images predicted as belted kingfisher, while the red bellies are highlighted in the images predicted as ringed kingfisher. This explains why the third image is mistaken – the red belly is easily observable through this view and is identical to a ringed kingfisher. Examples in Fig. 3 demonstrate that *sufficiency* and *necessity* provide semantically-complementary explanations to better support model behavior interpretation.

3.3. Deletion and Insertion Evaluation

We evaluated model prediction relevancy to highlighted regions using *deletion* and *insertion* experiments, following [19]. *Deletion* gauges prediction impact when key pixels are removed, while *insertion* tracks prediction changes as pixels are added back in order of importance. Using Area Under the Curve (AUC) for quantification [19, 39], superior

Dataset	Methods	VGG 16			Inception-v3			ResNet50		
		Deletion ↓	Insertion ↑	Overall ↑	Deletion ↓	Insertion ↑	Overall ↑	Deletion ↓	Insertion ↑	Overall ↑
ILSVRC	Grad-CAM[23]	0.1098	0.6112	0.5015	0.1276	0.6567	0.5291	0.1796	0.6889	0.5093
	Grad-CAM++[4]	0.1155	0.6033	0.4878	0.1309	0.6476	0.5167	0.1847	0.6799	0.4952
	SmoothGrad [25]	0.1136	0.6023	0.4887	0.1317	0.6465	0.5148	0.1849	0.6800	0.4951
	RISE [19]	0.1185	0.6188	0.5003	0.1404	0.6444	0.5040	0.1303	0.6932	0.5629
	Score-CAM [29]	0.1070	0.6382	0.5312	0.1309	0.6528	0.5219	0.2319	0.6218	0.3898
	CexCNN [6]	0.1161	0.6025	0.4864	0.1355	0.6543	0.5188	0.1886	0.6443	0.4557
	Group-CAM [39]	0.1138	0.6218	0.5080	0.1292	0.6545	0.5253	0.1794	0.6904	0.5110
	SUNY	0.1005	0.6468	0.5462	0.1215	0.6603	0.5388	0.1323	0.6988	0.5665
	SUNY-N	0.1057	0.6038	0.4981	0.1257	0.6453	0.5196	0.1374	0.6552	0.5178
	SUNY-S	0.1144	0.6389	0.5245	0.1309	0.6530	0.5221	0.2220	0.6922	0.4702
CUB	Grad-CAM[23]	0.0558	0.7617	0.7059	0.0963	0.7323	0.6360	0.0930	0.6452	0.5522
	Grad-CAM++[4]	0.0589	0.7541	0.6951	0.0950	0.7281	0.6331	0.0972	0.6407	0.5434
	SmoothGrad [25]	0.0594	0.7489	0.6895	0.0977	0.7244	0.6266	0.0974	0.6405	0.5431
	RISE [19]	0.0560	0.7583	0.7023	0.0855	0.7168	0.6314	0.0570	0.6567	0.5996
	Score-CAM[29]	0.0542	0.7575	0.7033	0.0901	0.7326	0.6424	0.0995	0.6351	0.5355
	CexCNN [6]	0.0630	0.7389	0.6760	0.1017	0.7283	0.6267	0.1014	0.6173	0.5159
	Group-CAM [39]	0.0606	0.7521	0.6915	0.0971	0.7290	0.6318	0.0926	0.6458	0.5532
	SUNY	0.0518	0.7591	0.7073	0.0842	0.7361	0.6519	0.0562	0.6645	0.6083
	SUNY-N	0.0537	0.7497	0.6960	0.0854	0.7165	0.6311	0.0667	0.6443	0.5776
	SUNY-S	0.0555	0.7577	0.7022	0.0894	0.7328	0.6434	0.0939	0.6577	0.5638

Table 1. Comparative evaluation between SUNY and baselines w.r.t. the *deletion*, *insertion*, and *overall* AUC, where lower *deletion*, higher *insertion*, and higher *overall* indicate a better explanation. The first and second best performances are marked in green and blue, respectively. SUNY-N and SUNY-S are not included for performance ranking.

model explanations are reflected by lower *deletion*, higher *insertion* and higher *overall (insertion-deletion)* scores. As shown in Table 1, SUNY consistently equals or exceeds baselines in most comparisons.

3.4. Saliency Attack

Researchers have proposed a series of local adversarial attack approaches [5, 7, 30, 36] guided by saliency maps, which is to fool CNN models by perturbing a small image region highlighted by saliency maps. These methods require the saliency maps to be “minimal and essential” [5]. Inspired by these insights, we propose an evaluation metric, $Attack_{score} = \frac{FlipRate}{AvgAttackSize}$, to validate whether SUNY explanations can detect the most important regions w.r.t. the model’s decision. After applying Gaussian noise to the saliency regions, we check any decision changes: $Flip = 1$ if $argmax(p(I)) \neq argmax(p(I'))$. To validate whether the region is “minimal”, we include *AvgAttackSize*, which is the average size of all saliency maps. The results reported in Table 2 proves that SUNY are better at highlighting the most important image region corresponding to the model’s decision.

3.5. Localization Evaluation

This section evaluates saliency map localization using the energy-based pointing game [29], aiming to measure the localization ability of saliency maps using the ground-truth bounding box of the target class, *bbox*. The input image is binarized with *bbox* by assigning the inside and outside regions with 1 and 0, respectively. Then, we apply the Hadamard product between the binarized input and the

Dataset	Methods	Saliency $Attack_{score} \uparrow$		
		VGG16	Inception-v3	ResNet50
ILSVRC	Grad-CAM[23]	0.9615	1.0435	0.7674
	Grad-CAM++[4]	0.9991	0.9821	0.8751
	SmoothGrad [25]	1.0449	0.9675	0.8776
	RISE [19]	0.9928	0.7353	1.0259
	Score-CAM[29]	0.5326	0.9673	0.3378
	CexCNN[6]	1.6341	1.0653	0.6393
	Group-CAM[39]	1.1556	1.0200	0.8020
	SUNY	2.0452	1.9874	1.5619
	SUNY-N	1.6344	1.0885	1.0726
	SUNY-S	0.5434	0.9685	0.5564
CUB	GradCam[23]	0.5969	0.5985	0.4694
	GradCam++[4]	0.6670	0.5950	0.5257
	SmoothGrad [25]	0.7783	0.5929	0.5260
	RISE [19]	0.5063	0.3860	1.1286
	Score-CAM[29]	1.2215	0.5989	0.8027
	CexCNN[6]	1.2673	0.5898	0.4171
	Group-CAM[39]	0.6742	0.5951	0.4991
	SUNY	2.8111	1.0475	1.7747
	SUNY-N	1.5863	0.7658	1.2083
	SUNY-S	1.2317	0.5884	0.8238

Table 2. Comparative evaluation between SUNY and baselines w.r.t. saliency attack scores (higher is better). (First and second best performances. SUNY-N and SUNY-S are not included for performance ranking.)

saliency map, the summary of which can quantify how much “energy” falls into *bbox*. The performance is measured by $Proportion = \frac{\sum map[i,j]_{(i,j) \in bbox}}{\sum map[i,j]}$. Table 3 reveals that SUNY, by simultaneously leveraging N and S features, outperforms other methods in terms of localization ability.

3.6. Sanity Check

The sanity check [1] verifies if a visual explanation method reliably reflects the model’s behavior. We conduct cascading randomization to the model’s weights from the top to

Dataset	Methods	Proportion (%) ↑		
		VGG16	Inception-v3	ResNet50
ILSVRC	Grad-CAM[23]	57.68	66.35	59.84
	Grad-CAM++[4]	61.31	65.93	61.74
	SmoothGrad[25]	62.18	65.78	61.75
	RISE[19]	58.93	59.26	59.48
	Score-CAM[29]	64.25	65.94	66.72
	CexCNN[6]	65.24	66.33	57.39
	Group-CAM[39]	62.70	66.17	60.68
	<i>SUNY</i>	65.61	66.71	68.02
CUB	Grad-CAM[23]	43.06	40.05	39.02
	Grad-CAM++[4]	45.45	40.45	41.25
	SmoothGrad[25]	47.12	40.34	41.28
	RISE[19]	37.28	34.74	36.32
	Score-CAM[29]	49.68	40.67	47.42
	CexCNN[6]	37.13	41.38	41.22
	Group-CAM[39]	43.53	41.08	40.36
	<i>SUNY</i>	49.97	41.96	43.21

Table 3. Comparative evaluation w.r.t. energy-pointing games’ proportion (higher is better). (First and second best performances.)

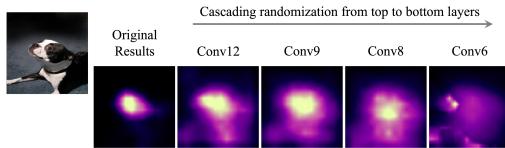


Figure 4. Sanity check of *SUNY*. The first column is the original heatmap visual explanation, and the following columns show results after randomizing specific layers.

the bottom layer successively and generate explanations every time after the randomization. If saliency maps are consistent across models with different parameters, the method does not pass the sanity check. Fig. 4 indicates *SUNY* pass the sanity check.

4. Conclusion

We design *SUNY*, a framework that offers bidirectional visual explanations of CNNs, integrating both *necessity* and *sufficiency* aspects. Qualitative assessments confirm that *SUNY* generates deeper, more meaningful visualizations, illustrating the added value of combining *necessity* and *sufficiency*. Furthermore, *SUNY* also passes the sanity check and quantitatively outperform seven other visual explanation methods in deletion and insertion evaluation, saliency attack, and localization evaluations across multiple CNN architectures and benchmark datasets.

Acknowledgement

This work is supported in part by NIBIB under grant no. P41 EB032840. H.-T. Lin is partially supported by the Ministry of Science and Technology in Taiwan via MOST 111-2918-I-002-006 and 112-2628-E-002-030.

References

- [1] Julius Adebayo, Justin Gilmer, Michael Muellly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. *arXiv preprint arXiv:1810.03292*, 2018. 4
- [2] Yutong Bai, Jieru Mei, Alan L Yuille, and Cihang Xie. Are transformers more robust than cnns? *Advances in Neural Information Processing Systems*, 34:26831–26843, 2021. 1
- [3] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6541–6549, 2017. 8
- [4] 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. 1, 2, 3, 4, 5
- [5] Zeyu Dai, Shengcui Liu, Ke Tang, and Qing Li. Saliency attack: Towards imperceptible black-box adversarial attack. *arXiv preprint arXiv:2206.01898*, 2022. 4
- [6] Hichem Debbi. Causal explanation of convolutional neural networks. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 633–649. Springer, 2021. 3, 4, 5
- [7] Xiaoyi Dong, Jiangfan Han, Dongdong Chen, Jiayang Liu, Huanyu Bian, Zehua Ma, Hongsheng Li, Xiaogang Wang, Weiming Zhang, and Nenghai Yu. Robust superpixel-guided attentional adversarial attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12895–12904, 2020. 4
- [8] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In *Proceedings of the AAAI conference on artificial intelligence*, pages 3681–3688, 2019. 1
- [9] Guy Grinfeld, David Lagnado, Tobias Gerstenberg, James F Woodward, and Marius Usher. Causal responsibility and robust causation. *Frontiers in Psychology*, 11:1069, 2020. 1
- [10] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)*, 51(5):1–42, 2018. 1
- [11] David Gunning. Explainable artificial intelligence (xai). *Defense Advanced Research Projects Agency (DARPA), nd Web*, 2, 2017. 1
- [12] Sergiu Hart. Shapley value. In *Game theory*, pages 210–216. Springer, 1989. 7
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3
- [14] Peter Lipton. Contrastive explanation. *Royal Institute of Philosophy Supplements*, 27:247–266, 1990. 1
- [15] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017. 2, 7

- [16] Gr  oire Montavon, Alexander Binder, Sebastian Lapuschkin, Wojciech Samek, and Klaus-Robert M  ller. Layer-wise relevance propagation: an overview. *Explainable AI: interpreting, explaining and visualizing deep learning*, pages 193–209, 2019. 3
- [17] Daniel Omeiza, Skyler Speakman, Celia Cintas, and Komminist Weldermariam. Smooth Grad-CAM++: An enhanced inference level visualization technique for deep convolutional neural network models. *arXiv preprint arXiv:1908.01224*, 2019. 1
- [18] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 8
- [19] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*, 2018. 1, 3, 4, 5
- [20] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *Advances in Neural Information Processing Systems*, 34:12116–12128, 2021. 1
- [21] Marco Tulio Ribeiro, Sameer Singh, and Carlos 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*, pages 1135–1144, 2016. 1
- [22] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 1, 3
- [23] 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. 1, 2, 3, 4, 5
- [24] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3
- [25] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Vi  gas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017. 2, 3, 4, 5
- [26] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014. 3
- [27] Erik   rumbelj and Igor Kononenko. Explaining prediction models and individual predictions with feature contributions. *Knowledge and information systems*, 41(3):647–665, 2014. 2
- [28] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 3
- [29] 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*, 2019. 1, 2, 3, 4, 5
- [30] Jie Wang, Zhaoxia Yin, Jing Jiang, and Yang Du. Attention-guided black-box adversarial attacks with large-scale multi-objective evolutionary optimization. *International Journal of Intelligent Systems*, 2022. 4
- [31] Zeyu Wang, Yutong Bai, Yuyin Zhou, and Cihang Xie. Can cnns be more robust than transformers? *arXiv preprint arXiv:2206.03452*, 2022. 1
- [32] David S Watson, Limor Gulchin, Ankur Taly, and Luciano Floridi. Local explanations via necessity and sufficiency: Unifying theory and practice. In *Uncertainty in Artificial Intelligence*, pages 1382–1392. PMLR, 2021. 1
- [33] Peter Welinder, Steve Branson, Takeshi Mita, Catherine Wah, Florian Schroff, Serge Belongie, and Pietro Perona. Caltech-ucsd birds 200. 2010. 1, 3
- [34] James Woodward. Sensitive and insensitive causation. *The Philosophical Review*, 115(1):1–50, 2006. 1
- [35] Ting-Wei Wu, Hua Zhang, Wei Peng, Fan L  u, and Pin-Jing He. Applications of convolutional neural networks for intelligent waste identification and recycling: A review. *Resources, Conservation and Recycling*, 190:106813, 2023. 1
- [36] Tao Xiang, Hangcheng Liu, Shangwei Guo, Tianwei Zhang, and Xiaofeng Liao. Local black-box adversarial attacks: A query efficient approach. *arXiv preprint arXiv:2101.01032*, 2021. 4
- [37] Yu Yang, Seungbae Kim, and Jungseock Joo. Explaining deep convolutional neural networks via latent visual-semantic filter attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8333–8343, 2022. 8
- [38] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014. 1, 8
- [39] Qinglong Zhang, Lu Rao, and Yubin Yang. Group-cam: group score-weighted visual explanations for deep convolutional networks. *arXiv preprint arXiv:2103.13859*, 2021. 1, 2, 3, 4, 5
- [40] Xiaorui Zhang, Jie Zhou, Wei Sun, and Sunil Kumar Jha. A lightweight cnn based on transfer learning for covid-19 diagnosis. *Computers, Materials & Continua*, 72(1), 2022. 1
- [41] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Object detectors emerge in deep scene cnns. *arXiv preprint arXiv:1412.6856*, 2014. 8
- [42] 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. 1

A. Appendix

This appendix section is organized as follows:

- Sec. A.1 presents the flexibility of *SUNY* to consider either a feature map or a model filter in a convolutional layer as a cause for analysis.
- Sec. A.2 provides a comparison between *SUNY* and previous SHAP image analyses [15].
- Sec. A.3 provides a more detailed description of *SUNY* implementation.
- Sec. A.4 provides more *SUNY* examples for qualitative evaluations.
- Sec. A.5 presents the training setting for our experiments.

A.1. A General Framework

As mentioned at the end of Section 2.1, our proposed framework provides the flexibility to consider either a feature map or a model filter in a convolutional layer as a cause for analysis. We differentiate between these two options as *SUNY*-feature and *SUNY*-filter in the Appendix, and provide an overview of the *SUNY*-filter framework in Fig. A1.

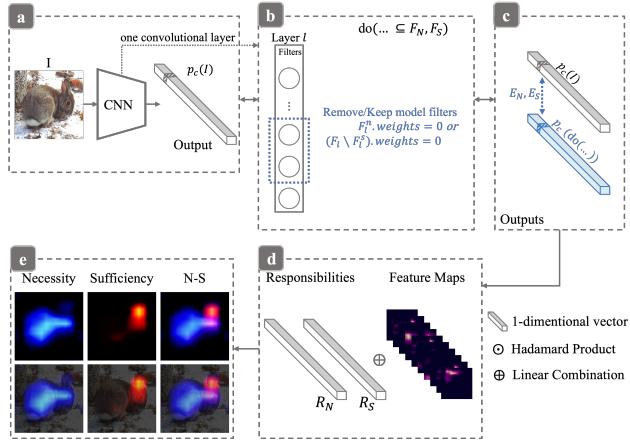


Figure A1. Overview of *SUNY*-filter framework.

A.2. Comparison with SHAP Methods

For our requirement *G1*, i.e. “to measure the importance of each individual cause in a group of coordinating causes,” shapely value [12] provides us a rational optimal solution to quantify the marginal contributions across diverse coalitions:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [v(S \cup \{i\}) - v(S)], \quad (\text{A1})$$

where ϕ_i is the Shapley value for an element i , F is the set of all elements, S is a subset of F that does not include element i , $v(S)$ is the value function that gives the total payoff that the subset of elements S can obtain by themselves. In our work, $E_N(\cdot)$ and $E_S(\cdot)$ (Eqns. (1) and (2),

respectively) can be seen as two value functions with different causal semantics. In practical applications, however, Equation A1 encounters a widely recognized issue of computational complexity. The primary reason for this issue is attributable to the excessive size of the total set S . Considering the unique attributes of the CNN model as applied to image processing, we have implemented the following modifications:

- **Constraining the scope of the subset S .** Rather than blindly covering all subsets and elements, we tend to focus on more necessary (sufficient) ones. We retain individuals with a higher $E_N(i)$ ($E_S(i)$) to form F_N (F_S) to constrain the scopes of the subset S and the element i .
- **Combine the model to select element i effectively.** In previous SHAP image analysis [15], images are divided into uniformly sized patches, each treated as element i for Shapley value calculation. This method limits the use of smaller patches due to high computational demands, leading to SHAP saliency maps that lack fine-grained importance distinctions, as illustrated in the Fig. A2. Our method, when integrated with the model, allows for the selection of either feature maps or filters as the object of analysis, thereby facilitating the provision of fine-grained visualization results.

Furthermore, we provide value functions (Eqns.(1) and (2)) with distinct causal interpretations to aptly characterize the significance derived from N and S .



Figure A2. SHAP[15] explanation example.

A.3. Algorithm Details

We present the implementation of *SUNY* in Alg. 1. Following the definitions in Sec. 2.1 in the main paper, *SUNY* provides causality-driven CNN visual explanations regarding input features or model filters as hypothesized causes, respectively, represented by the cause type E in Alg. 1 with values “feature” or “filter”. The corresponding explanations are *SUNY*-feature and *SUNY*-filter. We first define $p_c(\cdot)$ as a function to calculate the model’s prediction probability w.r.t. a class c for the input denoted by \cdot , as shown in line 1 of Alg. 1. Next, in line 3, we construct F_N and F_S , where F_N is the set of hypothetical single causes f_n with relatively higher **N Effect**, $E_N(f_n)$, (refer to Eqn. (1) in Sec. 2.1 of the main paper). Similarly, F_S contains all hypothetical single causes f_s with higher **S Effect**, $E_S(f_s)$, (refer to Eqn. (2) in the main paper). We then calculate **N-S Responsibilities** for every single cause in F_N and F_S following lines 5 - 20 (refer to Eqn. (3), (4)

Algorithm 1 SUNY: Causal Explanation of CNN

Require: Image I , model M , layer l , class c , cause type E
Ensure: N-S saliency maps: N_{map}, S_{map}

- 1: $p_c(\cdot) = \text{Softmax}(M(\cdot))[c]$ \triangleright prediction probability on c
- 2: $A_l \leftarrow M_l(I)$ \triangleright feature maps of the layer l
- 3: $F_N, F_S \leftarrow \text{getHypCauses}(I, M, l, E)$
- 4: $R_N, R_S \leftarrow \text{zeros}(A_l.\text{shape}[0])$ \triangleright initialize responsibilities
- 5: **if** E is “feature” **then**
- 6: **for** A_l^n **in** F_N **do**
- 7: $mask \leftarrow \text{norm}(\text{upsample}(A_l^n))$
- 8: Compute $R_N(A_l^n)$ based on Eqn.(3)
- 9: **for** A_l^s **in** F_S **do**
- 10: $mask \leftarrow \text{norm}(\text{upsample}(A_l^s))$
- 11: Compute $R_S(A_l^s)$ based on Eqn.(4)
- 12: **else if** E is “filter” **then**
- 13: **for** F_l^n **in** F_N **do**
- 14: $M^n \leftarrow \text{pruneFilters}(M, F_l^n)$
- 15: $p_c^n(\cdot) = \text{Softmax}(M^n(\cdot))[c]$
- 16: Compute $R_N(F_l^n)$ based on Eqn.(3)
- 17: **for** F_l^s **in** F_S **do**
- 18: $M^s \leftarrow \text{pruneFilters}(M, (F_l \setminus F_l^s))$
- 19: $p_c^s(\cdot) = \text{Softmax}(M^s(\cdot))[c]$
- 20: Compute $R_S(F_l^s)$ based on Eqn.(4)
- 21: $N_{map} = \text{norm}(\text{upsample}(\text{Relu}(\sum R_N^i A_l^i)))$
- 22: $S_{map} = \text{norm}(\text{upsample}(\text{Relu}(\sum R_S^i A_l^i)))$
- 23: **return** N_{map}, S_{map}

in the main paper). Specifically, for SUNY-feature (lines 6 - 11), we upsample and normalize the feature maps, A_l^n and A_l^s , and use the generated $mask$ as a feature extractor to intervene on input features from the image I . The intervention $do(F \setminus F_*)$ (removing F_*) is realized by the Hadamard product, $(I \odot (1 - mask))$. Similarly, $do(F_*)$ (keeping F_*) is implemented as $(I \odot mask)$. SUNY-filter (lines 13 - 20) removes and keeps the hypothesized causes by filter pruning, which means setting the corresponding filters’ weights to zero. Line 14 and line 18 correspond to $do(F \setminus F_l^n)$ and $do(F_l^s)$, respectively. After calculating **N-S Responsibilities** for all single causes, in lines 21 and 22, we obtain small saliency maps by $\text{Relu}(\sum R_N^i A_l^i)$, $\text{Relu}(\sum R_S^i A_l^i)$ and then upsample and normalize them to get the final saliency maps. The operation $norm$ represents the min-max normalization w.r.t. each single map, $norm(X) = \frac{X - \min(X)}{\max(X) - \min(X)}$, and $upsample$ represents the bilinear interpolation.

A.4. Additional Heatmap Examples

In this section, we provide more examples of semantic evaluations, where Fig. A3 and Fig. A4 present comparisons between SUNY and other methods, and Fig. A5 include more

causal explanation examples to demonstrate the usefulness of *necessity* and *sufficiency* ((a)(b)(c)) and textual explanation examples (as discussed below) to further examine more semantically-meaningful explanations with SUNY.

Textual explanations with SUNY. Recent research indicates that filters in a convolutional layer act as concept detectors [38, 41] and describes each filter with text [3, 37]. To examine visual-textual explanations with SUNY, we adopt [3] to automatically name the filters in a convolutional layer and present the top filters based on **N-S Responsibilities** (R_N, R_S) provided by SUNY, where higher R_N, R_S mean the corresponding filters are more necessary, sufficient, respectively. Fig. A5(d) shows some examples of visual and textual explanations with SUNY-filter, where we sort filters in descending order of the normalized R_N and R_S respectively, and present the top filters’ textual explanations. For each image, SUNY visualizations highlight regions that are *Sufficient* or *Necessary* for a specific prediction, and the complementary textual explanations for the top N and S filters further explain why the corresponding regions are essential for the class. By discussing the textual explanations with SUNY, we extend the usefulness of our technique in the ability to provide both intuitive and semantically-meaningful explanations.

A.5. Experiment Details

We implement SUNY and the seven aforementioned visual explanation methods in Python using PyTorch framework [18]. Specifically, we run the official or publicly available code of other methods as the results at the same data scale are unavailable. The platform is equipped with two NVIDIA RTX 3090 GPUs. Unless explicitly stated, the comparisons discussed in this section are conducted following the same settings, including –

(1) Images are resized and converted to the RGB format, transformed to tensors, and normalized to the range of $[0, 1]$.

(2) For every single input image across three datasets, visual explanation results are generated to explain every model with respect to (w.r.t.) every ground-truth class.

(3) For images with multiple ground-truth classes, the overall performance is measured by the mean quantification score across all classes.

Bird Species Classifier Training Details.

We trained bird species classifiers using CUB-200-2011 (CUB) with pre-trained VGG16, ResNet50, and Inception_v3, respectively. All models are retrieved from the TorchVision model zoo. We use the cross entropy loss and the SGD optimizer with a momentum of 0.9. The learning rate is 0.001 for VGG16 and ResNet50, and 0.0001 for Inception_v3.

Saliency Attack Details.

In the saliency attack evaluation, the image regions highlighted by a specific visual explanation method are cor-

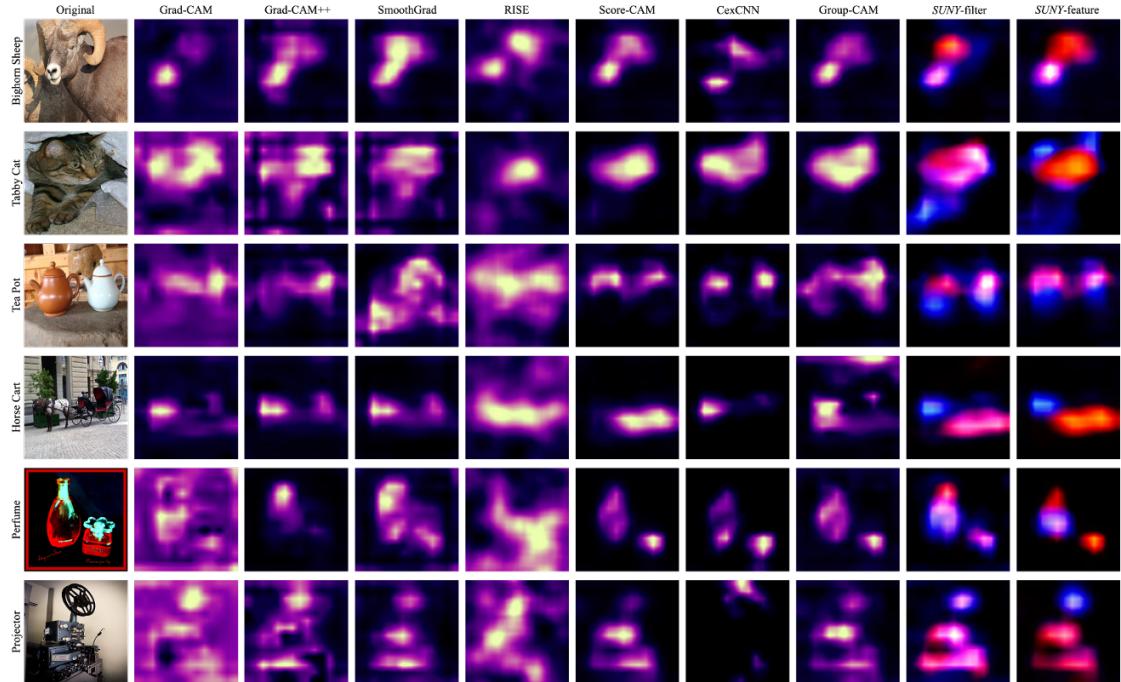


Figure A3. Visual comparison of saliency maps from different methods based on a VGG16 trained on ILSVRC. Each row is corresponding to one image from the ILSVRC validation set that is classified correctly by the model. The specified class for generating explanations is the predicted class, which is shown at the beginning of each row.

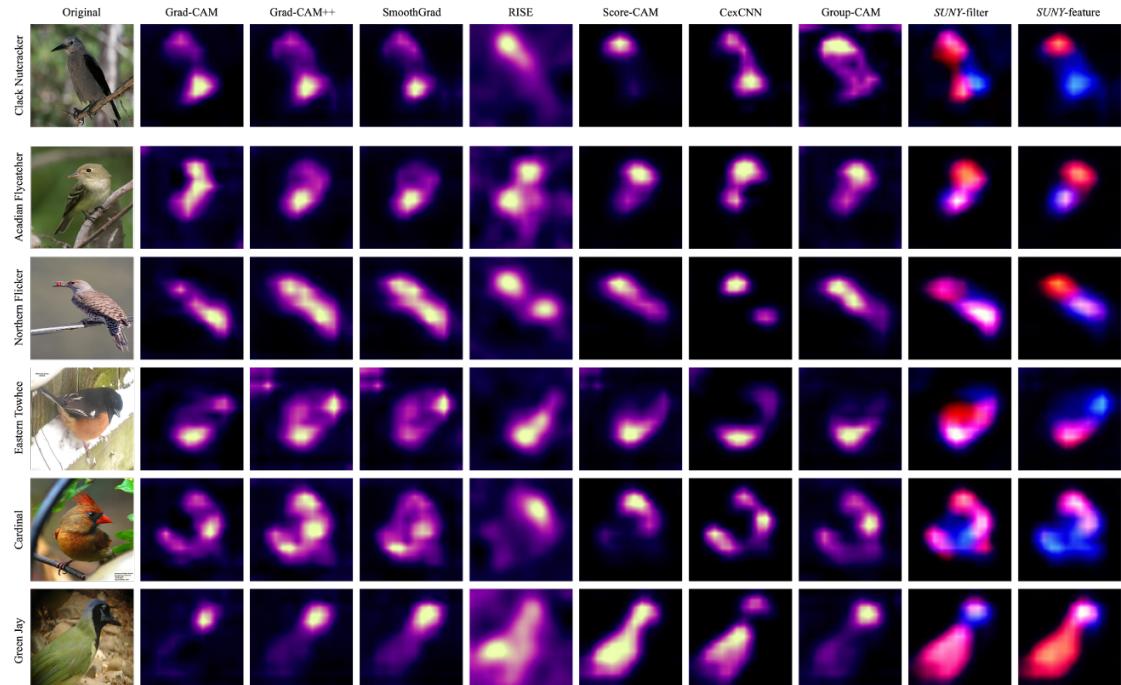


Figure A4. Visual comparison of saliency maps from different methods based on a VGG16 trained on CUB. Each row is corresponding to one image from the CUB validation set that is classified correctly by the model. The specified class for generating explanations is the predicted class, which is shown at the beginning of each row.

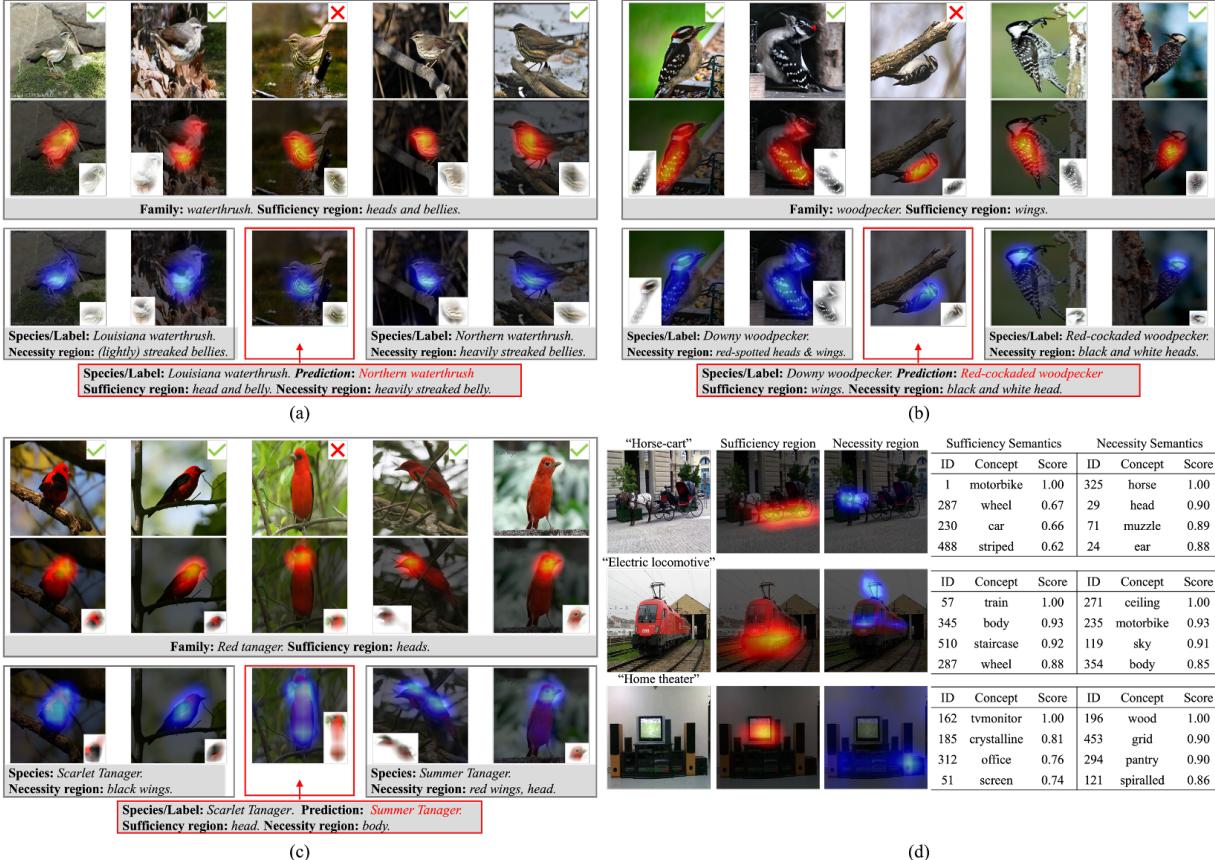


Figure A5. (a)-(c) Semantic evaluation of SUNY explanations for a VGG16 trained on CUB for bird species classification. The bird images in the first row are from four bird species belonging to two families and the correct/incorrect predictions are marked by \checkmark and \times , respectively. For the two images marked by \times , the model mistakes the actual species with the other species under the same family. Each column corresponds to one image; the second and third rows are *sufficiency* and *necessity* heatmaps, respectively. The small image in the bottom corner of each heatmap presents the highlighted image portion. (d) Examples of visual and textual explanations with SUNY-filter for a VGG16 trained on ILSVRC. For textual explanations, we provide top filters corresponding to the highest R_S and R_N , respectively.

rupted using random Gaussian noise with mean $\mu = 0$, variance $\sigma^2 = 0.03$. Specifically, we generate a noise matrix *noise* of the same size as the original input image I , where $\text{noise}[i, j] \sim \mathcal{N}(0, 0.03)$. For each visual explanation method, we have its saliency map as a matrix, where $\text{map}[i, j] \in [0, 1]$. The final image that we feed into the model is generated by $I' = I + \text{noise} \odot \text{map}$. Recalling our paper, the size of the region highlighted by the saliency map is calculated by $\text{Size}_{\text{map}} = \text{AttackSize} = \frac{\sum_{i,j} (1_{\text{true map}[i,j] \neq 0})}{\text{map.h} \times \text{map.w}}$.

In Fig. A6, we provide examples corresponding to one original image, where we add noise to the saliency regions and test whether such noise can fool the model. The original image shown in Fig. A6 is classified correctly. From the \checkmark and \times marks over the images with noise, we can find that the saliency attack is successful for most explanations (wrong classifications) except Score-CAM and Cex-

CNN (correct classifications). Moreover, for the explanations corresponding to successful saliency attacks, SUNY-filter and SUNY-feature have the smallest Size_{map} , indicating the most “minimal and essential” attacks.

The entire SUNY codebase and experiment setup will be made publicly available upon the publication of this paper.

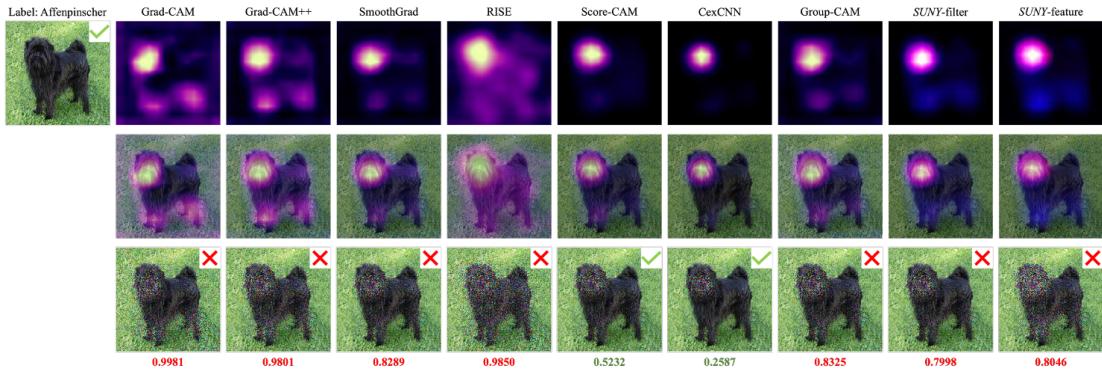


Figure A6. Saliency attack results based on different visual explanations. The original image on the top left is from the validation set of ILSVRC with the label *Affenpinscher*, which is classified correctly by the model. For the image grid, from the first to third row: heatmap, heatmap on the image, image with noise on the heatmap-highlighted region. The **✓** and **✗** indicate the model's correct/incorrect prediction of the corresponding image. Each column corresponds to one visual explanation. The number below each column is the size of the region highlighted by the saliency map ($Size_{map}$), where **red** indicates a **successful** saliency attack and **green** indicates a **failed** attack (For **successful** attacks, lower $Size_{map}$ is better).