

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

# Appendix

## Table of Contents

---

<b>A</b>	<b>Notation</b>	<b>15</b>
<b>B</b>	<b>Algorithm Details</b>	<b>15</b>
<b>C</b>	<b>Additional Experiments</b>	<b>16</b>
C.1	Qualitative Evaluation on the VGG16 with ImageNet . . . . .	16
C.2	<a href="#">Additional Experiment of Motivation Figure 1 on TwoMoons Dataset</a> . . . . .	19
C.3	Qualitative and Quantitative Evaluation on the IMDB Dataset . . . . .	19
C.4	<a href="#">Additional Quantitative Evaluation on ConvNeXt</a> . . . . .	21
C.5	<a href="#">Spurious Correlation Analysis via DCAPM Attribution</a> . . . . .	22
C.6	<a href="#">Entropy Analysis on Image Domain</a> . . . . .	23
C.7	Quantitative Evaluation on the Flowers and Oxford-Pet Datasets . . . . .	24
C.8	Optimization Progress Visualization . . . . .	24
C.9	Empirical Verification of Gradient Properties in ReLU Networks . . . . .	25
C.10	Empirical Validation of Spurious “ <i>Dummy Interactions</i> ” . . . . .	26
C.11	Validation of the Proposed Dummy Consistency Metric . . . . .	27
C.12	Noise Robustness Test . . . . .	28
<b>D</b>	<b>Experimental Details</b>	<b>29</b>
D.1	Datasets and Models . . . . .	29
D.2	Benchmark Methods and Implementations . . . . .	30
D.3	Hyperparameters for Attribution Computation . . . . .	31
D.4	Experimental Details for Figure 1 . . . . .	31
D.5	Non-Sensitivity Metric and Its Modification . . . . .	31
<b>E</b>	<b>Ablation Study</b>	<b>32</b>
E.1	Ablation Study: The Impact of Nullifying Dummy Interactions . . . . .	32
E.2	Ablation Study on Hyperparameters . . . . .	32
E.3	Ablation Study on the Cyclical Learning Rate . . . . .	33
<b>F</b>	<b>Implementation Code</b>	<b>34</b>
<b>G</b>	<b>Computational Efficiency and Scalability</b>	<b>36</b>
G.1	Computational Complexity . . . . .	36
G.2	Computational Resources and Runtimes . . . . .	36
<b>H</b>	<b>Mathematical Definitions Related to the Shapley Value</b>	<b>37</b>
H.1	The Assumptions of Aumann-Shapley Value . . . . .	37
H.2	The Axioms of Shapley-Shubik Value and Aumann-Shapley Value . . . . .	37
<b>I</b>	<b>Proofs of Theoretical Properties</b>	<b>38</b>
I.1	Proof of Efficiency . . . . .	38
I.2	Proof of Dummy/Null Player Axiom . . . . .	39
I.3	Proof of the Symmetry Axiom . . . . .	40
I.4	Proof of the Additivity Axiom . . . . .	40
I.5	Proof of the Positivity Axiom . . . . .	41
I.6	Proof of Dummy Consistency . . . . .	42
<b>J</b>	<b>Derivation of the DCAPM Attribution formula</b>	<b>43</b>

756	J.1	Derivation of the Derivative of the Path Function in Eq. 3 . . . . .	44
757	J.2	Derivation of the Dummy Consistent Attribution Path Method . . . . .	45
758			
759	<b>K</b>	<b>Concluding Remarks</b>	<b>45</b>
760	K.1	Limitation and Future Directions . . . . .	45
761	K.2	Large Language Model Usage Statement . . . . .	45
762			

## A NOTATION

Table 2 summarizes the mathematical notations used throughout this paper, defining key symbols for our attribution framework and optimization-based path method.

Table 2: **Table of notation.**

Notation	Description
$N$	The set of all features
$N \setminus \{i\}$	The set of all features except for feature $i$
$\mathcal{X}(N)$	The space of all possible inputs with $N$ features
$\mathbf{x} \in \mathbb{R}^{ N }$	An $ N $ -dimensional input vector
$x_i$	The value of the $i$ -th feature of the input vector $\mathbf{x}$
$\mathbf{x}_{-i}$	The input vector $\mathbf{x}$ with the $i$ -th feature excluded
$f(\mathbf{x}_{-i})$	The function value of the remaining features $\mathbf{x}_{-i}$ with respect to the dummy feature $i$
$f^{R_i^x}$	The reduced function on the given function $f$ with respect to the remaining features $\mathbf{x}_{-i}$
$\mathbf{x}'$	A reference point (the pseudo end point of a path)
$\mathbf{x}_{\text{dum}}$	An estimated baseline point (the end point of the path) composed of dummy features
$\mathcal{A}^\gamma(\mathbf{x}, \mathbf{x}'; f)$	An attribution method of feature $\mathbf{x}$ to the prediction of model $f$ with a path function $\gamma$
$\mathcal{A}_i^\gamma(\mathbf{x}, \mathbf{x}'; f)$	The $i$ -th attribution value of feature $\mathbf{x}$ to the prediction of model
$\gamma(t)$	The path function at step $t$ between input $\mathbf{x}$ and reference $\mathbf{x}'$
$t, \hat{t}$	The discrete time step of the path integration and its final step.
$\mathbf{m}^{(t)}$	The mask parameter at optimization step $t$ , typically with smaller dimensions than the input
$u(\cdot)$	The upsampling function that maps the low-resolution mask $\mathbf{m}$ to the input resolution
$M^{(t)}$	The upsampled mask $u(\mathbf{m}^{(t)})$
$\widetilde{M}_{\text{dum}}$	The estimated set of dummy features
$\mathcal{L}$	The objective function for dummy feature selection
$\mathcal{W}$	The Jacobian matrix of the upsampling function $u(\cdot)$
$\lambda$	The scaling factor for the L1 regularization term
$\tau_g$	<a href="#">The threshold for identifying dummy features and dummy interactions</a>
$\tau_{\text{dum}}$	<a href="#">The threshold for identifying dummy feature set</a>
$\eta$	The learning rate (w/o cyclical learning rate)
$\eta^{(t)}$	The cyclical learning rate at iteration $t$
$C^{(t)}$	The active set of non-dummy features at iteration $t$
$\mathcal{I}_{i,j}$	The Shapley interaction value between features $i$ and $j$

## B ALGORITHM DETAILS

The proposed method is detailed in Algorithm 1. For notational simplicity, we use  $\nabla_{\mathbf{m}^{(t)}} \mathcal{L}$  to denote the gradient  $d\mathcal{L}(\mathbf{x}, \mathbf{x}', \mathbf{m}^{(t)})/d\mathbf{m}^{(t)}$ . A clipping operation is applied after each update to constrain the mask values within the valid interval  $[0, 1]$ . The optimization process halts when any of the following

810 three termination criteria are met: **(1) Non-Negativity Condition:** The condition terminates the  
811 process if the function output  $f(\mathbf{x}^{(t)})$  is no longer positive (Line 6), ensuring the path has reached a  
812 region of non-positive contribution. **(2) Mask Convergence (Saturation):** The condition checks  
813 for the convergence of the mask itself. It halts the optimization if the change between consecutive  
814 mask updates falls below a predefined threshold, using a standard check that accounts for both  
815 absolute and relative tolerance (Line 19). **(3)  $\epsilon$ -Monotonicity Safeguard:** The third criterion is a  
816 safeguard that enforces  $\epsilon$ -monotonicity along the path. It terminates the process if a step unexpectedly  
817 increases the function value by more than  $\epsilon$  (Line 20), preventing the accumulation of attribution  
818 from non-monotonic segments.

---

### Algorithm 1 Dummy Consistent Attribution Path Method

---

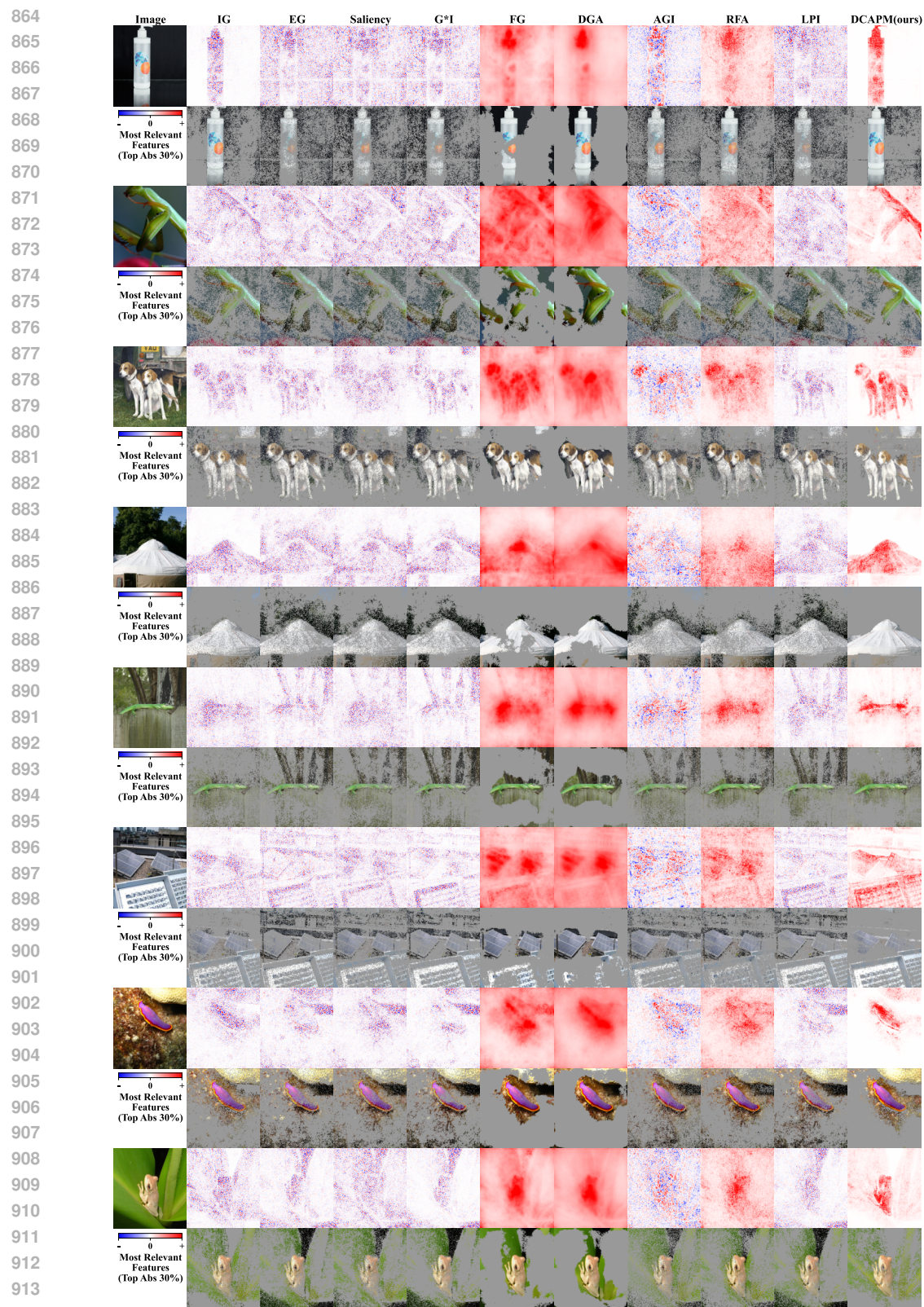
819 **1: Input:** Input  $\mathbf{x}$   
820 **2: Output:** Attribution  $\mathcal{A}(\mathbf{x})$   
821 **3: Parameter:** Mask  $\mathbf{m}$   
822 **4: Definitions:** Global reference point  $\mathbf{x}'$ ; **dummy feature threshold**  $\tau_g$ ; small value  $\epsilon$ ; the base  
823 learning rate  $\eta_{\min}$ ; the maximum learning rate  $\eta_{\max}$ ; the number of step size  $T_{\text{cycle}}$ , and the initial  
824 learning rate  $\eta^0$ ; path state and iteration step  $t$ ; maximum number of iterations  $T_{\text{fin}}$ ; upsampling  
825 function  $u : \mathbf{m} \in [0, 1]^{h \times w} \mapsto \mathcal{W}\mathbf{m} \in [0, 1]^{H \times W}$ ; bilinear interpolation weight matrix  $\mathcal{W}$ .  
826 **5: Initialize:**  $\mathbf{m}^{(0)} \leftarrow \mathbf{1}^{h \times w}$ ;  $x'_i \leftarrow \min_{\mathbf{x} \in D} x_i$ ;  $\mathcal{A} \leftarrow \mathbf{0}^{H \times W}$ ;  $C \leftarrow \emptyset$ ;  $\tau_g \leftarrow 0.1$ ,  $\epsilon \leftarrow 10^{-4}$ ;  
827  $\gamma(0) \leftarrow \mathbf{x} \odot u(\mathbf{m}^{(0)}) + \mathbf{x}' \odot (1 - u(\mathbf{m}^{(0)}))$ ;  $\nabla_{\mathbf{m}^{(0)}} \mathcal{L} \leftarrow \frac{d\mathcal{L}(\mathbf{x}, \mathbf{x}', \mathbf{m}^{(0)})}{d\mathbf{m}^{(0)}}$ ;  $t \leftarrow 0$ ;  $T_{\text{fin}} \leftarrow 600$ ;  
828  $\eta_{\min} \leftarrow 10^{-3}$ ;  $\eta_{\max} \leftarrow 0.01$ ;  $T_{\text{cycle}} \leftarrow 10$ ;  $\eta^{(0)} \leftarrow 10^{-3}$ .  
829 **6: while**  $f(\gamma(t)) > 0$  **and**  $t < T_{\text{fin}}$  **do** ▷ Non-negativity condition  
830  $\nabla_{\mathbf{m}^{(t)}} \mathcal{L} \leftarrow \frac{d\mathcal{L}(\mathbf{x}, \mathbf{x}', \mathbf{m}^{(t)})}{d\mathbf{m}^{(t)}}$   
831 **7:**  $q_\tau = \text{sort}(|\nabla_{\mathbf{m}^{(t)}} \mathcal{L}|)[\lceil \tau_g \cdot (h \times w - 1) \rceil]$  ▷  $\tau_g$ -quantile value (ascending order)  
832 **8:** **if**  $t > 1$  **then**  
833  $C = \{i \mid |\nabla_{m_i^{(t-1)}} \mathcal{L}| > q_\tau \text{ and } \nabla_{m_i^{(t)}} \mathcal{L} > \tau_g\}$  ▷ Masking dummy interactions  
834 **9:** **else**  
835  $C = \{i \mid \nabla_{m_i^{(t)}} \mathcal{L} > \tau_g\}$   
836 **10:** **end if**  
837 **11:** **for**  $i$  **in**  $C$  **do**  
838  $m_i^{(t+1)} = m_i^{(t)} - \eta^{(t)} \cdot \text{sign}(\nabla_{m_i^{(t)}} \mathcal{L})$  ▷ Dummy mask optimization  
839  $m_i^{(t+1)} = \text{Clip}(m_i^{(t+1)}, 0, 1)$   
840 **12:** **end for**  
841  $\gamma(t+1) \leftarrow \mathbf{x} \odot u(\mathbf{m}^{(t+1)}) + \mathbf{x}' \odot (1 - u(\mathbf{m}^{(t+1)}))$   
842 **13:** **if**  $|m_i^{(t)} - m_i^{(t+1)}| \leq \epsilon \cdot (1 + |m_i^{(t+1)}|) \quad \forall i$  ▷ Trivial difference condition  
843 **14:** **or**  $f(\gamma(t)) - f(\gamma(t+1)) < \epsilon$  **then** ▷  $\epsilon$ -Monotonicity condition  
844 **15:** **Terminate**  
845 **16:** **end if**  
846  $\mathcal{A}_i += (\gamma_i(t) - \gamma_i(t+1)) \cdot \frac{\partial f(\gamma(t+1))}{\partial \gamma_i(t+1)} \cdot \mathcal{W}_i \cdot \eta^{(t)} \cdot \nabla_{\mathbf{m}^{(t+1)}} \mathcal{L}$  ▷ Aggregate attributions  
847  $\eta^{(t+1)} = \eta_{\min} + (\eta_{\max} - \eta_{\min}) \cdot \max(0, 1 - \lfloor \frac{t - T_{\text{cycle}}}{T_{\text{cycle}}} \rfloor)$  ▷ Cyclical learning rate  
848  $t += 1$   
849 **17: end while**  
850 **18: return**  $\mathcal{A} \leftarrow \mathcal{A}/t$

---

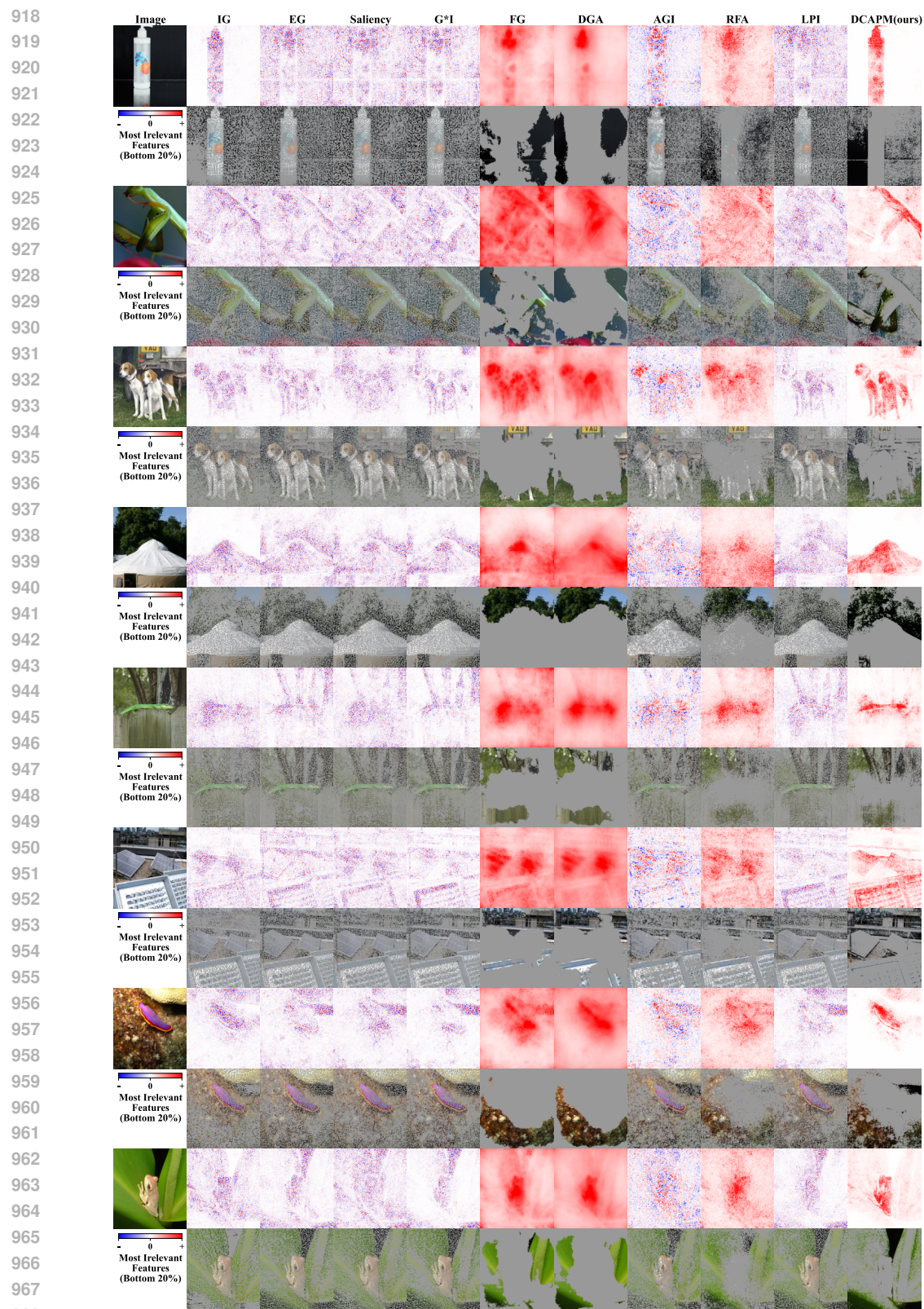
## C ADDITIONAL EXPERIMENTS

### C.1 QUALITATIVE EVALUATION ON THE VGG16 WITH IMAGENET

861 A qualitative evaluation on the ImageNet validation set visually demonstrates our method’s superior  
862 ability to distinguish between signal and noise. As shown in Figures 5 and 6, our approach produces  
863 sparse attribution maps. This stands in contrast to existing axiomatic methods, which yield diffuse  
and noisy results, and provides a cleaner separation of relevant and irrelevant features.

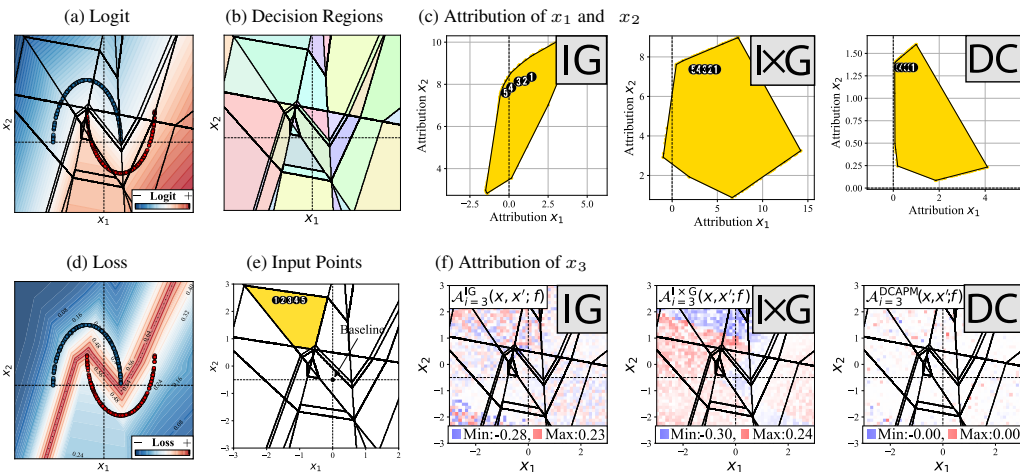


915 **Figure 5: Qualitative comparison of the most salient features (top 30%).** The top row for each  
 916 method shows the full attribution map, while the bottom row highlights only the top 30% of absolute  
 917 attribution values. Our method's top attributions are sharply focused on the foreground object,  
 whereas competing methods include significant background noise.



918  
 919  
 920  
 921  
 922  
 923  
 924  
 925  
 926  
 927  
 928  
 929  
 930  
 931  
 932  
 933  
 934  
 935  
 936  
 937  
 938  
 939  
 940  
 941  
 942  
 943  
 944  
 945  
 946  
 947  
 948  
 949  
 950  
 951  
 952  
 953  
 954  
 955  
 956  
 957  
 958  
 959  
 960  
 961  
 962  
 963  
 964  
 965  
 966  
 967  
 968  
 969 **Figure 6: Qualitative comparison of the least salient features (bottom 20%).** This visualization  
 970 highlights the features identified as most irrelevant (i.e., potential dummy features). Our method  
 971 correctly identifies only background regions as irrelevant, while competing methods erroneously  
 include noise or even parts of the foreground object in their bottom attributions.

972 C.2 ADDITIONAL EXPERIMENT OF MOTIVATION FIGURE 1 ON TWOMOONS DATASET



973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

Figure 7: **Additional Experiment of Figure 1 on TwoMoons.** A noise feature  $x_3$  is added to the 2D binary classification task. **(a, b, d, e) Model Behavior:** Logit/loss outputs and decision boundaries for signal features  $x_1$  and  $x_2$ . Points ①–⑤ (yellow region) correspond to a single-feature-axis-shift scenario. **(c) Dummy Consistency Violation:** IG shows the attributions for  $x_1$  and  $x_2$  are distorted by  $x_3$ . However, DCAPM (Ours) and non-path-based I×G remain undistorted. **(f) Dummy Player Violation:** IG and I×G assign non-zero attribution to the irrelevant  $x_3$ , whereas DCAPM correctly assigns it almost zero attribution.

1000 C.3 QUALITATIVE AND QUANTITATIVE EVALUATION ON THE IMDB DATASET

1001 To assess the generalizability of our method beyond CNNs, we conducted a comprehensive evaluation  
1002 on Transformer-based architectures. We compare our method, DCAPM, against Integrated Gradients  
1003 (IG) on BERT models fine-tuned for sentiment classification on the IMDB movie review dataset.

1004 **Experimental Setup** For qualitative analysis, we visualize the L2 norm of attribution values, where  
1005 greener shades indicate higher input importance. To ensure a fair comparison, the visualized samples  
1006 (#8 and #18) were specifically chosen from the overlapping region of the two methods' entropy  
1007 distributions (Figure 4), where they exhibit nearly identical entropy values. For quantitative analysis,  
1008 we measured the average entropy of attributions over the first 500 samples of the IMDB test set.  
1009 Detailed model and hyperparameter configurations are provided in Appendix D.

1010  
1011 **Results and Analysis** As shown in Tables 3 and 4, the qualitative results reveal a clear distinction  
1012 between the two methods. Our method (DCAPM) produces more focused attributions, concentrating  
1013 on semantically relevant tokens for a 'Negative' prediction such as "pains", "horrible", and "blame,"  
1014 while effectively suppressing attribution on irrelevant tokens. In contrast, IG yields noisier attributions,  
1015 often assigning high importance to functionally meaningless tokens like "[CLS]", "[SEP]", and  
1016 commas. This qualitative finding is supported by our quantitative analysis. As shown in Figure 4,  
1017 DCAPM achieves a significantly lower average entropy (4.1293) compared to IG (5.7015). This  
1018 confirms that our method produces sparser and more concentrated attributions, while IG tends to  
1019 distribute attribution more diffusely across a wider range of tokens, even if the BERT architectures  
1020 consist of GELU activation function. Collectively, these results confirm the robustness and scalability  
1021 of our method across three key dimensions:

1022 **i) Robustness to Model Depth** Our method's effectiveness is not constrained by model depth. The  
1023 qualitative results demonstrate that DCAPM provides more semantically meaningful attributions than  
1024 IG on both the 12-layer BERT-Base and the deeper 24-layer BERT-Large models. This is consistent  
1025 with our findings on CNNs of varying depths (e.g., VGG16, ResNet18, and InceptionV3), confirming  
the stability of our approach.

1026 **ii) Generalization to Transformer Architectures** The experiments on BERT-Base and BERT-  
 1027 Large confirm that our method, initially evaluated on CNNs, generalizes effectively to Transformer-  
 1028 based architectures. In both cases, DCAPM was qualitatively superior to the IG baseline in identifying  
 1029 the input features most relevant to the model’s sentiment classification.  
 1030

1031 **iii) Scalability with Model Size** DCAPM scales effectively to models of substantially different  
 1032 sizes. The quality of our results on the 334M-parameter BERT-Large model is comparable to that on  
 1033 the 110M-parameter BERT-Base model. This indicates that our method’s performance advantage is  
 1034 maintained even as the model size and parameter count increase threefold.  
 1035

1036 Table 3: **Qualitative comparison of IG and our method (DCAPM) on a sample (#8) from the**  
 1037 **IMDb test set, using the BERT-Base sentiment classifier.** Attribution maps visualize the L2 norm  
 1038 of attribution values, where greener shades indicate higher input importance.  
 1039

Method	Prediction	Entropy	Word Importance
IG	Negative	5.9269	#CLS it actually pains me to say it , but this movie was horrible on every level . the blame does not lie entirely with van dam ##me as you can see he tried his best , but let ' s face it , he ' s almost fifty , how much more can you ask of him ? i find it so hard to believe that the same people who put together und ##is ##puted 2 ; arguably the best ( western ) martial arts movie in years , created this . everything from the plot , to the dial ##og , to the editing , to the overall acting was just horribly put together and in ma #SEP
DCAPM	Negative	5.3265	#CLS it actually pains me to say it , but this movie was horrible on every level . the blame does not lie entirely with van dam ##me as you can see he tried his best , but let ' s face it , he ' s almost fifty , how much more can you ask of him ? i find it so hard to believe that the same people who put together und ##is ##puted 2 ; arguably the best ( western ) martial arts movie in years , created this . everything from the plot , to the dial ##og , to the editing , to the overall acting was just horribly put together and in ma #SEP

1079 Legend: ■ Zero Importance ■ High Importance

1080 Table 4: **Qualitative comparison of IG and our method (DCAPM) on a sample (#18) from the**  
 1081 **IMDb test set, using the BERT-Large sentiment classifier.** Attribution maps visualize the L2 norm  
 1082 of attribution values, where greener shades indicate higher input importance.  
 1083

1084	Method	Prediction	Entropy	Word Importance
1085	IG	Negative	5.5144	#CLS i just finished watching this
1086				movie and am disappointed to say that
1087				i didn ' t enjoy it a bit . it is so
1088				slow slow and un ##int ##eres ##ting
1089				. this kid from harry potter plays a
1090				shy teenager with an rude mother , and
1091				then one day the rude mother tells the
1092				kid to find a job so that they could
1093				accommodate an old guy apparently having
1094				no place to live has started to live
1095				with his family and therefore the kid
1096				goes to work for a old lady . and
1097				this old lady who is living all alone
1098	teaches him about girls , driving #SEP			
1099				
1100	DCAPM	Negative	5.1440	#CLS i just finished watching this
1101				movie and am disappointed to say that
1102				i didn ' t enjoy it a bit . it is so
1103				slow slow and un ##int ##eres ##ting
1104				. this kid from harry potter plays a
1105				shy teenager with an rude mother , and
1106				then one day the rude mother tells the
1107				kid to find a job so that they could
1108				accommodate an old guy apparently having
1109				no place to live has started to live
1110				with his family and therefore the kid
1111				goes to work for a old lady . and
1112				this old lady who is living all alone
1113	teaches him about girls , driving #SEP			

1121 Legend: ■ Zero Importance ■ High Importance

1122  
1123  
1124  
1125  
1126 C.4 ADDITIONAL QUANTITATIVE EVALUATION ON CONVNEXT

1127  
1128 In this section, we quantitatively evaluate our proposed method, DCAPM, against baseline ap-  
 1129 proaches—DGA and Integrated Gradients (IG)—focusing on attribution quality (DiffID) and compu-  
 1130 tational cost. Experiments were conducted using the ConvNext-T model (Liu et al., 2022), a modern  
 1131 architecture pretrained on ImageNet-1K with 28M parameters. We test on 500 randomly selected  
 1132 samples from ImageNet with seed 42. DCAPM achieved the best attribution fidelity with a mean  
 1133 DiffID score of 0.4801, surpassing both DGA (0.4737) and IG (0.2306). DCAPM offered a superior  
 balance of quality and efficiency, maintaining a reasonable runtime.

Table 5: **Quantitative comparison with baseline methods.** The mean and standard deviation for DiffID scores (higher is better) and runtime in seconds are reported. DCAPM achieves the best DiffID performance with reasonable computational cost.

Method	DiffID ( $\uparrow$ )	Runtime (s) ( $\downarrow$ )
IG	0.2306 $\pm$ 0.1245	<b>0.97 <math>\pm</math> 0.41</b>
DGA	0.4737 $\pm$ 0.1989	253.33 $\pm$ 404.29
<b>DCAPM (Ours)</b>	<b>0.4801 <math>\pm</math> 0.2133</b>	134.27 $\pm$ 165.51

### C.5 SPURIOUS CORRELATION ANALYSIS VIA DCAPM ATTRIBUTION

In this section, we demonstrate the utility of DCAPM for model debugging. An effective attribution method should faithfully reveal spurious correlations learned by the model (Adebayo et al., 2020; Schramowski et al., 2020).

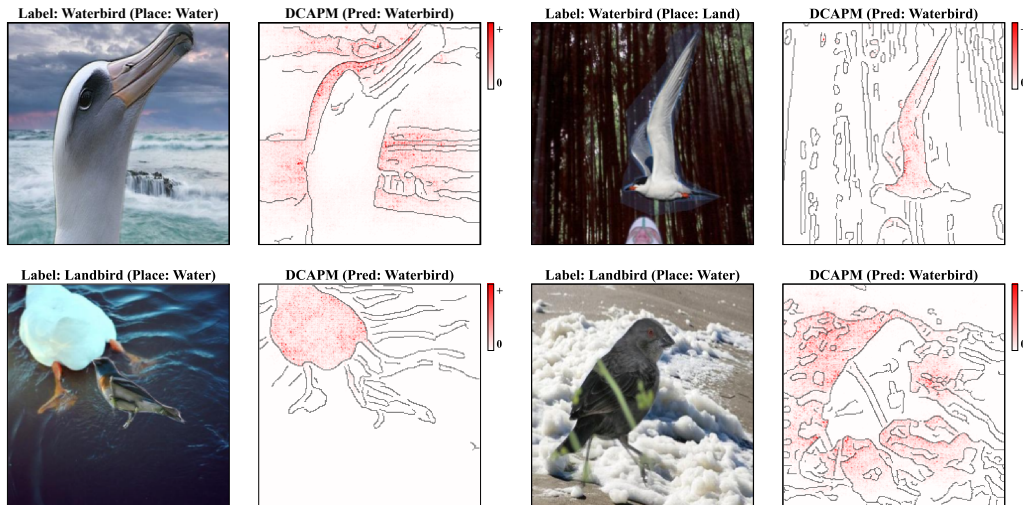


Figure 8: **Spurious correlation analysis via attribution maps.** (Top Left) **Aligned Correct:** The model highlights both the waterbird and the water, with a stronger focus on the water background, indicating a reliance on spurious correlations (lazy prediction). (Top Right) **Conflict Correct:** The model successfully overrides bias, exclusively highlighting the waterbird object without background noise. (Bottom) **Conflict Incorrect:** The model succumbs to bias by focusing on background features rather than the target object, such as a duck as a background object (Left) or the sea texture (Right), leading to misclassification.

**Experimental Setup** In this experiment, we conducted a spurious correlation analysis using the Waterbirds dataset<sup>3</sup>. This dataset is designed for a binary classification task to distinguish between Waterbirds and Landbirds. Each image contains background (Place) information, consisting of either Water or Land. Consequently, the entire dataset comprises four groups: Landbird on Land, Landbird on Water, Waterbird on Land, and Waterbird on Water. For model training, we utilized the ResNet50 architecture (He et al., 2016) and followed the fine-tuning protocol provided in the group\_DRO (Sagawa et al., 2020)<sup>4</sup>. The final test accuracy of the trained model is 77.22%.

**Hypothesis** We employ attribution methods to distinguish whether the model relies on spurious (place) or core (bird) features. Our hypotheses are: **(1) Aligned Group and Correct Prediction (Lazy Prediction):** The model likely exploits the background as an easy shortcut rather than focusing on the object. **(2) Conflict Group and Correct Prediction (Bias Override):** To predict correctly

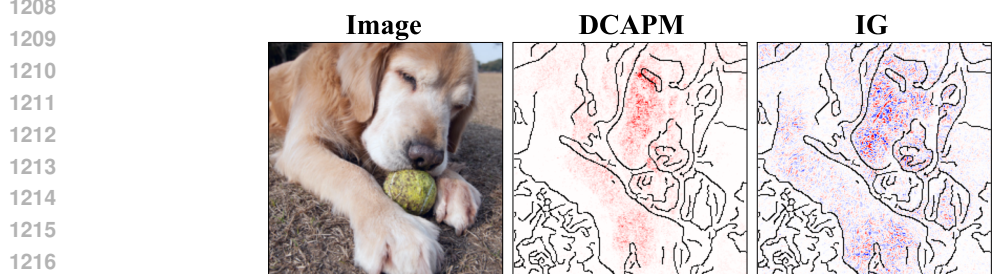
<sup>3</sup><https://huggingface.co/datasets/grodino/waterbirds>

<sup>4</sup>[https://github.com/kohpangwei/group\\_DRO/tree/master](https://github.com/kohpangwei/group_DRO/tree/master)

1188 against a conflicting background, the model suppresses spurious signals and focuses on the bird. (3)  
 1189 **Conflict Group and Incorrect Prediction (Succumb to Bias):** Misclassification results from the  
 1190 model focusing on the misleading background, confirming reliance on spurious correlations.  
 1191

1192 **Qualitative Analysis Results** To verify these hypotheses, we extracted and analyzed attribution  
 1193 maps for each case. As shown in Figure 8, the experimental results align with our hypotheses, and  
 1194 DCAPM successfully filters out noisy dummy features while clearly highlighting the spurious correla-  
 1195 tions learned by the model, demonstrating its utility for precise model debugging. In the **Conflict**  
 1196 **Correct** case (Top Left), we observe that the influence of the background (Land) is minimized, while  
 1197 the shape of the Waterbird object is clearly highlighted. Conversely, in the **Conflict Incorrect** cases  
 1198 (Bottom), high attribution values are distributed across the background region, visually demonstrat-  
 1199 ing that the model was misled by background information, resulting in misclassification. Finally, in the  
 1200 **Aligned Correct** case (Top Right), significant attribution occurs not only on the object but also on  
 1201 the background. This confirms that the model utilizes background information as a primary basis for  
 1202 prediction (spurious correlation) during the training process.

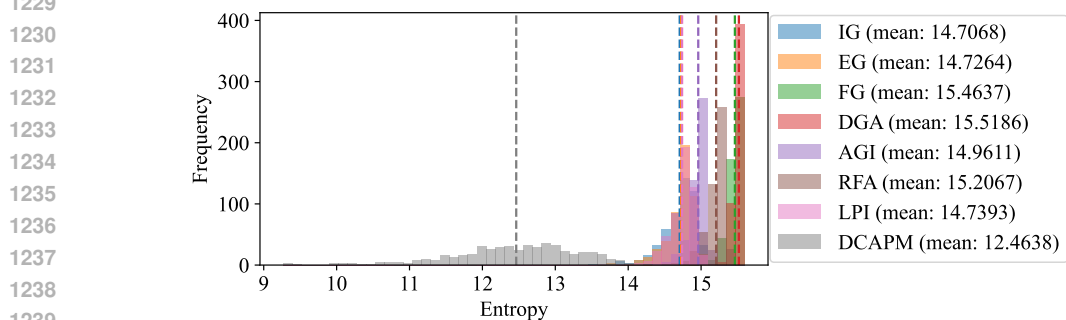
1203 **Additional Spurious Correlation Analysis Results** Additionally, spurious correlations can be  
 1204 seen in the imagenet dataset. As shown in Figure 9, DCAPM not only successfully eliminates noisy  
 1205 attributions caused by dummy features and spurious Shapley interactions but also clearly highlights  
 1206 the spurious correlation features relevant to the decision, proving its capability to distinguish actual  
 1207 learned biases from mathematical noise.



1218 Figure 9: **Spurious correlation analysis:** DCAPM reveals the spurious correlation of the ImageNet2012 training dataset, which frequently associates a ‘Golden Retriever’ image with tennis balls.  
 1219  
1220

1221 **C.6 ENTROPY ANALYSIS ON IMAGE DOMAIN**  
 1222

1223 In this section, we evaluate the sparsity of attribution maps by analyzing the entropy distribution  
 1224 across 500 randomly sampled images from the ImageNet validation set. As illustrated in Figure 10,  
 1225 all baseline methods exhibit relatively high entropy distributions, clustering around means ranging  
 1226 from 14.71 to 15.52. In contrast, DCAPM achieves a significantly lower mean entropy of 12.46. This  
 1227 distinct distribution shift confirms that our method yields the most sparse explanations, effectively  
 1228 concentrating attribution on relevant features while nullifying irrelevant noise.



1240 Figure 10: **Entropy analysis:** comparing baselines and DCAPM using VGG16 on 500 random  
 1241 samples from ImageNet2012 validation set.

## C.7 QUANTITATIVE EVALUATION ON THE FLOWERS AND OXFORD-PET DATASETS

To further validate the generalizability of our proposed method, we conducted a quantitative evaluation on the Flowers and Oxford-Pet datasets using a fine-tuned VGG16 model (see Appendix D for configuration details).

The results, summarized in Table 6, demonstrate that our method achieves superior performance in terms of robustness and reliability. Specifically, our approach obtains the best (lowest) scores for both Non-Sensitivity and Consistency across both datasets, significantly outperforming all baselines. For instance, on the Flowers dataset, our method achieves a Consistency score of 0.00027, an order of magnitude improvement over the next-best method, RFA (0.00232). While other methods such as FG and DGA excel on the DiffID metric, our approach remains highly competitive. These findings confirm that our method consistently generates more stable and faithful explanations across diverse data domains.

Table 6: **Quantitative evaluation of sensitivity and consistency on the Flowers and Oxford-Pet datasets.** Scores were computed using a fine-tuned VGG16 model on 500 randomly sampled instances from the validation set (fixed seed of 42). Higher is better for DiffID ( $\uparrow$ ), while lower is better for Non-Sensitivity ( $\downarrow$ ) and Consistency ( $\downarrow$ ). The **first** and **second**-best scores in each column are highlighted.

	DiffID ( $\uparrow$ )		Non-Sensitivity ( $\downarrow$ )		Dummy Consistency ( $\downarrow$ )	
	Flowers	OxfordPet	Flowers	OxfordPet	Flowers	OxfordPet
IG	0.014 $\pm$ 0.012	0.004 $\pm$ 0.003	0.00168 $\pm$ 0.00	0.00158 $\pm$ 0.00	0.20642 $\pm$ 0.05	0.19507 $\pm$ 0.05
EG	0.022 $\pm$ 0.019	0.005 $\pm$ 0.004	0.00167 $\pm$ 0.00	0.00160 $\pm$ 0.00	0.13257 $\pm$ 0.05	0.08331 $\pm$ 0.05
Saliency	0.011 $\pm$ 0.009	0.003 $\pm$ 0.002	0.00167 $\pm$ 0.00	0.00160 $\pm$ 0.00	0.08691 $\pm$ 0.05	0.04385 $\pm$ 0.04
I $\times$ G	0.018 $\pm$ 0.019	0.004 $\pm$ 0.003	0.00170 $\pm$ 0.00	0.00163 $\pm$ 0.00	0.07141 $\pm$ 0.06	0.02281 $\pm$ 0.03
FG	<b>0.061 <math>\pm</math> 0.046</b>	<b>0.027 <math>\pm</math> 0.022</b>	0.00131 $\pm$ 0.00	0.00094 $\pm$ 0.00	0.00913 $\pm$ 0.01	0.00430 $\pm$ 0.00
DGA	<b>0.061 <math>\pm</math> 0.044</b>	<b>0.028 <math>\pm</math> 0.024</b>	0.00138 $\pm$ 0.00	<b>0.00101 <math>\pm</math> 0.00</b>	0.07214 $\pm$ 0.03	0.03326 $\pm$ 0.02
AGI	0.023 $\pm$ 0.025	0.006 $\pm$ 0.005	0.00171 $\pm$ 0.00	0.00154 $\pm$ 0.00	0.17358 $\pm$ 0.06	0.16858 $\pm$ 0.06
RFA	0.044 $\pm$ 0.039	0.023 $\pm$ 0.021	0.00101 $\pm$ 0.00	<b>0.00071 <math>\pm</math> 0.00</b>	<b>0.00232 <math>\pm</math> 0.00</b>	<b>0.00196 <math>\pm</math> 0.00</b>
<b>DCAPM</b>	<b>0.067 <math>\pm</math> 0.044</b>	<b>0.015 <math>\pm</math> 0.022</b>	<b>0.00018 <math>\pm</math> 0.00</b>	<b>0.00017 <math>\pm</math> 0.00</b>	<b>0.00027 <math>\pm</math> 0.00</b>	<b>0.00031 <math>\pm</math> 0.00</b>

## C.8 OPTIMIZATION PROGRESS VISUALIZATION

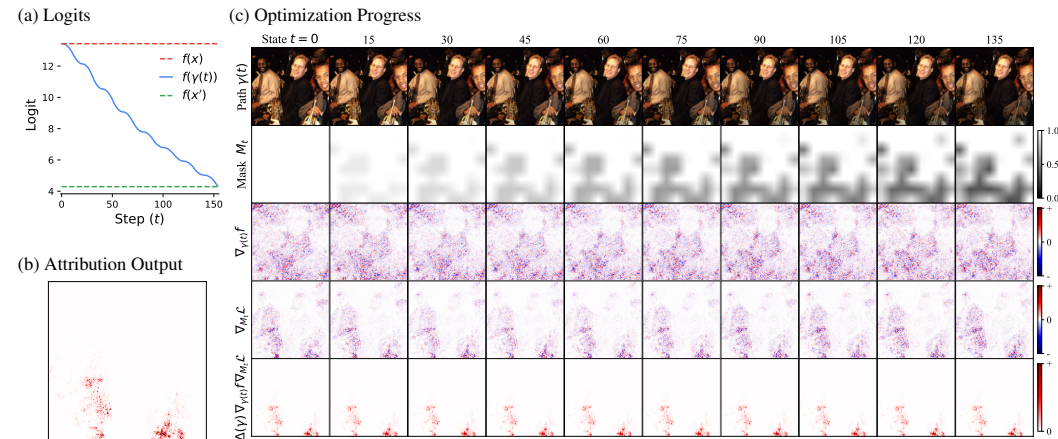


Figure 11: **The optimization process of target class ‘French Horn’.** As the mask is optimized, the dummy players corresponding to the target class approach 1, while the remaining players approach 0, leading the function value to converge toward  $f(x')$ . The target class is the predicted class at the first iteration.

C.9 EMPIRICAL VERIFICATION OF GRADIENT PROPERTIES IN RELU NETWORKS

To empirically verify our claim regarding *shattered gradients* and *rare zero gradients* in piecewise-linear networks, we conducted an experiment to analyze the gradient characteristics of a standard ReLU-based DNN.

**Experimental Setup** We used a pretrained VGG16 model and the COCO 2017 dataset (Lin et al., 2014), leveraging its provided object segmentation masks to isolate background pixels. Our experiment is based on the premise that for a given target object classification, pixels in the image background can serve as a practical proxy for irrelevant or *dummy* features. We then computed the gradients of the model’s classification output with respect to these background pixels. Our analysis focuses on two key aspects: (1) the distribution of these gradient values and (2) the sparsity of gradients across different layers of the network.

**Empirical Verification of “Rare Zero and Noisy Near-Zero Gradients”** Our premise was that the gradients in irrelevant background regions would not be exactly zero but would instead form a noisy, near-zero distribution. The statistics gathered from different layers of the VGG16 model, presented in Table 7, strongly support this claim.

As shown in Table 7, the mean of the background gradient distribution across all layers is extremely close to zero. However, the standard deviation is consistently non-zero. This empirically demonstrates that identifying irrelevant features by simply checking for ‘gradient == 0’ is infeasible in practice. Instead, the gradients form a *noisy near-zero* distribution, justifying the need for methods robust to such noise.

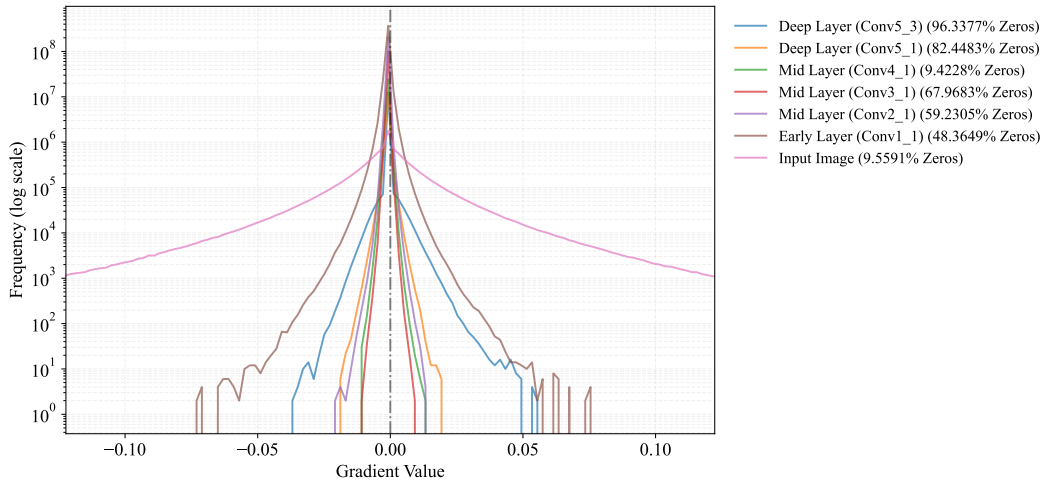


Figure 12: Distribution of masked gradients corresponding to background pixels in the input image, analyzed at different layers within a VGG16 model.

Table 7: Statistical properties of background gradients across different VGG16 layers.

Layer	Entropy	Mean (approx.)	Std. Dev. (approx.)	Min	Max
Deep (Conv5_3)	0.365	$2.88 \times 10^{-5}$	0.00104	-0.036	0.057
Deep (Conv5_1)	1.194	$1.28 \times 10^{-5}$	0.00063	-0.017	0.020
Mid (Conv4_1)	<b>2.680</b>	$-7.46 \times 10^{-6}$	0.00045	-0.010	0.015
Mid (Conv3_1)	1.233	$1.21 \times 10^{-6}$	0.00017	-0.009	0.009
Mid (Conv2_1)	1.172	$7.94 \times 10^{-7}$	0.00027	-0.019	0.014
Early (Conv1_1)	0.878	$5.18 \times 10^{-7}$	0.00087	-0.073	0.077
<b>Input Image</b>	<b>1.826</b>	<b><math>2.70 \times 10^{-6}</math></b>	<b>0.02367</b>	<b>-1.051</b>	<b>0.957</b>

1350 **Empirical Verification of “Shattered Gradients” via Sparsity Analysis** Our second premise was  
 1351 that gradients become less sparse and more *shattered* as they propagate backward to earlier layers  
 1352 (i.e., closer to the input). We verified this by measuring the percentage of exact zero gradients at  
 1353 different network depths, with results presented in Table 8.

1354 The results show a clear, general trend: the percentage of zero gradients is very high in the deep  
 1355 layers near the output, which is expected due to the pruning nature of ReLU activations. However,  
 1356 this sparsity dramatically *\*\*decreases in earlier layers\*\**, with the input layer gradients having only  
 1357  $\sim 9.6\%$  exact zeros. This demonstrates that as the gradient signal propagates backward, it becomes  
 1358 increasingly dense and *shattered*, affecting nearly all input features, including the irrelevant ones in  
 1359 the background.

1361 Table 8: **Sparsity of gradients across different VGG16 layers**, measured by the percentage of exact  
 1362 zeros.

Layer	Location	Percentage of Exact Zeros
Deep (Conv5_3)	Close to Output	96.34%
Deep (Conv5_1)	Close to Output	82.45%
Mid (Conv4_1)	Middle of Network	9.42%
Mid (Conv3_1)	Middle of Network	67.97%
Mid (Conv2_1)	Middle of Network	59.23%
Early (Conv1_1)	Close to Input	48.36%
<b>Input Image</b>	Input	<b>9.56%</b>

1374 In conclusion, these empirical results validate our claims regarding the nature of gradients in  
 1375 piecewise-linear networks. They provide a strong motivation for our work, which aims to develop a  
 1376 method that can reliably handle both *noisy near-zero* and *shattered* gradients.

### 1378 C.10 EMPIRICAL VALIDATION OF SPURIOUS “Dummy Interactions”

1380 To empirically validate our claim that feature interactions are captured when an attribution path  
 1381 traverses different linear regions, and that this can lead to non-zero attribution for dummy features,  
 1382 we conducted the following experiment.

1384 **Experimental Setup** We trained a three-layer Multi-Layer Perceptron (MLP) with a (4, 4, 2)  
 1385 architecture and two ReLU activation layers. The model was trained for binary classification on the  
 1386 TwoMoons synthetic dataset, on which it achieved 100% test accuracy. We analyzed the behavior  
 1387 of the gradient field along a straight-line path from an input point of  $[-1.0, 2.4, 0.2]$  to a baseline of  
 1388  $[0, -0.5, 0]$ . The input vector was augmented with a third dimension,  $x_3$ , which serves as a dummy  
 1389 feature independent of the classification task. Inspired by the approximation method for the Shapley  
 1390 interaction value from Chang et al. (2025), we then measured the change in the partial derivatives,  
 1391  $\Delta(\nabla f)$ , between adjacent discrete steps ( $t$  and  $t - 1$ ) along this path. The detailed method is as  
 1392 followed:

1393 **Methodology** Let  $f : \mathbb{R}^{N=3} \rightarrow \mathbb{R}$  be a differentiable function,  $\mathbf{x} \in \mathbb{R}^{N=3}$  be an input, and  
 1394  $\mathbf{x}' \in \mathbb{R}^{N=3}$  be a baseline. We define a straight-line path between them, discretized into  $T$  steps,  
 1395 where a point on the path at step  $t \in \{0, \dots, T\}$  is given by  $\mathbf{x}^t = \mathbf{x}' + \frac{t}{T}(\mathbf{x} - \mathbf{x}')$ .

1397 Inspired by the difference-of-gradients approach to approximating the Shapley interaction value from  
 1398 Chang et al. (2025), we define a metric to isolate the interaction between two features,  $i$  and  $j$ , during  
 1399 a single path step from  $t - 1$  to  $t$ . Unlike prior work that perturbs a feature to its baseline value, we  
 1400 measure the effect of the infinitesimal movement of feature  $j$  on the gradient of feature  $i$ .

1401 To achieve this, we define a counterfactual point  $\bar{\mathbf{x}}^{t-1}$  which is identical to  $\mathbf{x}^{t-1}$  except that feature  $j$   
 1402 has moved to its position at step  $t$ :

$$1403 \bar{\mathbf{x}}^{t-1} := (x_1^{t-1}, \dots, x_{j-1}^{t-1}, x_j^t, x_{j+1}^{t-1}, \dots, x_n^{t-1})$$

The interaction contribution of  $j$  on  $i$  during the step from  $t - 1$  to  $t$ , denoted as  $I_{i,j}^{(t)}$ , is then defined as the change in the partial derivative of  $f$  with respect to  $x_i$  between the original point  $\mathbf{x}^{t-1}$  and the counterfactual point  $\bar{\mathbf{x}}^{t-1}$ , scaled by the displacement of  $x_i$ :

$$I_{i,j}^{(t)} := \left( \frac{\partial f}{\partial x_i} \Big|_{\mathbf{x}=\mathbf{x}^{t-1}} - \frac{\partial f}{\partial x_i} \Big|_{\mathbf{x}=\bar{\mathbf{x}}^{t-1}} \right) \cdot (x_i^{t-1} - x_i^t)$$

Here, the term within the parenthesis isolates the change in the gradient of feature  $i$  caused solely by the movement of feature  $j$  in that step. Summing  $I_{i,j}^{(t)}$  over all steps  $t = 1, \dots, T$  would yield the total interaction between  $i$  and  $j$  along the path.

In this experiment, we omit the displacement term  $(x_i^{t-1} - x_i^t)$ . Because we use a straight-line path, this term simplifies to a constant scaling factor for each feature, which can be disregarded for the purpose of identifying non-zero interactions.

**Results** The change in partial derivatives was zero for most steps, indicating the path remained within a single linear region. However, at specific intervals where the path crossed linear region boundaries, we observed significant, non-zero changes. The results are denoted as  $\{t : [\Delta(\partial f/\partial x_1), \Delta(\partial f/\partial x_2), \Delta(\partial f/\partial x_3)]\}$ :

- Step 5: [5.6, -1.2, -0.0]
- Step 7: [0.1, 0.5, 0.1]

Critically, at Step 7, the traversal of a linear region boundary resulted in a non-zero change in the partial derivative with respect to the dummy feature ( $x_3$ ). This empirically demonstrates how spurious interactions, arising from the piecewise-linear structure of the network, can be assigned to irrelevant features. A path integral method, by aggregating these local changes, would consequently assign a non-zero attribution to the dummy feature. This experiment validates our claim that path-dependence in rectifier networks is a primary cause of axiomatic violations and motivates the need for a path-finding method robust to these dummy interactions.

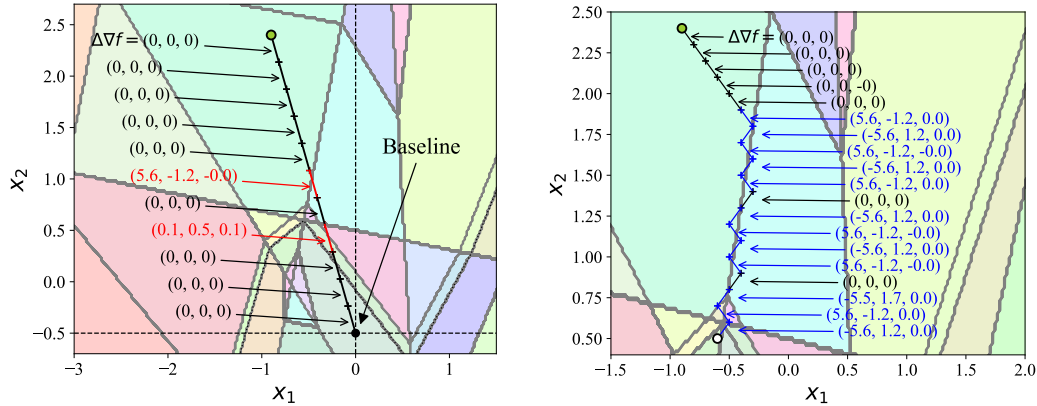


Figure 13: **Demonstration of a dummy interaction caused by traversing linear region boundaries.** A straight-line path (left) crossing the model’s decision boundaries (background heatmap) induces a spurious, non-zero interaction effect on the irrelevant dummy feature  $x_3$ . In contrast, DCAPM (right) finds a path that effectively reduces the spurious feature interaction on  $x_3$ .

### C.11 VALIDATION OF THE PROPOSED DUMMY CONSISTENCY METRIC

We validate our proposed metric through theoretical grounding and quantitative analysis. Theoretically, the metric is a direct operationalization of the Dummy Consistency Axiom. Quantitatively, we demonstrate that our metric captures a novel dimension of attribution quality, distinct from existing sensitivity measures. To do this, we performed a correlation analysis using the attribution results from ten methods on 500 ImageNet samples with a VGG16 model. As shown in Table 9, we report the Pearson and Spearman correlations between our Dummy Consistency metric and four established

1458 sensitivity metrics: Insertion/Deletion Games (Petsiuk et al., 2018), DiffID (Yang et al., 2023), and  
 1459 Non-Sensitivity (Nguyen & Martínez, 2020). The results show a statistically significant but low  
 1460 positive correlation between our metric and the others. This finding is critical: it empirically confirms  
 1461 that Dummy Consistency is a distinct and complementary property, not a redundant one. Therefore,  
 1462 improving consistency does not require a trade-off with sensitivity; rather, it contributes to a more  
 1463 holistically robust attribution.  
 1464

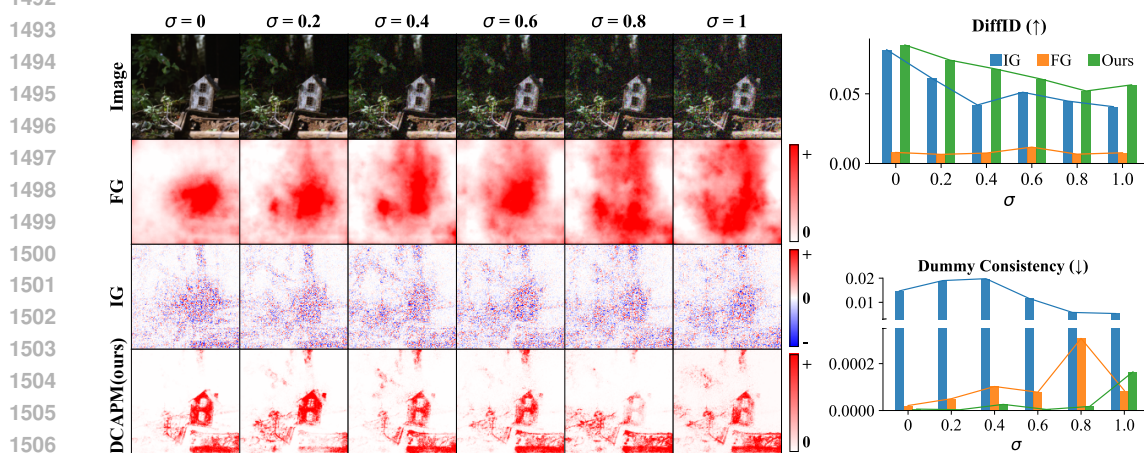
1465 **Table 9: Correlation between our Dummy Consistency metric and four established sensitivity**  
 1466 **metrics.** Pearson and Spearman coefficients are reported, with **p-values** in parentheses. Our metric is  
 1467 shown to be largely orthogonal to existing sensitivity measures. ‘1-DConsistency’ denotes the value  
 1468 of ‘1-(normalized Dummy Consistency score)’, used to align the interpretation direction (i.e., higher  
 1469 is better). (\*\*\*: p-value < 0.001).

	Insertion ( $\uparrow$ )	1-Deletion ( $\uparrow$ )	DiffID ( $\uparrow$ )	1-Nonsensitivity ( $\uparrow$ )
	<b>Pearson Correlation</b>			
	0.226 (***)	0.093 (***)	0.167 (***)	0.049 (***)
	<b>Spearman Rank Correlation</b>			
1-DConsistency ( $\uparrow$ )	0.273 (***)	0.107 (***)	0.322 (***)	0.019 (***)

### 1478 C.12 NOISE ROBUSTNESS TEST

1479  
 1480 Single-instance attribution methods such as Saliency Map (Simonyan et al., 2013) and FG (Srinivas  
 1481 & Fleuret, 2019) can be sensitive to input noise, potentially undermining the reliability of their  
 1482 explanations. As a path-based method, our approach is designed to offer greater robustness against  
 1483 such perturbations.

1484 Figure 14 presents a qualitative and quantitative evaluation of this robustness. We compare our  
 1485 method against IG and FG on a single ‘birdhouse’ sample, to which we incrementally add Gaussian  
 1486 noise with increasing standard deviation ( $\sigma$  from 0 to 1). On the DiffID sensitivity metric, our  
 1487 method demonstrates robustness comparable to IG, maintaining competitive performance even at  
 1488 high noise levels. The primary advantage of our approach, however, is observed in the dummy  
 1489 consistency metric. While IG and FG produce high MSE scores due to spurious attributions on the  
 1490 noisy background pixels, our method maintains a substantially lower MSE, highlighting its ability to  
 1491 consistently focus on true signal features.



1508  
 1509 **Figure 14: Noise robustness evaluation for the ‘Birdhouse’ class, comparing our method against**  
 1510 **IG and FG.** (a) Visual comparison of attribution maps under increasing levels of additive noise. (b)  
 1511 DiffID scores (sensitivity), where higher is better. (c) Dummy Consistency scores (MSE), where  
 lower is better.

## D EXPERIMENTAL DETAILS

This section provides the experimental details for this study, including the benchmark methods, datasets, and hyperparameters employed across all experiments.

### D.1 DATASETS AND MODELS

#### D.1.1 DATASETS

Our evaluation is conducted on four benchmark datasets: three for computer vision and one for natural language processing.

- **Flowers**<sup>5</sup>: A dataset of 4,317 object-centered images resized to  $3 \times 224 \times 224$ , containing five classes: *daisy*, *dandelion*, *rose*, *sunflower*, and *tulip*. The dataset was randomly split the data into training and validation sets using a 9:1 ratio, with a fixed random seed of 42 for reproducibility.
- **Oxford-IIIT Pet** (Parkhi et al., 2012): An object-centered dataset with 37 classes specifying species (*cat* or *dog*) and breed. It consists of 3,680 training and 3,669 testing images, all resized to  $3 \times 224 \times 224$ .
- **ImageNet-2012** (Deng et al., 2009): The full dataset comprising 1.28M training images and 50,000 testing images across 1,000 classes. All images were resized to  $3 \times 224 \times 224$  and normalized using a mean of (0.485, 0.456, 0.406) and a standard deviation of (0.229, 0.224, 0.225).
- **IMDb** (Maas et al., 2011): A sentiment analysis dataset containing 50,000 movie reviews, evenly split into training and testing sets with balanced positive and negative labels.

Table 10: **Dataset and model summary.** An asterisk (\*) denotes models that were fine-tuned and accuracy scores computed by us. Train/Test splits are indicated by a colon.

	Flowers	Oxford-IIIT Pet	ImageNet-2012	IMDb
Image Size	$3 \times 224 \times 224$	$3 \times 224 \times 224$	$3 \times 224 \times 224$	N/A (Text)
# of Labels	5	37	1,000	2
Dataset Size	4,317	3,680 : 3,669	1.28M : 50,000	25,000 : 25,000
Model (hidden dim.)	VGG16* (1024)	VGG16* (1024)	VGG16 (4096) InceptionV3 (2048) ResNet18 (512)	BERT-Base (768) BERT-Large*(1024)
Accuracy	0.9305*	0.9021*	V: 0.7159, I: 0.6954, R: 0.6976	Base: 0.8909 Large: 0.9464*

#### D.1.2 MODELS

We used the following pretrained models in our experiments:

- **Vision Models:** We employed standard **VGG16** (Simonyan & Zisserman, 2015), **InceptionV3** (Szegedy et al., 2016), and **ResNet18** (He et al., 2016) models. The classifier head dimensions were 4096 (VGG16), 2048 (InceptionV3), and 512 (ResNet18), respectively. For the Flowers and Oxford-IIIT Pet datasets, we fine-tuned VGG16 models with a reduced classifier dimension of 1024.
- **Transformer Models:** We evaluated two **BERT** (Devlin et al., 2019) models sourced from the transformers library. The first is **BERT-base** model (textattack/bert-base-uncased-imdb<sup>6</sup>), a 12-layer model with a 768-dimensional hidden state and 110 million parameters. The second is a fine-tuned **BERT-large** model (jigarcpatel/fine-tuned-bert-large-uncased-IMDb-dataset<sup>7</sup>), a 24-layer model with a 1024-dimensional hidden state and 336 million parameters. For the BERT-large experiments, the maximum sequence length was set to 256 tokens to ensure computational tractability.

<sup>5</sup><https://www.kaggle.com/datasets/alxmmaev/flowers-recognition>

<sup>6</sup><https://huggingface.co/textattack/bert-base-uncased-imdb>

<sup>7</sup><https://huggingface.co/jigarcpatel/fine-tuned-bert-large-uncased-IMDb-dataset>

- **Synthetic Data Models:** For the 2D data experiments shown in Figures 1 and 13 (**Circles** and **TwoMoons** (Buitinck et al., 2013)), we trained a simple MLP with two ReLU-activated hidden layers, as detailed in Table 11. The models yielded test accuracies of 1.000 and 0.9045 for the Circles and TwoMoons datasets, respectively.

Table 11: **MLP architecture for the TwoMoons and Circles datasets.**

Layer Type	Output Shape
Fully Connected + ReLU	4
Fully Connected + ReLU	4
Fully Connected	2

## D.2 BENCHMARK METHODS AND IMPLEMENTATIONS

To evaluate our approach, we compare it against a comprehensive set of benchmark methods chosen for their relevance to attribution consistency in Piecewise Linear Neural Networks (PLNNs). We group these methods into the following categories: **Gradient-based Methods:** FullGrad (FG) (Srinivas & Fleuret, 2019), Distilled Gradient Aggregation (DGA) (Jeon et al., 2022), Saliency Maps (Simonyan et al., 2013), and Gradient $\times$ Input(G\*I) (Shrikumar et al., 2017). **Path and Integration-based Methods:** Adversarial Gradient Integration (AGI) (Pan et al., 2021), Recalibrating Feature Attributions (RFA) (Yang et al., 2022), and Local Path Integration (LPI) (Yang et al., 2023). **Axiomatic Methods:** Integrated Gradients (IG) (Sundararajan et al., 2017) and Expected Gradients (EG). For all experiments, we used the official open-source implementations of each method with their default hyperparameter settings.

- **Integrated Gradients (IG)** (Sundararajan et al., 2017) is an axiomatic feature attribution technique grounded in Aumann-Shapley values. It approximates the Shapley value by integrating gradients along a straight-line path from a single baseline input to the original input.
- **Expected Gradients (EG)** (Erion et al., 2021) builds upon IG by reformulating the path integral as an expectation over multiple reference values sampled from a background distribution. This converges to attributions that sum to the difference between the model’s expected output and its current output for a given input.
- **FullGrad (FG)** (Srinivas & Fleuret, 2019) avoids path integration to mitigate counterintuitive attribution issue that can arise when paths cross linear region boundaries. It combines local input gradients with bias gradients from throughout the network, which are heuristically upsampled to the input dimension. By doing so, it sacrifices the efficiency axiom but compensates with bias terms. FG introduces the *weak dependence* property, which, unlike *input invariance*, allows two distinct inputs to have different attributions even if they reside in the same, potentially disconnected, linear region.
- **Distilled Gradient Aggregation (DGA)** (Jeon et al., 2022) is a variant of FG that iteratively refines attributions through a masking process. It employs two masks: one to filter out weakly contributing features (noise) and another to regularize features with excessively high contributions. The final attribution is an aggregation of these masked features, using FG as the base explainer.
- **Adversarial Gradient Integration (AGI)** (Pan et al., 2021) is a path-based method that finds both a path and multiple baselines adversarially. Instead of a straight line, AGI defines a non-linear path by following the gradient ascent direction on the model’s loss surface until the prediction for the target class is nullified. This process is stabilized using techniques from adversarial attacks, such as Projected Gradient Descent (PGD).
- **Recalibrating Feature Attributions (RFA)** (Yang et al., 2022) extends AGI by defining criteria for valid path interpolations for aggregating positive attributions. To achieve this, RFA posits that a valid baseline should represent the absence of the input signal, and the path from such a baseline should follow the gradient ascent direction in the model’s logit surface. This validity condition is applied on a feature-wise basis across multiple baselines sampled from the training data.
- **Local Path Integration (LPI)** (Yang et al., 2023) also aggregates gradients over multiple baselines, but constrains the selection of these baselines to the same linear region as the input to consider the

---

property of piecewise linear neural networks (Chu et al., 2018). However, its method for estimating this region—K-Means clustering in the output logit space—is an indirect heuristic. Identifying points that share the exact same linear function is practically infeasible due to the extremely high number of tiny linear regions (polytopes) in DNNs.

### D.3 HYPERPARAMETERS FOR ATTRIBUTION COMPUTATION

The hyperparameters for the attribution methods used in our experiments were set as follows.

**Vision Models (VGG16, ResNet18, InceptionV3)** For vision models, the dummy-mask optimization process, applied in both the attribution and evaluation phases of our method, was configured as follows. We set the mask resolution to  $h \times w = 8 \times 8$  and employed a triangular cyclical learning rate policy (Smith, 2017) with a base rate of 0.0001, a maximum rate of 0.01, and a cycle step size of 10. The L1 regularization factor ( $\lambda$ ) was set to 0.0001, and the threshold for identifying dummy interactions ( $\tau$ ) was 0.1. The optimization was run for a maximum of 600 steps, using a constant  $\epsilon$  of  $1 \times 10^{-4}$  within a termination condition for numerical stability while satisfying the  $\epsilon$ -monotonicity axiom.

**Language Models (BERT-Base, BERT-Large)** For experiments involving BERT models, we used a distinct set of hyperparameters for our method (DCAPM). The mask size was adapted to match the model’s hidden dimension: 768 for BERT-base and 1024 for BERT-large. The optimization was configured with a learning rate of 0.001, an L1 regularization factor ( $\lambda$ ) of 0.001, and a dummy interaction threshold ( $\tau$ ) of 0.1, running for 600 iterations with  $\epsilon = 1 \times 10^{-4}$ . For these language model experiments, both DCAPM and the IG benchmark used a zero-vector baseline.

### D.4 EXPERIMENTAL DETAILS FOR FIGURE 1

This experiment is designed to demonstrate the vulnerability of standard attribution methods to the introduction of a single, semantically irrelevant feature on the rectified DNNs. We show that both path-based methods (e.g., IG) and single-instance gradient-based methods (e.g., FG) produce distorted attributions, violating core axioms, Dummy Player and Dummy Consistency axioms, when presented with such an input.

**Experimental Setup** The experimental setup reproduces the 2D Circles classification task from (Jeon et al., 2022, Figure 4). We trained a three-layer Multi-Layer Perceptron (MLP) with two ReLU activation layers, which achieved 100% test accuracy on this task. To the original 2D input space, we introduce an irrelevant dummy feature axis ( $x_3$ ). Values for this dummy feature were sampled from a uniform distribution over the range of the original dataset features.

**Key Observation** The central finding, illustrated in Figure 1, is that the introduction of this single dummy feature causes both IG and FG to produce distorted and counter-intuitive attributions for the true signal features ( $x_1, x_2$ ). This demonstrates a clear violation of the Dummy Player and Dummy Consistency axioms and motivates the need for attribution methods that are robust to such axiomatic failures.

### D.5 NON-SENSITIVITY METRIC AND ITS MODIFICATION

The *Non-Sensitivity* metric, originally proposed by (Nguyen & Martínez, 2020), serves as a quantitative measure for the Dummy Player axiom, as it evaluates the alignment between a feature’s attribution and its functional relevance to the model’s output. In this work, we compare the Non-Sensitivity metric with our proposed metric. However, due to the unavailability of an official open-source implementation, we first reintroduce the original definition and then detail our own modified implementation.

**Definition 5** (Non-Sensitivity (Nguyen & Martínez, 2020)). *Let  $A_0 = \{i \in \{1, \dots, N\} \mid a_i = 0\}$  denote the set of feature indices assigned zero attribution. Further, let  $X_0 = \{i \in \{1, \dots, N\} \mid \mathbb{E}[\|f(x) - f(x_{-i})\|^2] = 0\}$ . be the set of indices corresponding to features on which the model  $f$  is*

not functionally dependent. Then, the Non-Sensitivity metric is defined as:

$$\text{NonSens}(a, f) = |A_0 \Delta X_0| \quad (12)$$

where  $|\cdot|$  denotes the cardinality of a set, and  $\Delta$  represents the symmetric difference (XOR) between the two sets.

The original metric penalizes an attribution method if a feature with zero attribution ( $i \in A_0$ ) is actually functionally relevant ( $i \notin X_0$ ), or vice-versa. However, this formulation presents practical challenges, as exact zero attributions and perfectly invariant outputs are rare in real-world DNNs. Furthermore, the raw score is not normalized, making comparisons across different feature spaces difficult. To address these limitations, we propose a more practical, normalized version of the Non-Sensitivity metric that incorporates a tolerance threshold and operates on feature partitions.

**Definition 6 (Modified Non-Sensitivity).** Let  $a_i$  denote the attribution scores sorted in ascending order by their absolute values, and  $x_i$  be the corresponding input features. We define a uniform partition  $\mathcal{P} = \{S_k\}_{k=1}^K$ , where each subset  $S_k = \{x_{(\mathcal{P}_k)}, \dots, x_{(\mathcal{P}_{k+1})}\} \subset N$  contains a consecutive segment of features from the ordered list. Given a tolerance threshold  $\epsilon > 0$ , we define dummy players with respect to the attribution and the model’s functional behavior as follows:

$$A_0 = \{S_k \subset N \mid \exists i \in S_k, \|a_i\| < \epsilon\}, \quad (13)$$

$$X_0 = \{S_k \subset N \mid \|f(x) - f(x_{-S_k})\| < \epsilon\}, \quad (14)$$

where  $f(x_{-S_k})$  denotes the function output when features in  $S_k$  are replaced by their global reference point values, keeping the remaining features unchanged. The modified Non-Sensitivity metric is then defined as:

$$\text{Modified Non-Sensitivity} = \frac{|A_0 \Delta X_0|}{|\mathcal{P}|}, \quad (15)$$

The metric is normalized by  $K = |\mathcal{P}|$ , the total number of partitions. A lower score on this metric signifies a better alignment between the attribution scores and the model’s functional dependencies.

## E ABLATION STUDY

### E.1 ABLATION STUDY: THE IMPACT OF NULLIFYING DUMMY INTERACTIONS

Our proposed method consists of two key optimization phases: (1) dummy feature selection and (2) dummy interaction nullification. To isolate the contribution of each phase, we perform an ablation study on VGG16 using 500 ImageNet samples. We first evaluate performance using only dummy feature selection and then measure the additional gains from incorporating the interaction nullification term. The results in Table 12 demonstrate that while feature selection alone provides a strong baseline, the explicit nullification of dummy interactions is crucial for achieving the best performance across both sensitivity and dummy consistency metrics.

Table 12: **Ablation study evaluating the impact of nullifying dummy interactions.** The metrics of Sensitivity and Dummy Consistency are reported on VGG16 using 500 samples from the ImageNet validation set.

Phase Configuration	DiffID(↑)	NonSens(↓)	DumCons(↓)
(1) Dummy Feature Selection Only	.0334 ±.05	.00190 ±.00	.005985 ±.03
<b>(1) + (2) With Dummy Interaction Nullification</b>	<b>.0362 ±.06</b>	<b>.00023 ±.00</b>	<b>.000398 ±.00</b>

### E.2 ABLATION STUDY ON HYPERPARAMETERS

The proposed method involves three key hyperparameters: the L1 regularization scaling factor,  $\lambda$ , the threshold for identifying dummy interactions,  $\tau$ , and mask size  $h, w$  for the optimization parameter. We conducted ablation studies to analyze their impact on performance.

**Analysis of Regularization Factor  $\lambda$**  Ablation results for  $\lambda$  are presented in Table 13. The results indicate that  $\lambda = 0.0001$  provides the best trade-off, achieving optimal sensitivity scores (DiffID and NonSens) without significantly compromising the dummy consistency score. While the regularization term plays a role in fine-tuning performance, its overall impact is modest. This suggests that the effectiveness of our method is primarily driven by its structural design—the dummy feature selection-based path-finding—rather than by hyperparameter optimization. Based on this, we selected  $\lambda = 0.0001$  for all experiments in this paper.

Table 13: **Ablation study on the scaling factor  $\lambda$  of the regularization term.** Sensitivity and Consistency scores are reported for our method with varying values of  $\lambda$ , evaluated using the VGG16 model and 500 samples from the ImageNet validation set.

$\lambda$	DiffID( $\uparrow$ )	Non-Sensitivity( $\downarrow$ )	Dummy Consistency( $\downarrow$ )
0	.0299 $\pm$ .05	.00056 $\pm$ .00	.000472 $\pm$ .00
0.1	.0334 $\pm$ .05	.00190 $\pm$ .00	<b>.000385 <math>\pm</math>.00</b>
<b>0.0001</b>	<b>.0362 <math>\pm</math>.06</b>	<b>.00023 <math>\pm</math>.00</b>	.000398 $\pm$ .00

**Analysis of Dummy Interaction Threshold  $\tau_g$**  For the threshold  $\tau_g$ , the results in Table 14 show that our method is largely robust to its value. The performance metrics exhibit minimal variation across the tested range, confirming the stability of our approach with respect to this parameter. Based on this finding, we selected  $\tau_g = 0.1$  for all main experiments.

Table 14: **Ablation study on the threshold  $\tau_g$  used to identify dummy interactions based on gradient magnitudes.** Sensitivity and Consistency scores are reported for varying values of  $\tau_g$ , evaluated using the VGG16 model and 500 samples from the ImageNet validation set.

$\tau_g$	DiffID( $\uparrow$ )	Non-Sensitivity( $\downarrow$ )	Dummy Consistency( $\downarrow$ )
0.001	.0359 $\pm$ .06	.00022 $\pm$ .00	.000474 $\pm$ .00
0.01	.0359 $\pm$ .06	.00022 $\pm$ .00	.000474 $\pm$ .00
<b>0.1</b>	<b>.0362 <math>\pm</math>.06</b>	.00023 $\pm$ .00	<b>.000398 <math>\pm</math>.00</b>
0.2	.0348 $\pm$ .06	<b>.00019 <math>\pm</math>.00</b>	.000465 $\pm$ .00

**Ablation Study on Mask Size** We conducted an ablation study to evaluate the impact of mask size ( $h = w \in \{4, 8, 16, 32\}$ ), holding all other hyperparameters constant as specified in Appendix D.3. The results, presented in Table 15, indicate that mask size is a sensitive hyperparameter, revealing a trade-off between the evaluation metrics. While larger mask sizes (16 and 32) yield optimal scores on individual metrics (DiffID and Dummy Consistency, respectively), a mask size of 8 provides the most effective balance. Specifically, the 8x8 mask achieves a substantially better Non-Sensitivity score than other configurations while maintaining highly competitive performance on the DiffID metric. Therefore, based on this trade-off, we selected a mask size of 8x8 for all main experiments.

Table 15: **Ablation study evaluating the impact of the mask size.** The metrics of Sensitivity and Consistency are reported on VGG16 using 500 samples from the ImageNet validation set.

Mask Size	DiffID( $\uparrow$ )	Non-Sensitivity( $\downarrow$ )	Dummy Consistency( $\downarrow$ )
4	.03397 $\pm$ .05	.00190 $\pm$ .00	.00025 $\pm$ .00
<b>8</b>	.03619 $\pm$ .06	<b>.00023 <math>\pm</math>.00</b>	.00040 $\pm$ .00
16	<b>.03647 <math>\pm</math>.04</b>	.00036 $\pm$ .00	.00019 $\pm$ .00
32	.03381 $\pm$ .07	.00046 $\pm$ .00	<b>.00017 <math>\pm</math>.00</b>

### E.3 ABLATION STUDY ON THE CYCLICAL LEARNING RATE

We conducted an ablation study to evaluate the effect of employing a cyclical learning rate, a technique intended to help the optimization process avoid poor local minima. The results, presented in Table 16,

1782 reveal a clear trade-off between our primary evaluation metrics. While the cyclical learning rate  
 1783 improved the Dummy Consistency score, this came at the cost of a decrease in the DiffID sensitivity  
 1784 score. Given that maintaining high sensitivity is a critical requirement for a reliable attribution  
 1785 method, we concluded that the standard learning rate policy offered a superior overall performance  
 1786 balance. Consequently, all main experiments presented in this paper were conducted without the  
 1787 cyclical learning rate policy.

1788  
 1789 **Table 16: Ablation study evaluating the impact of the Cyclical Learning Rate.** The metrics of  
 1790 sensitivity and dummy consistency are reported on the VGG16 using 500 samples from the ImageNet  
 1791 validation set.

Phase Configuration	DiffID( $\uparrow$ )	Non-Sensitivity( $\downarrow$ )	Dummy Consistency( $\downarrow$ )
<b>Without Cyclic LR</b>	<b>.0362 <math>\pm</math>.06</b>	<b>.00023 <math>\pm</math>.00</b>	<b>.000398 <math>\pm</math>.00</b>
With Cyclic LR	.0349 $\pm$ .06	.00023 $\pm$ .00	.000218 $\pm$ .00

## 1798 F IMPLEMENTATION CODE

1799 Below is the Python implementation of the our proposed algorithm, consistent attribution path  
 1800 method(DCAPM), and consistency metric.

```

1803 import torch
1804
1805 class DCAPM(object):
1806     def __init__(self,
1807                 model,                # (func) Model function to be explained
1808                 preprocess,          # (func) Predefined function to approximate reference point
1809                 im_size = 224,       # (int) Input image size
1810                 m_size = 8,          # (int) Mask optimization parameter size
1811                 max_iter = 600,      # (int) Maximum number of iteration
1812                 base_lr = 0.0001,    # (float) Base learning rate(LR) for cyclical LR
1813                 max_lr = 0.01,      # (float) Maximum learning rate for cyclical LR
1814                 step_size = 10,     # (int) Step size for cyclical LR
1815                 eps = 1e-4,          # (float) Infinitesimal value
1816                 tau = 0.1,           # (float) Threshold for dummy localization
1817                 lambda_reg = 1e-4,   # (float) Scaling factor for L1 regularization
1818                 ):
1819
1820         self.model, self.preprocess = model, preprocess
1821         self.im_size, self.m_size = im_size, m_size
1822         self.max_iter, self.eps, self.tau = max_iter, eps, tau
1823         self.base_lr, self.max_lr, self.step_size = base_lr, max_lr, step_size
1824
1825     def get_cyclic_lr(step, base_lr, max_lr, step_size):
1826         cycle = step // (2 * step_size)
1827         x = step % (2 * step_size)
1828         scale = 1.0 - abs(x - step_size) / step_size
1829         lr = base_lr + (max_lr - base_lr) * scale
1830         return torch.tensor(lr)
1831
1832     def get_mask(self, image):
1833         prev_masked_image = image.clone()
1834         input_size = (self.im_size, self.im_size)
1835         grid = torch.ones((1, 1, self.m_size, self.m_size), requires_grad=True)
1836         baseline = self.preprocess(torch.zeros_like(image).detach().cpu())
1837         masked_attrs = torch.zeros((1, 3, self.im_size, self.im_size))
1838
1839         scores, masks, images, masked_attrs, FLAG = [], [], [], [], False
1840         C = torch.ones_like(grid)
1841         for self.iteration in range(self.max_iter):
1842             lr = self.get_cyclic_lr(step=self.iteration, self.base_lr,
1843                                   self.max_lr, self.step_size)
  
```

---

```

1836
1837     # Resize the mask parameter to the original image size.
1838     mask = torch.nn.functional.interpolate(grid, input_size, mode='bilinear',
1839                                           align_corners=False)
1840     mask.data = torch.clamp(mask.data, 0, 1)
1841     masked_image = image * mask + baseline * (1 - mask)
1842
1843     # Obtain the logit value from the masked image.
1844     with torch.no_grad():
1845         if self.iteration == 0:
1846             original_score = self.model(image)
1847             target = torch.argmax(original_score, -1)
1848             masked_score = self.model(masked_image)[range(masked_image.shape[0]), target]
1849
1850     # Terminate based on predefined conditions.
1851     if torch.allclose(mask.reshape(-1).detach(),
1852                      prev_mask.reshape(-1), rtol=self.eps, atol=self.eps) or
1853        prev_score - masked_score.item() < eps or
1854        self.iteration == self.max_iter - 1 or masked_score < 0:
1855         break
1856
1857     reg_loss = lambda_reg * torch.norm(mask, 1)
1858     loss = masked_score - reg_loss
1859     loss.backward(retain_graph=True)
1860
1861     # Compute the difference of gradients
1862     with torch.no_grad():
1863         gradient = grid.grad
1864         threshold = torch.quantile(torch.abs(grad), self.tau, interpolation='lower')
1865         C_tmp = gradient > threshold
1866         grid = grid - gradient.sign() * C_tmp * C * lr
1867         grid.data = torch.clamp(grid.data, 0, 1)
1868         grid.requires_grad = True
1869         C = grad > threshold
1870
1871     # Compute the gradient of the loss function w.r.t. the mask.
1872     mask_grad = lr * torch.autograd.grad(
1873         outputs=loss, inputs=mask, retain_graph=True)[0].detach()
1874
1875     # Compute the gradient of the self.model w.r.t. the masked image.
1876     fn_grad = torch.autograd.grad(
1877         outputs=masked_score,
1878         inputs=masked_image,
1879         retain_graph=True)[0].detach()
1880
1881     delta = prev_masked_image - masked_image
1882     self.masked_attrs += fn_grad * mask_grad * delta
1883     prev_masked_image = masked_image.clone().detach()
1884     prev_mask = mask.clone().detach()
1885     prev_score = masked_score.item()
1886
1887     return mask
1888
1889     def get_attribution(self, image):
1890         _ = self.get_mask(image)
1891
1892         # Normalize by the number of iterations.
1893         attribution = self.masked_attrs.detach() / (self.iteration + 1)
1894
1895         return attribution
1896
1897     def consistency(
1898         orig_attribution, # (tensor) Original attribution map to be evaluated.
1899         image,           # (tensor) Original image to be evaluated.
1900         explainer,       # (func) Target explainer to be evaluated.
1901         minmax_scale,    # (func) Minmax scaling function.
1902         tau=0.1,        # (float) Threshold for dummy consistency.

```

---

```

1890 ):
1891     ''' Dummy players estimation-based dummy consistency measure. '''
1892
1893     _, C, H, W = image.shape
1894     HW = H * W
1895     step = (HW - 1) // (n_steps - 1)
1896
1897     # Obtain the approximated global reference point.
1898     substrate_fn = torch.zeros_like
1899     baseline = explainer.preprocess(substrate_fn(image))
1900
1901     # Get sorted attribution indices
1902     if orig_attribution.dim() == 4 and orig_attribution.shape[1] > 1:
1903         orig_attribution = orig_attribution.mean(dim=1)
1904         min_, max_ = orig_attribution.min(), orig_attribution.max()
1905         orig_attribution = minmax_scale(orig_attribution, min_, max_)
1906         threshold = torch.quantile(orig_attribution, tau, interpolation='lower')
1907         dummy_mask = orig_attribution > threshold
1908         masked_image = image * dummy_mask + baseline * (1 - dummy_mask)
1909
1910     # Compute reduced attribution
1911     reduced_attr = explainer.get_attribution(masked_image).mean(dim=1)
1912     reduced_attr = minmax_scale(reduced_attr, min_, max_)
1913     reduced_attr = reduced_attr.view(1, HW)
1914
1915     # Compute conditioned attribution
1916     condition_attr = orig_attribution.view(1, HW)
1917
1918     # Calculate MSE only on the dummy cases (abs attribution < tau).
1919     mask = torch.abs(condition_attr) < threshold
1920     valid_count = mask.sum(dim=-1)
1921     mse = torch.pow(condition_attr * mask - reduced_attr * mask, 2).sum(dim=-1) /
1922         torch.where(valid_count > 0, valid_count, 1)
1923     mse = torch.nan_to_num(mse, nan=0.0)
1924
1925     return mse.mean()

```

## 1922 G COMPUTATIONAL EFFICIENCY AND SCALABILITY

### 1925 G.1 COMPUTATIONAL COMPLEXITY

1927 In this section, we analyze the computational complexity of the proposed algorithm with respect to  
1928 the input feature size,  $N$ . The time complexity is dominated by four primary operations performed  
1929 at each of the  $T_{\text{fin}}$  iterations: gradient computation, mask updates, attribution aggregation, and the  
1930 nullification of dummy interactions. Since each of these operations is  $\mathcal{O}(N)$ , the total time complexity  
1931 per sample is  $\mathcal{O}(T_{\text{fin}} \cdot N)$ , scaling linearly with the number of input features. The space complexity  
1932 is also linear, requiring  $\mathcal{O}(N)$  memory. This includes storage for the current and previous masks, the  
1933 current and previous gradients, and the accumulated attribution values. As our proposed evaluation  
1934 metric reuses the attribution computation process, its complexity is identical.

### 1937 G.2 COMPUTATIONAL RESOURCES AND RUNTIMES

1939 All experiments were conducted on a single NVIDIA RTX A6000 GPU. We report the following  
1940 empirical runtimes: For vision models, computing attribution for a single  $224 \times 224 \times 3$  image  
1941 required approximately **22.02 seconds** for 600 iterations. The total memory usage was approximately  
1942 2,270 MB per sample, which includes the classifier model parameters and intermediate gradients. For  
1943 a language model BERT-Large model with a maximum sequence length of 512 tokens, the attribution  
computation required approximately **27 seconds** for 100 iterations on the same hardware.

---

## 1944 H MATHEMATICAL DEFINITIONS RELATED TO THE SHAPLEY VALUE

### 1946 H.1 THE ASSUMPTIONS OF AUMANN-SHAPLEY VALUE

1947 Here, we provide the formal mathematical definitions for the key assumptions and axioms referenced  
1948 in our theoretical framework.

1949 **Definition 7** (Monotonicity (Non-decreasing)). A function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , where  $\mathcal{X} \subseteq \mathbb{R}^{|N|}$ , is  
1950 **non-decreasing** if for any two input vectors  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ :

$$1951 \text{ If } x_i \geq y_i \text{ for all } i \in N, \text{ then } f(\mathbf{x}) \geq f(\mathbf{y}).$$

1952 **Definition 8** (Non-negativity of Function Output). A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  has a **non-negative output**  
1953 if its range is a subset of the non-negative real numbers, denoted  $\mathbb{R}_+$ :

$$1954 \forall \mathbf{x} \in \mathcal{X}, f(\mathbf{x}) \geq 0.$$

1955 **Definition 9** (Differentiability and Continuity of Value Function). A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is  
1956 **continuously differentiable** (denoted as  $f \in C^1(\mathcal{X})$ ) if all of its first-order partial derivatives exist  
1957 and are themselves continuous functions for all points in the domain  $\mathcal{X}$ .

$$1958 \frac{\partial f}{\partial x_i}(\mathbf{x}) \text{ exists and is continuous for all } i \in N \text{ and for all } \mathbf{x} \in \mathcal{X}.$$

1959 **Definition 10** (Non-negativity of Input Domain). The input domain  $\mathcal{X}$  of a function  $f$  is **non-negative**  
1960 if every component of every vector in the domain is non-negative. This is equivalent to the domain  
1961 being a subset of the non-negative orthant of  $\mathbb{R}^{|N|}$ .

$$1962 \mathcal{X} \subseteq [0, \infty)^{|N|}$$

1963 This means for any vector  $\mathbf{x} = (x_1, x_2, \dots, x_{|N|}) \in \mathcal{X}$ , the condition  $x_i \geq 0$  holds for all  $i \in N$ .

### 1968 H.2 THE AXIOMS OF SHAPLEY-SHUBIK VALUE AND AUMANN-SHAPLEY VALUE

1969 The classical Shapley value,  $\phi(v) = (\phi_i(v))_{i \in N}$ , is uniquely characterized by four axioms for  
1970 cooperative games  $(N, v)$ , where  $N$  is a set of features and  $v$  is the value function (Shapley, 1953).  
1971 The Aumann-Shapley value (Aumann & Shapley, 1974) extends this framework to continuous  
1972 domains. In this extension, the Efficiency, Additivity (or more generally, Linearity), and Symmetry  
1973 axioms are preserved. The Aumann-Shapley framework also introduces additional axioms, such as  
1974 Positivity or Monotonicity, which are not present in the original discrete formulation but are related  
1975 to Additivity (Linearity) and essential for ensuring a fair attribution allocation in continuous feature  
1976 spaces.

1977 **Definition 11** (Efficiency). The allocation  $\sigma(v)$  is **efficient** if  $\sum_{i \in N} \sigma_i(v) = v(N)$ .

1978 **Definition 12** (Dummy Player). Player  $i$  is a **dummy** in  $(N, v)$  if  $v(S \cup \{i\}) - v(S) =$   
1979  $v(\{i\})$ ,  $\forall S \subseteq N \setminus \{i\}$ .

1980 **Definition 13** (Null Player). Player  $i$  is a **null player** in  $(N, v)$  if  $v(S \cup \{i\}) = v(S)$ ,  $\forall S \subseteq N \setminus \{i\}$ .

1981 **Definition 14** (Symmetry). Players  $i$  and  $j$  are **symmetric** in  $(N, v)$  if  $v(S \cup \{i\}) = v(S \cup$   
1982  $\{j\})$ ,  $\forall S \subseteq N \setminus \{i, j\}$ .

1983 **Definition 15** (Additivity). For any two games  $(N, v_1)$  and  $(N, v_2)$ , the attribution method is **additive**  
1984 if  $\sigma_i(v_1 + v_2) = \sigma_i(v_1) + \sigma_i(v_2)$ ,  $\forall i \in N$ , where  $(v_1 + v_2)(S) = v_1(S) + v_2(S)$  for every coalition  
1985  $S \subseteq N$ .

1986 **Definition 16** (Linearity). Any two games  $v_1$  and  $v_2$  are **linear** if  $\sigma_i(N, \alpha v_1 + \beta v_2) = \alpha \sigma_i(N, v_1) +$   
1987  $\beta \sigma_i(N, v_2)$  for each  $i$ , where  $\alpha$  and  $\beta$  are nonnegative real numbers.

1988 **Definition 17** (Positivity). For monotonic games (i.e.,  $v(S) \leq v(T)$  whenever  $S \subseteq T$ ), the value  
1989 assigned to each player is **nonnegative**:  $\sigma_i(v) \geq 0 \quad \forall i$ .

#### 1992 H.2.1 DEFINITION OF $\epsilon$ -MONOTONICITY

1993 The Aumann-Shapley (AS) value is theoretically grounded for functions that are strictly monotonic  
1994 or, more generally, absolutely continuous. However, this assumption is rarely met by real-world  
1995 functions, which are often non-monotonic. The concept of  $\epsilon$ -monotonicity (Aumann & Shapley,  
1996 1974) was introduced to address this issue by relaxing the strict condition of perfect monotonicity,  
1997 thereby extending the applicability of the AS framework for more general functions.

1998 **Definition 18** ( $\epsilon$ -Monotonicity). A set function  $v : 2^{|N|} \rightarrow \mathbb{R}$  is defined as  $\epsilon$ -monotonic if its  
 1999 Downward Variation is less than or equal to  $\epsilon$ . The Downward Variation of a function  $v$  in the space  
 2000 of bounded variation (BV) is the supremum of the sum of the decreases in the function's value between  
 2001 successive sets, taken over all possible chains  $\phi = (S_0, S_1, \dots, S_k)$  where  $S_0 \subset S_1 \subset \dots \subset S_k = I$ ,  
 2002 with  $I$  denoting the grand coalition:

$$2003 \text{Downward Variation}(v) = \sup_{\phi} \sum_{i=0}^{k-1} \max(v(S_i) - v(S_{i+1}), 0) \leq \epsilon$$

2004  
 2005  
 2006 In essence,  $\epsilon$ -monotonicity formalizes the idea of a function being *almost monotone*, allowing for  
 2007 small, controlled deviations from a strictly increasing path. If  $\epsilon = 0$ , this definition reduces precisely  
 2008 to perfect monotonicity, where no decreases in function value are permitted.

2009  
 2010 The theoretical power of this concept is established by Aumann & Shapley (1974, Proposition B.1),  
 2011 which proves that any function  $v$  in a subspace  $Q$  of BV can be decomposed into the difference  
 2012 of two  $\epsilon$ -monotonic functions ( $v = v_1 - v_2$ , where  $v_1, v_2 \in Q^\epsilon$ ). This decomposition is highly  
 2013 significant because it provides a robust analytical framework for handling non-monotonic functions.  
 2014 It allows the monotonicity-violating (decreasing) portions of a function to be represented as the  
 2015 negative component of the decomposition ( $-v_2$ ), enabling the principled application of the AS value  
 2016 to functions that are not strictly monotonic.

## 2017 I PROOFS OF THEORETICAL PROPERTIES

### 2018 I.1 PROOF OF EFFICIENCY

2019  
 2020 **Theorem 1** (Approximate Efficiency). The attribution method proposed in Algorithm 1 satisfies the  
 2021 Efficiency axiom approximately, provided that the optimization process successfully terminates.

2022  
 2023 *Proof.* The proof consists of demonstrating the necessary and sufficient conditions for our method to  
 2024 approximate the Efficiency axiom,  $\sum_i \mathcal{A}_i(\mathbf{x}) \approx f(\mathbf{x}) - f(\mathbf{x}')$ .

2025  
 2026  
 2027 ( $\rightarrow$ ) **Necessity.** For the Efficiency axiom to hold, even approximately, two conditions are necessary.  
 2028 First, the attribution  $\mathcal{A}_i$  computed at each step must correspond to the incremental change in the  
 2029 function's value,  $f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)})$ . This correspondence is ideally captured by a first-order Taylor  
 2030 approximation:

$$2031 f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}) \approx \sum_{i=1}^{|N|} \left. \frac{\partial f}{\partial x_i} \right|_{x=x^{(t)}} \cdot (x_i^{(t+1)} - x_i^{(t)})$$

2032  
 2033  
 2034 Second, the optimization path must terminate at a point  $\mathbf{x}^{(\hat{t})}$  that is functionally equivalent to the  
 2035 baseline, such that  $f(\mathbf{x}^{(\hat{t})}) \approx f(\mathbf{x}')$ . If either of these conditions fails, the total sum of attributions  
 2036 will diverge from  $f(\mathbf{x}) - f(\mathbf{x}')$ .

2037  
 2038 ( $\leftarrow$ ) **Sufficiency.** Our algorithm is designed to sufficiently meet the necessary conditions. The  
 2039 attribution update rule (Line 23 in Algorithm 1) is structurally based on the gradient and displacement  
 2040 terms of the path integral, ensuring it serves as a proxy for the change in function value. Furthermore,  
 2041 the termination criteria (Line 18-20) explicitly force the optimization to halt when the function output  
 2042 approaches a zero-information state (i.e.,  $f(\mathbf{x}^{(\hat{t})}) \approx f(\mathbf{x}')$ ).

2043  
 2044 Combining these, the total sum of attributions can be related to the total change in function value via  
 2045 a telescoping sum:

$$2046 \sum_{i=1}^{|N|} \mathcal{A}_i(\mathbf{x}) \approx \sum_{t=0}^{\hat{t}-1} \sum_{i=1}^{|N|} \left. \frac{\partial f}{\partial x_i} \right|_{x=x^{(t)}} \cdot (x_i^{(t+1)} - x_i^{(t)}) \approx \sum_{t=0}^{\hat{t}-1} (f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}))$$

2047  
 2048 By the property of telescoping sums, this simplifies to:

$$2049 \sum_{i=1}^{|N|} \mathcal{A}_i(\mathbf{x}) \approx f(\mathbf{x}^{(\hat{t})}) - f(\mathbf{x}^{(0)})$$

2052 Given that the sufficiency of our termination criteria ensures  $f(\mathbf{x}^{(t)}) \approx f(\mathbf{x}')$  and by definition  
 2053  $\mathbf{x}^{(t=0)} = \mathbf{x}$ , we conclude:

$$2054 \sum_{i=1}^{|N|} \mathcal{A}_i(\mathbf{x}) \approx f(\mathbf{x}') - f(\mathbf{x})$$

2057 Multiplying by  $-1$  yields the desired approximate efficiency property.  $\square$

## 2059 I.2 PROOF OF DUMMY/NULL PLAYER AXIOM

2061 **Theorem 2** (Algorithmic Satisfaction of the Dummy Player Axiom). *The proposed method assigns*  
 2062 *zero attribution to any feature that is a **true null player** throughout the entire optimization path.*  
 2063 *It correctly assigns non-zero attribution to features that are only **conditionally dummy**, thereby*  
 2064 *capturing their interactive effects.*

2065 *Proof.* The proof relies on analyzing the algorithm’s behavior based on the status of a feature’s  
 2066 dummy variable at each discrete step  $t$ . We consider two cases.

2068 **Case 1: True Null Player.** A feature  $i$  is a true null player if its corresponding mask gradient  
 2069 magnitude,  $|\nabla_{m_i^{(t)}} \mathcal{L}|$ , remains below the threshold  $q_\tau$  for all steps  $t \in [0, \hat{t}]$ . By the design of our  
 2070 algorithm (Lines 10 and 12 in Algorithm 1), a feature satisfying this condition is never included in  
 2071 the optimization set  $C$ .

2073 This has a direct consequence on its mask value,  $m_i$ , and its corresponding input value,  $x_i$ :

$$2074 i \notin C \implies m_i^{(t+1)} = m_i^{(t)} \implies x_i^{(t+1)} = x_i^{(t)} \quad \forall t \quad (16)$$

2076 In other words, the input value for a true null player remains static throughout the entire path.

2077 The attribution for feature  $i$ ,  $\mathcal{A}_i$ , is defined as the sum of its incremental contributions at each step:

$$2079 \mathcal{A}_i = \sum_{t=0}^{\hat{t}-1} \mathcal{A}_i^{(t)} \quad (17)$$

2082 where the incremental attribution  $\mathcal{A}_i^{(t)}$  is given by:

$$2084 \mathcal{A}_i^{(t)} = (x_i^{(t)} - x_i^{(t+1)}) \cdot \left( \frac{\partial f(\mathbf{x}^{(t+1)})}{\partial x_i^{(t+1)}} \cdot \dots \right) \quad (18)$$

2087 Since we have established that  $x_i^{(t+1)} = x_i^{(t)}$  for a true null player, the displacement term  $(x_i^{(t)} -$   
 2088  $x_i^{(t+1)})$  is zero for all  $t$ . Therefore, the contribution at each step is zero:

$$2089 \mathcal{A}_i^{(t)} = 0 \cdot (\dots) = 0 \quad \forall t \quad (19)$$

2091 Consequently, the total attribution for any true null player is zero, and the axiom is satisfied.

$$2093 \mathcal{A}_i = \sum_{t=0}^{\hat{t}-1} 0 = 0 \quad (20)$$

2097 **Case 2: Conditionally Dummy Player.** Consider a feature  $i$  whose dummy status changes along  
 2098 the path due to interactions—that is, it is identified as a dummy at some step  $t$  but as a non-dummy  
 2099 (signal) feature at another step  $t + a$ .

- 2101 • At any step where feature  $i$  is identified as a dummy, as established in Case 1, its attributional  
 2102 update  $\mathcal{A}_i^{(t)}$  for that step is zero.
- 2103 • At any step where feature  $i$  is identified as a signal feature, it is included in the set  $C$ . Its mask  
 2104 value  $m_i$  is updated, leading to a non-zero displacement  $(x_i^{(t)} - x_i^{(t+1)})$ . This results in a non-zero  
 2105 attributional update  $\mathcal{A}_i^{(t)}$  for that step.

2106 The final attribution,  $\mathcal{A}_i = \sum_t \mathcal{A}_i^{(t)}$ , will be non-zero if the feature acts as a signal feature for at  
 2107 least one step. This outcome does not violate the axiom. Instead, it demonstrates that the algorithm  
 2108 correctly identifies that feature  $i$  is not a true null player but an interacting feature whose relevance  
 2109 is conditional on the state of other features. The final non-zero attribution correctly reflects the  
 2110 contributions made during the specific path segments where its relevance was revealed, regardless of  
 2111 whether its status changes from dummy to signal, or vice versa.  $\square$

### 2112 I.3 PROOF OF THE SYMMETRY AXIOM

2113 To prove that our method satisfies the Symmetry axiom, we first introduce a key condition as a lemma,  
 2114 upon which our main theorem depends.

2115 **Lemma 1** (Gradient Symmetry Condition). *Let  $i$  and  $j$  be two features that are symmetric with*  
 2116 *respect to the model’s output function  $f$ , such that  $f(S \cup \{i\}) = f(S \cup \{j\})$  for all  $S \subseteq N \setminus \{i, j\}$ .*  
 2117 *The Gradient Symmetry Condition is satisfied if this functional symmetry implies symmetry in the*  
 2118 *gradient of our optimization loss function  $\mathcal{L}$  with respect to the corresponding low-resolution mask*  
 2119 *components,  $m_i$  and  $m_j$ , at all steps  $t$  along the optimization path. Formally:*

$$2120 f(S \cup \{i\}) = f(S \cup \{j\}) \implies \nabla_{m_i^{(t)}} \mathcal{L} = \nabla_{m_j^{(t)}} \mathcal{L}, \quad \forall t.$$

2121 **Theorem 3** (Conditional Satisfaction of the Symmetry Axiom). *If the Gradient Symmetry Condition*  
 2122 *(Lemma 1) holds, the attribution method proposed in Algorithm 1 satisfies the Symmetry axiom,*  
 2123 *yielding identical attributions for symmetric features, i.e.,  $\mathcal{A}_i = \mathcal{A}_j$ .*

2124 *Proof.* The proof proceeds by induction, showing that if two symmetric features start with identical  
 2125 mask values and are subject to identical gradients at every step (as guaranteed by Lemma 1), their  
 2126 entire optimization trajectories and final attributions will be identical.

2127 **Base Case (t=0):** The algorithm is initialized with a uniform mask, so  $m_i^{(0)} = m_j^{(0)} = 1$ .

2128 **Inductive Step:** Assume that at step  $t$ , the mask values for the symmetric features are identical:  
 2129  $m_i^{(t)} = m_j^{(t)}$ . By Lemma 1, the gradients at this step are also identical:  $\nabla_{m_i^{(t)}} \mathcal{L} = \nabla_{m_j^{(t)}} \mathcal{L}$ .

2130 Our mask update rule (Line 15 in Algorithm 1) is deterministic and symmetric:

$$2131 m_k^{(t+1)} = m_k^{(t)} - \eta^{(t)} \cdot \text{sign} \left( \nabla_{m_k^{(t)}} \mathcal{L} \right).$$

2132 Since both the initial values ( $m_k^{(t)}$ ) and the update terms ( $\nabla_{m_k^{(t)}} \mathcal{L}$ ) are identical for  $k = i$  and  $k = j$ ,  
 2133 it follows that the updated mask values will also be identical:

$$2134 m_i^{(t+1)} = m_j^{(t+1)}.$$

2135 By induction, the mask trajectories for features  $i$  and  $j$  are identical for all  $t$ .

2136 **Conclusion:** Identical low-resolution mask trajectories ( $m_i^{(t)} = m_j^{(t)}$ ) imply identical high-resolution  
 2137 mask trajectories ( $M_i^{(t)} = M_j^{(t)}$ ) and, consequently, identical input path trajectories ( $x_i^{(t)} = x_j^{(t)}$ ).  
 2138 As all components of the incremental attribution formula (Line 23) depend on these trajectories and  
 2139 the gradients—which are also symmetric under our assumption—the incremental attributions are  
 2140 identical at every step,  $\mathcal{A}_i^{(t)} = \mathcal{A}_j^{(t)}$ . Summing over all steps yields the final result:

$$2141 \mathcal{A}_i = \sum_{t=0}^{\hat{t}-1} \mathcal{A}_i^{(t)} = \sum_{t=0}^{\hat{t}-1} \mathcal{A}_j^{(t)} = \mathcal{A}_j.$$

2142 Thus, the Symmetry axiom is satisfied, conditional on the Gradient Symmetry Lemma.  $\square$

### 2143 I.4 PROOF OF THE ADDITIVITY AXIOM

2144 **Theorem 4** (Non-Additivity of the Proposed Method). *The proposed attribution method, which relies*  
 2145 *on a function-dependent dynamic path, does **not** satisfy the Additivity axiom. This is a deliberate*  
 2146 *design trade-off to prioritize the Dummy Consistency axiom.*

2160 *Proof.* Let  $f_1$  and  $f_2$  be two distinct model functions, and let  $f = f_1 + f_2$ . The Additivity axiom  
 2161 requires that  $\mathcal{A}_i(f) = \mathcal{A}_i(f_1) + \mathcal{A}_i(f_2)$  for all features  $i$ . We demonstrate that our method violates  
 2162 this condition due to its path-finding mechanism.

2163 **i) Functional Dependency of the Path.** The core of our method is an optimization process that finds  
 2164 a path,  $\gamma(t)$ , by minimizing a loss function  $\mathcal{L}$ . This loss function is a function of the model’s output,  
 2165  $f$ . Consequently, the path itself is functionally dependent on  $f$ . Let us denote the path generated for  
 2166 a function  $g$  as  $\gamma_g(t)$ .

2167 The optimization landscape of  $\mathcal{L}(f)$  is, in general, not a simple linear combination of the landscapes  
 2168 of  $\mathcal{L}(f_1)$  and  $\mathcal{L}(f_2)$ . Therefore, the optimal paths found for each function will be different:

$$2170 \gamma_f(t) \neq \gamma_{f_1}(t) \neq \gamma_{f_2}(t)$$

2172 **ii) Consequence for Attribution.** The attribution for a feature  $i$  is defined as an integral (or sum) of  
 2173 terms evaluated along its specific path:

$$2175 \mathcal{A}_i(g) = \sum_{t=0}^{i_g-1} \mathcal{A}_i^{(t)}(g, \gamma_g(t))$$

2178 Since the paths  $\gamma_f, \gamma_{f_1}, \gamma_{f_2}$  are different, the values of the mask  $\mathbf{m}^{(t)}$ , the input  $\mathbf{x}^{(t)}$ , and the gradients  
 2179  $\nabla \mathcal{L}$  and  $\partial f / \partial x_i$  will be different at each step  $t$  for each function. For example, the incremental  
 2180 attribution for function  $f$  is:

$$2182 \mathcal{A}_i^{(t)}(f, \gamma_f(t)) = (\gamma_f(t)_i - \gamma_f(t+1)_i) \cdot \left( \frac{\partial(f_1 + f_2)}{\partial x_i} \Big|_{\mathbf{x}=\gamma_f(t+1)} \cdot \dots \right)$$

2185 Even though the gradient term is additive ( $\partial(f_1 + f_2) = \partial f_1 + \partial f_2$ ), it is evaluated at a point  
 2186  $\gamma_f(t+1)$  that is unique to the function  $f$ . Because all other terms in the product, especially the path  
 2187 displacement ( $x_i^{(t)} - x_i^{(t+1)}$ ), are also dependent on the unique path  $\gamma_f$ , there is no basis to assume  
 2188 that the incremental attributions are additive. That is,

$$2189 \mathcal{A}_i^{(t)}(f) \neq \mathcal{A}_i^{(t)}(f_1) + \mathcal{A}_i^{(t)}(f_2)$$

2191 Therefore, the total attributions are **not** additive:

$$2193 \mathcal{A}_i(f) \neq \mathcal{A}_i(f_1) + \mathcal{A}_i(f_2)$$

2195 **Discussion of the Trade-off.** This violation is not a flaw, but a fundamental trade-off. Additivity  
 2196 is guaranteed by methods like Integrated Gradients that use a fixed, model-agnostic path (a straight  
 2197 line). Our method deliberately sacrifices this property to find a dynamic, model-specific path that  
 2198 is optimized to avoid the spurious attributions of dummy features and dummy interactions and  
 2199 satisfy the Dummy Player and the Dummy Consistency axioms in DNNs. We argue that for complex,  
 2200 piecewise-linear models like rectified DNNs, satisfying the Dummy Consistency is a more critical and  
 2201 pressing requirement for generating trustworthy explanations than strictly adhering to Additivity.  $\square$

## 2202 I.5 PROOF OF THE POSITIVITY AXIOM

2204 **Theorem 5** (Satisfaction of the Positivity Axiom). *For a monotonic game, where the function  $f$  is  
 2205 non-decreasing with respect to its features, the proposed attribution method satisfies the Positivity  
 2206 axiom, yielding a non-negative attribution for all features, i.e.,  $\mathcal{A}_i \geq 0$ .*

2208 *Proof.* The proof relies on a key property of our algorithm’s optimization path when applied to a  
 2209 monotonic function, which we first establish as a lemma.

2210 **Lemma 2** (Algorithmic Enforcement of Path Monotonicity). *The optimization path  $\gamma(t)$ , generated  
 2212 by Algorithm 1, is constructed to be  $\epsilon$ -monotonically non-increasing with respect to the function value  
 2213  $f$ . That is, for any step  $t$  that does not trigger termination, the condition  $f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}) \leq \epsilon$  is  
 maintained.*

2214 *Proof of Lemma 2.* The algorithm’s objective is to reduce the function value from  $f(\mathbf{x})$  towards the  
 2215 baseline value  $f(\mathbf{x}')$  by progressively masking features. The mask update rule (Line 21) is designed  
 2216 to move in a direction that minimizes the loss  $\mathcal{L}$ , which is proportional to  $f$ . Furthermore, the  
 2217 algorithm includes an explicit safeguard: the termination condition (Line 27) halts the process if a  
 2218 step violates  $\epsilon$ -monotonicity ( $f(\mathbf{x}^{(t+1)}) - f(\mathbf{x}^{(t)}) > \epsilon$ ). Thus, any path generated by the algorithm  
 2219 is, by construction,  $\epsilon$ -monotonically non-increasing.  $\square$

2220  
 2221 We now proceed with the main proof. The total attribution  $\mathcal{A}_i$  is the sum of incremental contributions  
 2222  $\mathcal{A}_i^{(t)}$  from each step of the path. For the Positivity axiom to hold, each incremental term  $\mathcal{A}_i^{(t)}$  must be  
 2223 non-negative.

2224 The full incremental attribution at step  $t$  for a feature  $i \in C$  is given by (Line 23 in Algorithm 1):

$$2225 \mathcal{A}_i^{(t)} = (\mathbf{x}_i^{(t)} - \mathbf{x}_i^{(t+1)}) \cdot \frac{\partial f(\mathbf{x}^{(t+1)})}{\partial \mathbf{x}_i^{(t+1)}} \cdot \mathcal{W}_i \cdot \eta^{(t)} \cdot \nabla_{\mathbf{m}^{(t)}} \mathcal{L}$$

2228 While simplified proofs often focus only on the displacement and gradient terms, here we analyze the  
 2229 sign of each component in the full product under the assumption of a monotonic game:

2231 **i) The Gradient Term** ( $\frac{\partial f}{\partial \mathbf{x}_i} \geq 0$ ): By the definition of a monotonic game, the function  $f$  is non-  
 2232 decreasing with respect to its features. For a differentiable function, this implies that the partial  
 2233 derivative is non-negative at all points along the path.

2234 **ii) The Mask Gradient Term** ( $\nabla_{\mathbf{m}^{(t)}} \mathcal{L} \geq 0$ ): The algorithm’s objective is to reduce the function  
 2235 value  $f$ . For a monotonic function, this can only be achieved by reducing the presence of features  
 2236 that contribute positively. A feature  $i$  contributes positively if increasing its mask value  $m_i$  increases  
 2237 the function output. This implies that for any feature selected for an update in the set  $C$ , its mask  
 2238 gradient must be non-negative, assuming the loss  $\mathcal{L}$  is directly proportional to  $f$ .

2239 **iii) The Displacement Term** ( $(x_i^{(t)} - x_i^{(t+1)}) \geq 0$ ): The mask update rule (Line 15) is  $m_i^{(t+1)} =$   
 2240  $m_i^{(t)} - \eta^{(t)} \cdot \text{sign}(\nabla_{m_i^{(t)}} \mathcal{L})$ . Since we established in the point above that  $\nabla_{m_i^{(t)}} \mathcal{L} \geq 0$ , its sign is  
 2241  $+1$ . The update rule simplifies to  $m_i^{(t+1)} \leq m_i^{(t)}$ , which guarantees that the mask value is non-  
 2242 increasing. A non-increasing mask value implies a non-increasing input feature value, ensuring that  
 2243 the displacement term is non-negative.

2244 **iv) Scalar Terms** ( $\mathcal{W}_i, \eta^{(t)}$ ): The bilinear interpolation weights  $\mathcal{W}_i$  and the learning rate  $\eta^{(t)}$  are  
 2245 non-negative by definitions.

2246 Since all five components of the incremental attribution product are non-negative, each incremental  
 2247 attribution  $\mathcal{A}_i^{(t)}$  must be non-negative. The total attribution, being a sum of these non-negative terms,  
 2248 is therefore also non-negative:

$$2249 \mathcal{A}_i = \sum_{t=0}^{\hat{t}-1} \mathcal{A}_i^{(t)} \geq 0.$$

2250 Hence, our method satisfies the Positivity axiom for monotonic games.  $\square$

## 2256 I.6 PROOF OF DUMMY CONSISTENCY

2257 **Theorem 6** (Algorithmic Satisfaction of the Dummy Consistency Axiom). *The attribution method*  
 2258 *proposed in Algorithm 1 satisfies the Dummy Consistency axiom.*

2260 *Proof.* The proof demonstrates that the axiom is satisfied by the explicit mechanics of the active set  
 2261 filtering of our algorithm with  $C^{(t)}$ , which is designed to nullify the effects of spurious interactions.  
 2262 We first establish a lemma that formalizes this filtering behavior.

2263 **Lemma 3** (Dummy Interaction Filtering via the Active Set). *The construction of the active set  $C^{(t)}$*   
 2264 *in Algorithm 1 is designed to filter out features whose gradients become unstable due to dummy*  
 2265 *interactions. Specifically, if an update to a signal feature at step  $t - 1$  causes a previously dummy*  
 2266 *feature  $j$  at step  $t - 1$  to exhibit a non-zero gradient at step  $t$ , the active set definition prevents  $j$  from*  
 2267 *being updated.*

2268 *Proof of Lemma 3.* The core of the dummy interaction problem is that a change in a signal feature  
 2269 can induce a non-zero gradient on a feature that should be irrelevant. Let’s consider a feature  $j$  that  
 2270 was a dummy at step  $t - 1$ , meaning its gradient magnitude was below the threshold:  $|\nabla_{m_j^{(t-1)}} \mathcal{L}| \leq \tau$ .  
 2271

2272 Now, assume that due to an update of other features in the set  $C^{(t-1)}$ , an interaction occurs, and the  
 2273 gradient for feature  $j$  becomes significant at step  $t$ :  $|\nabla_{m_j^{(t)}} \mathcal{L}| > \tau$ .  
 2274

2275 Our algorithm defines the active set for the next update at step  $t$  in Eq. 7:

$$2276 C^{(t)} = \left\{ i \in \mathcal{M} \mid |\nabla_{m_i^{(t-1)}} \mathcal{L}| > \tau \text{ and } |\nabla_{m_i^{(t)}} \mathcal{L}| > \tau \right\}, \quad \mathcal{M} = \{1, \dots, h \times w\}.$$

2277 For feature  $j$  to be included in this set, it must satisfy both conditions. By construction, the mask  
 2278 update rule (Line 15 in Algorithm 1) is only applied to features within the active set  $C^{(t)}$ . As feature  
 2279  $j$  is not in this set, its mask value  $m_j$  is not updated at this step. The active set has thus successfully  
 2280 filtered out the spurious interaction effect, preventing it from propagating into the attribution path.  $\square$   
 2281

2282 We now proceed with the main proof, leveraging Lemma 3.  
 2283

2284  
 2285 ( $\rightarrow$ ) **Necessity.** For the Dummy Consistency axiom to hold, a crucial condition is necessary: the  
 2286 process of removing a dummy feature  $i$  must not alter the final attribution of any other feature  $j$ .  
 2287 This requires that the path taken by the optimization process for the game without feature  $i$ , which  
 2288 produces  $\mathcal{A}_j(\mathbf{x}_{-i}, \mathbf{x}'_{-i}; f^{R_i^x})$ , must be functionally equivalent to the path for the full game, at least  
 2289 with respect to feature  $j$ .  
 2290

2291 ( $\leftarrow$ ) **Sufficiency.** Our algorithm is sufficient to meet this necessary condition. Consider a true dummy  
 2292 feature  $i$ . As established in the proof of the Dummy Player axiom, a true dummy is never included in  
 2293 the active set  $C^{(t)}$  at any step. Its mask value  $m_i$  remains constant, and its influence is never actively  
 2294 reduced by our optimization.

2295 Now, consider the definition of Dummy Consistency: comparing the attribution  $\mathcal{A}_j(\mathbf{x}, \mathbf{x}'; f)$  with  
 2296  $\mathcal{A}_j(\mathbf{x}_{-i}, \mathbf{x}'_{-i}; f^{R_i^x})$ .  
 2297

2298 **Scenario i** In the full game with input  $\mathbf{x}$ , our algorithm proceeds by identifying and reducing the  
 2299 mask values of signal features based on the active set criteria. As feature  $i$  is a dummy, its status  
 2300 never affects the decision to update the mask of any other feature  $j$ .  
 2301

2302 **Scenario ii** In the reduced game with input  $\mathbf{x}_{-i}$  and function  $f^{R_i^x}$ , the dummy feature  $i$  is absent  
 2303 from the start.

2304 Since the presence of the dummy feature  $i$  has no influence on the sequence of mask updates for  
 2305 any other feature  $j \neq i$  in the full game as established by the filtering mechanism in Lemma 3, the  
 2306 sequence of active sets  $C^{(t)}$ , the mask updates, and consequently the entire attribution path  $\gamma(t)$  will  
 2307 be identical in both scenarios.

2308 Because the attribution paths and all relevant gradients for any feature  $j$  are identical in both the full  
 2309 game and the reduced game, their final computed attributions must also be identical:

$$2310 \mathcal{A}_j(\mathbf{x}, \mathbf{x}'; f) = \mathcal{A}_j(\mathbf{x}_{-i}, \mathbf{x}'_{-i}; f^{R_i^x}) \quad \forall j \neq i.$$

2311 Thus, our method satisfies the Dummy Consistency axiom by its procedural design.  $\square$   
 2312  
 2313

## 2314 J DERIVATION OF THE DCAPM ATTRIBUTION FORMULA

2315  
 2316 The attribution for our method is grounded in the classical path integral formulation from Eq. 1.  
 2317 We adapt this general form by substituting the specific components derived from our proposed  
 2318 optimization-based path,  $\gamma(t)$ .  
 2319

2320 First, we restate the standard, continuous form of the path method in Eq. 1:

$$2321 \mathcal{A}_i^\gamma(\mathbf{x}, \mathbf{x}'; f) = \int_0^{\hat{t}} \frac{\partial f(\gamma(t))}{\partial \gamma_i(t)} \frac{d\gamma_i(t)}{dt} dt. \quad (21)$$

2322 The term  $\frac{d\gamma_i(t)}{dt}$  represents the instantaneous velocity of the  $i$ -th feature along the path. We will  
 2323 substitute our derived expression for this term from the following derivation result into the integral.  
 2324

### 2325 J.1 DERIVATION OF THE DERIVATIVE OF THE PATH FUNCTION IN EQ. 3

2327 Here, we provide the complete derivation for the derivative of the proposed path function with respect  
 2328 to the time step,  $\frac{d\gamma(t)}{dt}$ . The derivation relies on the definitions of the path function  $\gamma(t)$  and the update  
 2329 rule for the mask  $\mathbf{m}^{(t)}$  in Eq. 3. First, let's define the components of the equation:  
 2330

2331 **Attribution Path  $\gamma(t)$ :** The path is defined as a feature-wise interpolation between the input  $\mathbf{x}$  and  
 2332 the reference point  $\mathbf{x}'$ . The interpolation is controlled by a function  $u(\mathbf{m}^{(t)})$ , where each component  
 2333 of  $u$  is between 0 and 1.  $\odot$  denotes the element-wise multiplication.  
 2334

$$2335 \gamma(t) := u(\mathbf{m}^{(t)}) \odot \mathbf{x} + (1 - u(\mathbf{m}^{(t)})) \odot \mathbf{x}' \quad (22)$$

2337 **Update Rule for  $\mathbf{m}^{(t)}$ :** The mask variable  $\mathbf{m}^{(t)}$  is updated over time based on the gradient of a  
 2338 loss function  $\mathcal{L}$ . The derivative of mask with respect to the time is given by a sign-based gradient  
 2339 descent rule which is from PGD technique (Madry et al., 2018) with a small positive learning rate  $\eta$ :  
 2340

$$2341 \frac{d\mathbf{m}^{(t)}}{dt} := \eta \cdot \text{sign}(\nabla_{\mathbf{m}^{(t)}} \mathcal{L}) \quad (23)$$

2343 where the derivatives of loss function with respect to mask is abbreviated as  $\nabla_{\mathbf{m}^{(t)}} \mathcal{L} = \frac{d\mathcal{L}(\mathbf{x}, \mathbf{x}', \mathbf{m}^{(t)})}{d\mathbf{m}^{(t)}}$ .  
 2344

2345 **Jacobian of the Interpolation Function  $\mathcal{W}$ :**  $\mathcal{W}$  denotes the Jacobian matrix of the interpolation  
 2346 function  $u$  with respect to its input  $\mathbf{m}^{(t)}$ .  
 2347

$$2348 \mathcal{W} := \frac{\partial u(\mathbf{m}^{(t)})}{\partial \mathbf{m}^{(t)}} \quad (24)$$

2350 We start by differentiating the path  $\gamma(t)$  with respect to time  $t$ .  
 2351

$$2352 \frac{d\gamma(t)}{dt} = \frac{d}{dt} \left[ u(\mathbf{m}^{(t)}) \odot \mathbf{x} + (1 - u(\mathbf{m}^{(t)})) \odot \mathbf{x}' \right] \quad (a)$$

$$2354 = \frac{d}{dt} \left( u(\mathbf{m}^{(t)}) \odot \mathbf{x} \right) + \frac{d}{dt} \left( (1 - u(\mathbf{m}^{(t)})) \odot \mathbf{x}' \right) \quad (b)$$

$$2356 = \mathbf{x} \odot \frac{du(\mathbf{m}^{(t)})}{dt} + \mathbf{x}' \odot \frac{d(1 - u(\mathbf{m}^{(t)}))}{dt} \quad (c)$$

$$2358 = \mathbf{x} \odot \frac{du(\mathbf{m}^{(t)})}{dt} - \mathbf{x}' \odot \frac{du(\mathbf{m}^{(t)})}{dt} \quad (d)$$

$$2360 = (\mathbf{x} - \mathbf{x}') \odot \frac{du(\mathbf{m}^{(t)})}{dt} \quad (e)$$

$$2362 = (\mathbf{x} - \mathbf{x}') \odot \left( \frac{\partial u(\mathbf{m}^{(t)})}{\partial \mathbf{m}^{(t)}} \cdot \frac{d\mathbf{m}^{(t)}}{dt} \right) \quad (f)$$

$$2364 = (\mathbf{x} - \mathbf{x}') \odot \left( \mathcal{W} \cdot \left[ \eta \cdot \text{sign} \left( \frac{d\mathcal{L}(\mathbf{x}, \mathbf{x}', \mathbf{m}^{(t)})}{d\mathbf{m}^{(t)}} \right) \right] \right) \quad (g)$$

$$2366 = (\mathbf{x} - \mathbf{x}') \odot \mathcal{W} \cdot \eta \cdot \text{sign} \left( \frac{d\mathcal{L}(\mathbf{x}, \mathbf{x}', \mathbf{m}^{(t)})}{d\mathbf{m}^{(t)}} \right)$$

2371  
 2372 In the step (a) we apply the sum rule for differentiation. (b) since  $\mathbf{x}$  and  $\mathbf{x}'$  are constant with respect to  
 2373  $t$ , they can be treated as constants. (c) the derivative of  $(1 - u)$  is  $-du/dt$ . (d) factor out the common  
 2374 term using the distributive property of the Hadamard product. (f) substitute the definitions for the  
 2375 Jacobian ( $\mathcal{W}$ ) and the update rule for  $\mathbf{m}$ . (g) rearrange the scalar term eta for clarity. This completes  
 the full derivation from the definition of the path  $\gamma(t)$  to the final expression in Eq. 9.  $\square$

2376  
2377  
2378  
2379  
2380  
2381  
2382  
2383  
2384  
2385  
2386  
2387  
2388  
2389  
2390  
2391  
2392  
2393  
2394  
2395  
2396  
2397  
2398  
2399  
2400  
2401  
2402  
2403  
2404  
2405  
2406  
2407  
2408  
2409  
2410  
2411  
2412  
2413  
2414  
2415  
2416  
2417  
2418  
2419  
2420  
2421  
2422  
2423  
2424  
2425  
2426  
2427  
2428  
2429

## J.2 DERIVATION OF THE DUMMY CONSISTENT ATTRIBUTION PATH METHOD

Substituting our derived expression for the second derivative term  $\frac{d\gamma_i(t)}{dt}$  into the general path integral formula yields the specific attribution equation for our method:

$$\mathcal{A}_i^\gamma = \int_0^{\hat{t}} \frac{\partial f(\gamma(t))}{\partial \gamma_i(t)} \cdot [(x_i - x'_i) \cdot \mathcal{W}_i \cdot \eta \cdot \text{sign}(\nabla_{\mathbf{m}(t)} \mathcal{L})] dt \quad (25)$$

This formula represents the attribution in continuous time. To implement this computationally, we must discretize the integral into a sum over finite time steps. In this discrete setting, the infinitesimal term  $\frac{d\gamma_i(t)}{dt} dt$  is replaced by the finite displacement over a single step,  $\Delta\gamma_i(t) = \gamma_i(t) - \gamma_i(t+1)$ .

To establish this link formally, we analyze the structure of the incremental displacement. From the definition of our path function, the displacement between two consecutive steps is:

$$\begin{aligned} \Delta\gamma_i(t) &= \gamma_i(t) - \gamma_i(t+1) \\ &= [x_i \cdot \mathcal{W}_i(\mathbf{m}^{(t)}) + x'_i \cdot (1 - \mathcal{W}_i(\mathbf{m}^{(t)}))] - [x_i \cdot \mathcal{W}_i(\mathbf{m}^{(t+1)}) + x'_i \cdot (1 - \mathcal{W}_i(\mathbf{m}^{(t+1)}))] \\ &= (x_i - x'_i) \cdot (\mathcal{W}_i(\mathbf{m}^{(t)}) - \mathcal{W}_i(\mathbf{m}^{(t+1)})) \end{aligned} \quad (26)$$

This derivation reveals a critical insight: the per-step displacement  $(\gamma_i(t) - \gamma_i(t+1))$  is not a constant, but is in fact the total displacement  $(x_i - x'_i)$  scaled by the change in the mask over that step.

Our final attribution formula (Eq. 10) is a discrete implementation of the integral that incorporates this relationship. It replaces the total displacement factor  $(x_i - x'_i)$  inside the continuous integral (Eq. 25) with the more fine-grained, per-step displacement factor  $(\gamma_i(t) - \gamma_i(t+1))$  in the final summation:

$$\mathcal{A}_i^\gamma(\mathbf{x}, \mathbf{x}'; f) = \sum_{t=0}^{\hat{t}-1} \frac{\partial f(\gamma(t))}{\partial \gamma_i(t)} \cdot \underbrace{[(\gamma_i(t) - \gamma_i(t+1)) \cdot \mathcal{W}_i \cdot \eta \cdot \text{sign}(\nabla_{\mathbf{m}(t)} \mathcal{L})]}_{\text{Custom discrete weighting}} \quad (27)$$

This formulation ensures that the contribution at each step is weighted by the actual movement of features during that specific step, which is guided by our dummy interaction filtering mechanism.

## K CONCLUDING REMARKS

### K.1 LIMITATION AND FUTURE DIRECTIONS

This work opens several promising avenues for future research. First, while our gradient-based optimization is computationally efficient compared to the combinatorial computation approach of the original Shapley-Shubik value (Shapley & Shubik, 1954), it is susceptible to local minima. Exploring methods to improve global optimality, such as alternative optimization strategies or initialization techniques, is a key next step. Second, our current approach focuses on nullifying second-order interactions involving a group of dummy features. Extending this framework to detect and mitigate third or higher-order interactions among dummy or signal features presents an essential and challenging direction for future investigation. Ultimately, to ensure a comprehensive and fair assessment of attribution methods as we progress, we advocate for a joint evaluation protocol with other evaluation metrics. We propose that our Dummy Consistency metric be used in conjunction with established sensitivity measures to provide a more comprehensive understanding of an explanation’s axiomatic soundness and faithfulness.

### K.2 LARGE LANGUAGE MODEL USAGE STATEMENT

We utilized a large language model as a research and writing assistant throughout the preparation of this manuscript. Its applications included refining academic writing, proofreading, clarifying theoretical concepts, and assisting with code generation. The authors validated all results and retain full intellectual responsibility for the final manuscript.