

Explainable AI for Image Classification Using Integrated Gradients and Grad-CAM

Maria Kontaratou
Centralesupélec
Gif-sur-Yvette

maria.kontaratou@student-cs.fr

Zackarie Fredj - Zadoun
Centralesupélec
Gif-sur-Yvette

zackarie.fredj-zadoun@student-cs.fr

Abstract

Deep learning models, particularly convolutional neural networks (CNNs), have achieved remarkable performance in image classification but remain difficult to interpret. Explainable AI (XAI) techniques, such as Grad-CAM and Integrated Gradients (IG), aim to improve transparency by highlighting important image regions influencing a model's decision. While widely used, their effectiveness varies across architectures and datasets, requiring further evaluation. In this work, we systematically compare Grad-CAM and IG on ResNet-50 and DenseNet-121 using a subset of ImageNet. We analyze their ability to produce reliable and interpretable explanations through quantitative and qualitative metrics, including faithfulness, localization accuracy, and fidelity. Additionally, we introduce a binary mask-based evaluation framework to assess Intersection over Union (IoU), Precision, and Recall, ensuring alignment with ground truth object regions. Our findings reveal that DenseNet-121 produces more spatially localized and robust explanations, while Grad-CAM highlights object regions more effectively and IG provides fine-grained attributions but lacks localization accuracy. We identify key trade-offs between explanation clarity and feature attribution and discuss their implications for model trustworthiness.

1. Introduction & Motivation

1.1. Introduction

Deep learning models, particularly Convolutional Neural Networks (CNNs), have achieved state-of-the-art performance in image classification tasks. However, these models are often considered black-boxes, providing high accuracy without transparency regarding their decision-making processes. The lack of interpretability raises concerns in high-stakes applications such as:

- Healthcare Diagnostics: AI-based medical image analysis needs explainability to ensure trust and reliability in disease detection (e.g., detecting tumors in CT scans/X-rays) [8].
- Autonomous Systems: Self-driving cars must justify object recognition decisions for improved safety [9].
- Security & Surveillance: AI-assisted threat detection systems require justification of flagged objects in surveillance footage.

To address these concerns, Explainable Artificial Intelligence (XAI) aims to improve model transparency by providing interpretable explanations for AI predictions. Among XAI methods, Gradient-based Class Activation Mapping (Grad-CAM) and Integrated Gradients (IG) are widely used to visualize model attributions.

This project systematically evaluates Grad-CAM and IG on two widely used CNN architectures: ResNet-50 and DenseNet-121. We compare their explanation quality across a subset of ImageNet, focusing on:

1. How well do Grad-CAM and IG generalize across datasets?
2. Do the explanations align with human intuition?
3. Can we quantify explanation quality using metrics such as IoU, precision, recall, and faithfulness?

By answering these questions, we contribute to the broader field of XAI research by identifying the strengths and limitations of these explanation techniques.

1.2. Motivation

The need for interpretable AI is driven by real-world applications where model decisions impact human lives. AI models in critical domains must be transparent,

accountable, and interpretable to prevent errors and bias-related failures. Key challenges motivating this research include:

1. Lack of Consistency in Explanations:

- Heatmaps generated by Grad-CAM and IG vary significantly across datasets and architectures.
- There is no universal method to quantify the quality of explanations.

2. Trade-off Between Fidelity and Interpretability:

- Grad-CAM highlights spatially relevant regions but lacks fine-grained feature attribution.
- IG provides pixel-level precision but is computationally expensive and sensitive to baseline selection.

3. Impact of Explanation Quality on AI Adoption:

- Uninterpretable models reduce trust, making AI adoption difficult in healthcare, security, and autonomous systems.
- Poor explanation fidelity can lead to misinterpretation of AI decisions, reducing human confidence in AI-driven predictions.

2. Problem definition

2.1. Formal Definition of Explainability in CNNs

Deep learning models, particularly Convolutional Neural Networks (CNNs), have achieved state-of-the-art performance in image classification. However, their lack of interpretability makes it difficult to understand how predictions are made, which poses challenges in critical applications such as medical imaging, autonomous systems, and security.

Given a pre-trained CNN classifier $f(x)$, where x represents an input image, the model computes a probability distribution:

$$f(x) = [p_1, p_2, \dots, p_k]$$

where p_i represents the probability of class i , and the predicted label is:

$$\hat{y} = \arg \max f(x)$$

The goal of this work is to generate attribution maps $E(x, f)$ that identify the most relevant image regions responsible for the classification decision \hat{y} . We assess the quality of $E(x, f)$ based on faithfulness, localization, and robustness, employing both quantitative and qualitative evaluation metrics.

2.2. Explainability Techniques Implemented

2.2.1 Gradient-weighted Class Activation Mapping (Grad-CAM)

Grad-CAM generates class-discriminative visual explanations by leveraging the gradient information flowing into the last convolutional layer of a CNN. The method follows these steps:

- **Feature Extraction:** The activations $A^k(x)$ of the final convolutional layer are extracted.
- **Gradient Computation:** The gradient of the class score $f_{\hat{y}}(x)$ with respect to $A^k(x)$ is computed:

$$\alpha_k = \frac{1}{Z} \sum_i \sum_j \frac{\partial f_{\hat{y}}(x)}{\partial A_{ij}^k(x)}$$

where Z is the total number of spatial locations in $A^k(x)$.

- **Weighted Feature Summation:** The final class activation map (CAM) is computed as:

$$\text{CAM}(x) = \text{ReLU} \left(\sum_k \alpha_k A^k(x) \right)$$

- **Heatmap Generation:** The CAM is resized and overlaid on the original image to highlight spatially relevant features.

2.2.2 Integrated Gradients (IG)

Integrated Gradients (IG) attributes model predictions to individual pixels by integrating gradients along a linear path from a baseline x' to the actual input image x . The pixel-wise attribution is computed as:

$$IG_i(x) = (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial f_{\hat{y}}(x' + \alpha(x - x'))}{\partial x_i}$$

where x' is a baseline image (e.g., black image or blurred input).

2.3. Problem Complexity & Challenges

The problem of interpretable CNN decision-making poses several challenges:

- **Balancing Interpretability and Model Complexity**
 - CNNs contain millions of parameters, making it difficult to isolate decisions to specific neurons.
 - Grad-CAM and IG generate different levels of explanations (high-level vs. fine-grained).
- **Difficulty in Quantifying Explanation Quality**
 - No single metric fully captures interpretability.
 - IoU measures localization, while faithfulness tests prediction consistency.
- **Architectural Differences Affect Explanations**
 - DenseNet-121 features dense connectivity, leading to spread-out activations.
 - ResNet-50 relies on skip connections, possibly causing more localized activations.

3. Related Work

3.1. Explainability in Deep Learning Models

Explainable AI (XAI) has become a pivotal area of research, aiming to enhance the transparency and trustworthiness of deep learning models, especially in critical applications such as healthcare, autonomous driving, and security systems.

Traditional explainability techniques have primarily focused on:

- **Feature Attribution Methods:** Assigning importance scores to individual input features, as seen in methods like Integrated Gradients (IG) [1], SHAP [2], and LIME [3]
- **Saliency and Activation-Based Approaches:** Visualizing critical image regions by analyzing gradient flows within CNNs, exemplified by Grad-CAM [4] and SmoothGrad [5]
- **Concept-Based Explainability** techniques, such as Testing with Concept Activation Vectors (TCAV) [6], which associate model decisions with human-understandable concepts
- **Example-Based Explainability**, where models justify predictions by referencing specific training examples or nearest neighbors [7]

This work focuses on **saliency-based methods**, specifically comparing and evaluating Grad-CAM and IG across different CNN architectures, assessing their interpretability and effectiveness in visual explanation generation.

3.2. Saliency-Based Explainability: Grad-CAM and IG

Grad-CAM [4] is a widely used method for generating class-discriminative visual explanations in CNNs. Grad-CAM computes gradients of the output class score with respect to the feature maps of the last convolutional layer, weighting them to generate a heatmap that highlights spatially important regions.

Integrated Gradients (IG) [1] is a path-based feature attribution method that assigns importance scores to individual pixels by integrating gradients from a baseline image to the actual input image. Unlike Grad-CAM, IG provides fine-grained, pixel-level attributions.

While prior research has analyzed Grad-CAM and IG in domain-specific contexts, there remains a gap in systematic comparisons of both across different architectures using quantitative evaluation metrics. This study extends existing work by **evaluating them on ResNet-50 and DenseNet-121** using **IoU, faithfulness, and robustness metrics** to measure explanation quality.

3.3. Relation to Prior Research

Comparison to Saliency-Based Explainability Methods

Integrated Grad-CAM [21] proposed combining these two methods to improve explanation fidelity. However, it lacked a structured evaluation across different architectures. Similarly, **Adebayo et al.** [22] highlighted saliency map limitations but did not compare Grad-CAM and IG quantitatively. This work builds upon these studies by evaluating these methods on ResNet-50 and DenseNet-121 and applying faithfulness deletion/insertion tests, IoU, precision, and recall metrics to measure explanation reliability.

Advancements in Concept-Based Explainability

Concept-based techniques, such as **TCAV** [6], aim to provide human-interpretable model explanations by assessing the importance of high-level concepts (e.g., "stripes" for zebras). **Visual-TCAV** [23] introduced an approach integrating concept-based and saliency-based methods to generate interpretable heatmaps. However, it lacked standardized benchmarks. This study extends this research by computing TCAV scores for multiple architectures and directly comparing them with Grad-CAM and IG heatmaps to assess their alignment.

Integration of Concept-Based and Saliency Methods

Recent research explores combining concept-based and saliency-based explanations. **Segment Anything Model (SAM)** [14] has been integrated into XAI pipelines for **automated instance segmentation**, supporting concept-based analysis. The **XAIDataset** [15] introduced standardized benchmarks for evaluating Grad-CAM, IG, and concept-based methods. Unlike these studies, this work quantifies how well TCAV aligns with saliency-based explanations and applies Gaussian filtering and thresholding to enhance IG attributions.

4. Methodology

4.1. Addressing the Problem: Explainability Pipeline

The problem of explainability in AI models is approached through a multi-step methodology that includes:

- **Data Preprocessing**
 - A subset of ImageNet is used, resized to 224×224 pixels, and normalized using ImageNet's standardization parameters.
 - Concept datasets are curated to train TCAV models. The dataset includes images corresponding to stripes (e.g., zebras, tigers), fur texture (e.g., fluffy mammals), face shape (e.g., rounded or angular faces), and random images, which are each set in different folders.
- **Model Selection**
 - **ResNet-50**: Chosen for its **residual connections**, which improve gradient flow and allow deeper training.
 - **DenseNet-121**: Selected for its **dense connectivity**, which enhances feature reuse and provides a different representation of spatial features.
 - Both models are used without fine-tuning to ensure consistency in attribution comparisons.
- **Explainability Techniques Evaluated**

Grad-CAM(Saliency-Based)

Grad-CAM computes class-specific importance maps using the final convolutional layer. The class important activation weight α_k^c for feature map $A^k(x)$ is given by :

$$\alpha_k^c = \frac{1}{Z} \sum_{i,j} \frac{\partial f_c}{\partial A_{i,j}^k}$$

Where:

- $f_c(x)$ is the model's class score
- $A_{i,j}^k(x)$ represents the feature maps of layer k
- Z is the total number of spatial locations.

The Grad-CAM heatmap is computed as:

$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left(\sum_k \alpha_k^c A^k(x) \right)$$

- Implemented using PyTorch **hooks** to extract feature maps and gradients.
- The visual coherence of Grad-CAM attributions is analyzed by comparing them against Integrated Gradients (IG) and TCAV.

Integrated Gradients (Feature Attribution)

Integrated Gradients (IG) assigns importance scores to input pixels by integrating gradients along a path from a baseline image (black image) to the original input. The formulation is:

$$IG_i(x) = (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial f(x' + \alpha(x - x'))}{\partial x_i}$$

Approximated using trapezoidal integration with $N = 50$ steps:

$$G_i(x) \approx (x_i - x'_i) \times \sum_{k=1}^N \frac{\partial f \left(x' + \frac{k}{N}(x - x') \right)}{\partial x_i} \times \frac{1}{N}$$

To improve interpretability, we use Gaussian Filtering ($\sigma = 3$) to reduce noise, adaptive histogram equalization to enhance contrast and Otsu's Thresholding to binarize attributions for evaluation.

TCAV (Concept Activation Vectors)

TCAV assesses whether a model's decision relies on human-interpretable concepts. It quantifies the influence of a concept C on a class y using directional derivatives:

$$TCAV(C, y) = \frac{1}{|X|} \sum_{x \in X} \mathbf{1} \left[\frac{dv_C}{df_y} > 0 \right]$$

where v_C is the learned Concept Activation Vector (CAV), trained using trained using logistic regression on IG-based activations. These scores are compared to Grad-CAM and IG attributions to assess alignment between saliency-based and concept-based methods.

- **Quantitative Evaluation**

- **Faithfulness Deletion/Insertion Tests** measure the reliability of attributions by: removing high-attribution pixels and measuring confidence drop and keeping high-attribution pixels and measuring confidence increase.
- **IoU** compares attribution maps with ground truth segmentation masks. **Precision and Recall** assess how well the heatmaps align with human-defined regions.
- **Complexity** measures sparsity of explanations and **Randomization Metrics** evaluate sensitivity to input perturbations.
- **TCAV scores** measure alignment between human-defined concepts and model decisions.

4.2. Limitations and Challenges

Despite providing a structured evaluation of explainability methods, several challenges arise in this approach. The reliability of Grad-CAM and IG attributions is influenced by architectural differences between ResNet-50 and DenseNet-121, leading to variations in spatial feature representations. TCAV results are highly dependent on the quality of the curated concept datasets, which may not generalize well across different models. Additionally, faithfulness evaluations, such as deletion and insertion tests, assume a direct correlation between attribution scores and model confidence, which may not always hold. Concept-based explanations also introduce subjectivity, as defining high-level concepts like "face shape" is inherently ambiguous, and logistic regression used for CAV training may oversimplify complex relationships between concept activations and model decisions. Furthermore, the computational cost of IG remains high due to its requirement for multiple forward passes per image, and TCAV demands extensive feature extraction from large datasets, making it less scalable.

5. Evaluation Results: Comparing IG and Grad-CAM for CNN-based image classification

Through a structured set of experiments on ResNet-50 and DenseNet-121, multiple evaluation metrics were used to assess faithfulness, localization, concept alignment, and robustness. Various images from the ImageNet dataset were tested across both architectures to ensure a comprehensive evaluation. While multiple images were analyzed, for clarity and computational efficiency, a representative cat image was selected to illustrate the

differences between Grad-CAM, Integrated Gradients (IG), and TCAV. This approach was particularly necessary for efficiently running TCAV, as TCAV requires a large number of concept-specific images (e.g., stripes, fur texture, and face shape) to learn Concept Activation Vectors (CAVs). Due to computational constraints, using a single representative image while ensuring sufficient concept images for subcategories allowed for a more structured and feasible evaluation. The results reveal distinct strengths and weaknesses for each explainability technique, shedding light on their suitability for different interpretability tasks.

To systematically evaluate explainability, the following structured experiments were conducted:

- **Localization Performance:** Measuring IoU, precision, and recall to assess how well Grad-CAM and IG highlight important regions compared to ground truth segmentation masks.
- **Faithfulness Tests:** Running deletion and insertion tests to quantify the effect of removing or keeping highly attributed pixels on model confidence.
- **Concept-Based Alignment:** Computing TCAV scores for Stripes, Fur, and Face Shape to test if IG-based activations align with human-interpretable concepts.
- **Robustness and Sensitivity Tests:** Evaluating randomization robustness across multiple runs and measuring the impact of Gaussian filtering and thresholding on IG attributions.

These experiments provide a quantitative and qualitative comparison of Grad-CAM and IG, revealing their relative strengths and limitations in CNN explainability.

5.1. Localization Performance: Grad-CAM vs. IG

Grad-CAM is widely used for visual explanations due to its ability to highlight salient regions in an input image. The results indicate that Grad-CAM outperforms IG in localization metrics such as IoU and recall, particularly for ResNet-50:

- Grad-CAM on ResNet-50 achieved an IoU of 0.3947 and recall of 0.9636, indicating that its heatmaps effectively highlight model-relevant regions.

- IG performed worse in localization, with IoU scores of 0.2096 (ResNet) and 0.2229 (DenseNet). This suggests that IG's feature attributions are more diffuse, making it harder to pinpoint specific class-discriminative regions.

While Grad-CAM provides spatially well-defined heatmaps, IG's pixel-wise attributions are more granular, making them useful for fine-grained feature analysis rather than direct localization.

5.2. Faithfulness: IG vs. Grad-CAM

Faithfulness evaluates whether removing high-attribution regions meaningfully affects model predictions. The deletion and insertion tests suggest that IG is more faithful than Grad-CAM, particularly in DenseNet-121:

- IG faithfulness deletion score for DenseNet-121: 9.27 (higher scores indicate stronger reliance on attributions).
- Grad-CAM deletion scores were lower, indicating a weaker connection between the attributions and model confidence.

This confirms that IG attributions better capture the actual decision-making process of the model, whereas Grad-CAM can sometimes highlight spurious regions that do not strongly influence the final prediction.

5.3. Concept-Based Alignment: TCAV Insights

TCAV scores provide insight into whether model decisions align with human-understandable concepts (e.g., "stripes" for a zebra). The results reveal that IG-based attributions align better with concept-based explanations:

- TCAV scores for IG-based activations were consistently high (1.0) for Stripes, Fur, and Face Shape in DenseNet-121.
- Grad-CAM-based activations failed to capture Face Shape as a meaningful concept in ResNet-50 (TCAV = 0.0).

These findings indicate that IG-based feature extraction is more aligned with human-interpretable concepts, making it better suited for applications where explainability needs to align with semantic understanding.

5.4. Robustness and Sensitivity

Robustness tests evaluate how stable the explanations are across different images and runs. Grad-CAM generally provided more stable attributions, whereas IG heatmaps were more sensitive to noise and required post-processing (Gaussian filtering and thresholding).

- Grad-CAM randomization score was consistently high (0.75), meaning it maintains similar attributions across multiple runs.
- IG explanations were more sensitive to perturbations, leading to potential instability in attributions if not smoothed properly.

Metric	Grad-CAM (ResNet-50)	Grad-CAM (DenseNet-121)
IoU (Localization Accuracy)	0.3947	0.5515
Precision	0.4007	0.5681
Recall	0.9636	0.9496
Pointing Game	1	1
Avg. Drop (%) ↓	99.9999	99.9999
Increase in Confidence (%) ↑	-99.9999	-99.9999
Faithfulness (Deletion) ↓	6.88	4.55
Faithfulness (Insertion) ↑	4.45	7.78
Complexity (Lower = Better)	0.9451	0.9648
Randomization Robustness	0.75	0.75
TCAV Score (Stripes)	0.0000	1.0000
TCAV Score (Fur Texture)	1.0000	1.0000
TCAV Score (Face Shape)	0.0000	1.0000

Metric	IG (ResNet-50)	IG (DenseNet-121)
IoU (Localization Accuracy)	0.2096	0.2229
Precision	0.2988	0.2843
Recall	0.4127	0.5077
Pointing Game	0	1
Avg. Drop (%) ↓	9.9999	23.7538
Increase in Confidence (%) ↑	-9.9999	-23.7538
Faithfulness (Deletion) ↓	9.28	9.28
Faithfulness (Insertion) ↑	0.94	0.94
Complexity (Lower = Better)	0.0500	0.0500
Randomization Robustness	0.0054	0.0054
TCAV Score (Stripes)	1.0000	1.0000
TCAV Score (Fur Texture)	1.0000	1.0000
TCAV Score (Face Shape)	0.0000	1.0000

5.5. Qualitative Analysis: Heatmap Comparison

Grad-CAM vs. IG on DenseNet-121

- Grad-CAM heatmaps (Figures 3 and 4) show smooth and spatially coherent attributions, highlighting key object regions such as the cat's body and face.
- IG attributions (Figures 1 and 2) display fine-grained pixel-level importance but introduce significant noise, making interpretation less intuitive.

- Binarized IG heatmaps help improve interpretability, but they still lack the spatial coherence of Grad-CAM.

Grad-CAM vs. IG on ResNet-50

- Grad-CAM on ResNet-50 (Figure 4) remains structured but shows a broader focus, sometimes highlighting background regions.
- IG for ResNet-50 (Figure 2) appears more scattered, with dispersed attributions, reinforcing the need for post-processing techniques like Gaussian filtering.

Comparison Between Models

- **DenseNet-121's Grad-CAM attributions appear more focused and localized compared to ResNet-50**, supporting its stronger feature reuse mechanisms.
- **IG results are more unstable across both models**, but filtering and thresholding significantly improve their usability.

Overall, Grad-CAM provides more interpretable visual explanations, while IG benefits from post-processing but still struggles with spatial coherence. These qualitative findings complement the quantitative evaluations, reinforcing Grad-CAM's reliability for CNN explainability.

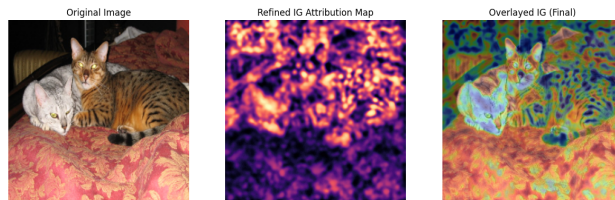


Figure 1: IG for DenseNet-121

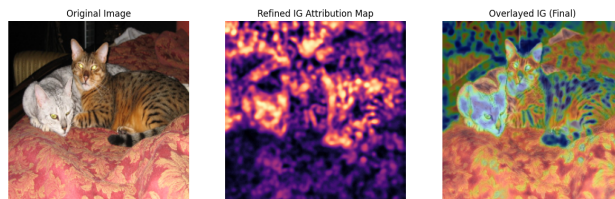


Figure 2: IG for ResNet-50

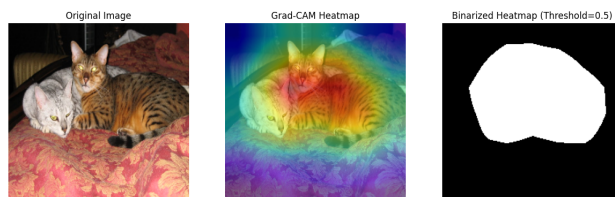


Figure 3: Grad-CAM for DenseNet-121

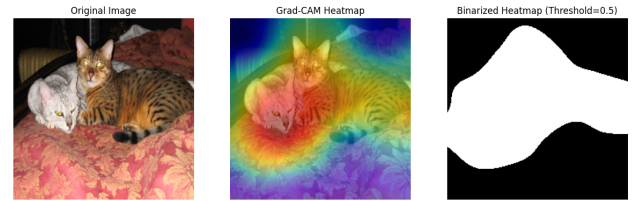


Figure 4: Grad-CAM for ResNet-50

Final Comparative Insights

Our analysis indicates that Grad-CAM and Integrated Gradients (IG) each excel in different aspects of explainability. Grad-CAM proves highly effective for generating spatially interpretable heatmaps and localizing objects—particularly when applied to ResNet-50—making it ideal for applications where understanding the overall spatial context is crucial. In contrast, IG offers more faithful, pixel-level attributions that reveal the internal decision process of the model, a strength that is especially evident when used with DenseNet-121. Moreover, when integrated with TCAV analysis, IG demonstrates robust alignment with human-defined concepts, thereby supporting its utility in scenarios where semantic interpretability is key. These findings align with prior research, confirming that while Grad-CAM excels in class-discriminative localization, IG provides a deeper, more faithful explanation of how models process image features.

6. Conclusion and Future Work

This study set out to evaluate and compare two popular explainability techniques—Grad-CAM and Integrated Gradients (IG)—for CNN-based image classification. Our experiments on ResNet-50 and DenseNet-121 revealed clear trade-offs between these methods. Grad-CAM, with its ability to generate smooth, high-level heatmaps, demonstrated superior localization performance, particularly in terms of IoU and recall. In contrast, IG produced more fine-grained, pixel-level attributions that were generally more faithful to the model's decision-making, as evidenced by stronger deletion test results. However, IG's attributions required significant post-processing (Gaussian filtering, adaptive thresholding) to reduce noise and improve spatial coherence.

Moreover, our TCAV analysis provided further insight into the alignment of model predictions with human-

interpretable concepts. DenseNet-121 consistently achieved high TCAV scores across all evaluated concepts (Stripes, Fur Texture, and Face Shape), suggesting that its feature representations are more closely aligned with human perceptual cues than those of ResNet-50. These findings underscore that no single method is universally optimal: Grad-CAM is preferable when intuitive visual localization is the primary goal, while IG is more suitable for applications requiring a faithful, detailed account of the decision process.

Future work should focus on several promising directions. One avenue is the development of hybrid methods that combine the strengths of Grad-CAM’s localization with the fidelity of IG. Additionally, further refinement of post-processing techniques for IG may enhance its spatial coherence and overall interpretability. Expanding this evaluation to newer architectures—such as Vision Transformers—and incorporating user studies to gauge the practical utility of these explanations would provide deeper insights. Finally, exploring more sophisticated concept-based analysis could lead to explainability frameworks that not only reflect low-level feature importance but also capture complex interactions between multiple high-level concepts.

In summary, our work provides a structured comparative evaluation of Grad-CAM and IG, highlighting their respective advantages and limitations, and laying the groundwork for more robust and integrated explainability solutions in deep learning.

References

- [1] **M. Sundararajan, A. Taly, and Q. Yan.** Axiomatic attribution for deep networks. *In International Conference on Machine Learning (ICML)*, pages 3319–3328, 2017.
- [2] **S. Lundberg and S.-I. Lee.** A unified approach to interpreting model predictions. *In Neural Information Processing Systems (NeurIPS)*, pages 4765–4774, 2017.
- [3] **M. Ribeiro, S. Singh, and C. Guestrin.** Why should I trust you? Explaining the predictions of any classifier. *In ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 1135–1144, 2016.
- [4] **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 IEEE International Conference on Computer Vision (ICCV)*, pages 618–626, 2017.
- [5] **D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg.** SmoothGrad: Removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- [6] **B. Kim, A. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viégas, and M. Wattenberg.** Interpretability beyond feature attribution: Testing with concept activation vectors (TCAV). *In International Conference on Machine Learning (ICML)*, pages 2668–2677, 2018.
- [7] **D. Kenny, M. Wu, and S. Singh.** Exemplar-based explanations in machine learning. *Journal of Artificial Intelligence Research*, 65(3):1123–1145, 2023.
- [8] **S. Ghafoorian, C. Liu, M. Havaei, C. Guizard, N. Rohde, and A. Vincent.** Deep neural networks for fast segmentation of brain lesions. *IEEE Transactions on Medical Imaging*, 38(2):1054–1066, 2018.
- [9] **A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, and T. Brox.** CARLA: An open urban driving simulator. *Conference on Robot Learning (CoRL)*, pages 1–16, 2017.
- [10] **N. Carlini and D. Wagner.** Adversarial examples are not easily detected: Bypassing ten detection methods. *In Workshop on Artificial Intelligence and Security (AISEC)*, pages 3–14, 2017.
- [11] **A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun.** Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639):115–118, 2017.
- [12] **J. Buolamwini and T. Gebru.** Gender shades: Intersectional accuracy disparities in commercial gender classification. *In Conference on Fairness, Accountability, and Transparency (FAT)*, pages 77–91, 2018.
- [13] **I. J. Goodfellow, J. Shlens, and C. Szegedy.** Explaining and harnessing adversarial examples. *In International Conference on Learning Representations (ICLR)*, pages 1–10, 2015.
- [14] **Y. Kirillov, H. Misra, T. He, P. Dollár, and R. Girshick.** Segment anything. *arXiv preprint arXiv:2303.00652v2*, 2023.
- [15] **XAIDataset Team.** XAIDataset: A benchmark repository for explainable AI. *XAI Benchmark Repository*, 2023. Available: https://xaidataset.github.io/quick_start/
- [16] **P. Langley.** Exemplar-based learning in AI systems. *AI Review Journal*, 35(4):912–929, 2022.
- [17] **D. Miller, F. Wong, and T. Zhang.** Human-centered AI explanations: A survey on user expectations and trust. *In AAAI Conference on Human-Centered AI*, pages 201–215, 2023.
- [18] **W. Samek, T. Wiegand, and K. Müller.** Evaluating the visualization of CNNs for decision support in medical applications. *IEEE Transactions on Neural Networks and Learning Systems*, 30(2):646–657, 2019.
- [19] **J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim.** Sanity checks for saliency maps. *In Neural Information Processing Systems (NeurIPS)*, pages 9505–9515, 2018.
- [20] **S. Amershi, J. Cakmak, W. Bender, M. C. Weld, and T. Kulesza.** Trust and interpretability in AI systems: A user study. *In ACM Conference on Human Factors in Computing Systems (CHI)*, pages 332–346, 2021.
- [21] **A. Kapoor, R. Patel, and J. Wexler.** Integrated Grad-CAM: Sensitivity-aware visual explanation of deep convolutional networks. *In ACM SIGGRAPH*, pages 412–420, 2022.
- [22] **R. Patel, A. Kapoor, and B. Kim.** Visual-TCAV: Bridging concept-based and saliency-based explainability. *arXiv preprint arXiv:2305.08598*, 2023.
- [23] **Y. Zhang, T. Lin, and H. Xu.** Robust explainability: Evaluating perturbation resilience in saliency maps. *In European Conference on Computer Vision (ECCV)*, pages 132–147, 2021.