

# UNLEARNING-BASED NEURAL INTERPRETATIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Gradient-based interpretations often require an anchor point of comparison to avoid saturation in computing feature importance. We show that current baselines defined using static functions—constant mapping, averaging or blurring—inject harmful colour, texture or frequency assumptions that deviate from model behaviour. This leads to accumulation of irregular gradients, resulting in attribution maps that are biased, fragile and manipulable. Departing from the static approach, we propose UNI to compute an (un)learnable, debiased and adaptive baseline by perturbing the input towards an *unlearning direction* of steepest ascent. Our method discovers reliable baselines and succeeds in erasing salient features, which in turn locally smooths the high-curvature decision boundaries. Our analyses point to unlearning as a promising avenue for generating faithful, efficient and robust interpretations.

## 1 INTRODUCTION

The utility of large models is hampered by their lack of explainability and robustness guarantees. Yet breakthroughs in language modelling (Meta, 2024; Anthropic, 2024; Jiang et al., 2023; Google, 2024; Achiam et al., 2023) and generative computer vision (Rombach et al., 2022; Liu et al., 2023; Deepmind, 2024; Brooks et al., 2024) yield promising high-stakes applications, spanning domains of healthcare, scientific discovery, law and finance. As such, being able to interpret these models has become a primary concern for researchers, policymakers and the general populace, with international calls for explainability, accountability and fairness in AI decision-making (European Commission, 2021; White House OSTP, 2022; Bengio et al., 2023). To this end, recent works focus on the 2 main directions of making models *inherently explainable* (Böhle et al., 2022; Brendel & Bethge, 2018; Koh et al., 2020; Bohle et al., 2021; Chen et al., 2019; Ross et al., 2017) and *post-hoc interpretable* (Bau et al., 2017; Kim et al., 2018; Zhou et al., 2018; Ghorbani et al., 2019b). Unfortunately, the former is marred by the status quo of proprietary models and prohibitive training costs. This motivates seeking robust attributions which reliably explain model predictions, to facilitate better risk assessment and trade-off calibration (Böhle et al., 2022; Doshi-Velez & Kim, 2017).

Post-hoc methods explain a black-box model’s output by attributing its decision back to predictive features of the input. They achieve this via leveraging components of the model itself (*e.g.* gradients and activations), or through approximation with a simpler, interpretable simulator. A desirable post-hoc explanation should exhibit *high faithfulness* – to be rationale-consistent (Yeh et al., 2019; Atanasova et al., 2020) with respect to a model’s decision function; *low sensitivity* – to yield reliably similar saliency predictions for input features in the same local neighbourhood (Alvarez Melis & Jaakkola, 2018; Ghorbani et al., 2019b); *low complexity* – the explanation should be functionally simpler and more understandable than the original black-box model (Bhatt et al., 2021).

Gradient-based saliency methods are widely used for feature attribution, due to their simplicity, efficiency and post-hoc accessibility. This can be further decomposed into 3 families: perturbative, backpropagative and path-based, which we detail in Section 6. Gradient-based attribution is intuitive since the first-order derivative reveals which features significantly influence the model’s classification decision. However, naively using local gradients yields unfaithful attributions due to saturation, where the non-linear output function flattens in vicinity of the input and zero gradients are computed (Sundararajan et al., 2017; 2016). To improve gradient-sensitivity, later methods introduce a baseline input for reference, and backpropagate the difference in activation scores on a path between the reference and image-of-interest (Shrikumar et al., 2016; Sundararajan et al., 2017). The baseline is chosen to be devoid of predictive features and far away from the saturated local neighbourhood. However, such methods accumulate gradient noise when interpolating from the baseline to the input,

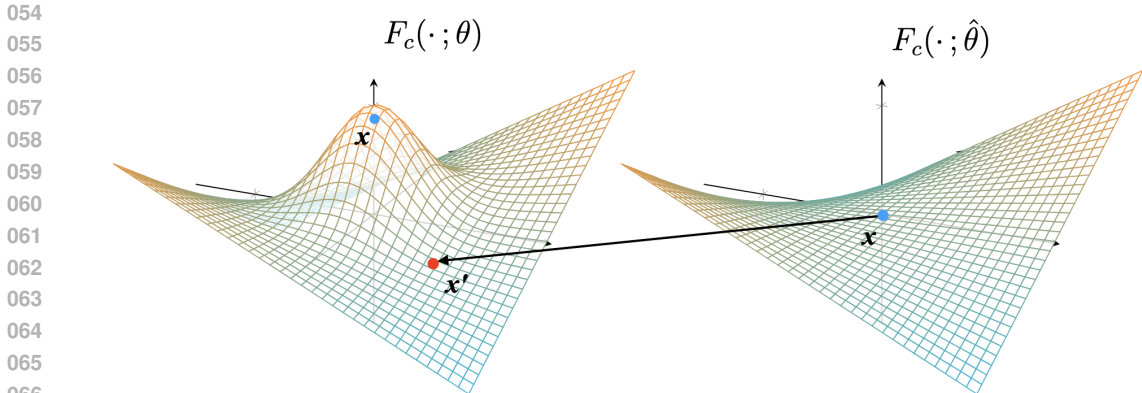


Figure 1: *Left*: Confidence of original model  $\theta$  at image  $x$  and baseline  $x'$ . *Right*: Confidence of *unlearned* model  $\hat{\theta}$  at image  $x$ . After unlearning in the model space  $\theta \mapsto \hat{\theta}$ , we optimise the baseline to match the unlearned input confidence, such that  $F_c(x'; \theta) \approx F_c(x; \hat{\theta})$ .

leading to high local sensitivity (Ancona et al., 2018). Consequently, attribution maps become disconnected, sparse and irregular, where the saliency scores fluctuate wildly between neighbouring pixels of the same object and are visually noisy (Adebayo et al., 2018). This noise accumulation has two root causes—a *poorly chosen baseline* and *high-curvature output manifold* along the path features. Previous works (Sturmfels et al., 2020; Xu et al., 2020) have sought better baselines by empirically comparing between using a black image, a gaussian noised image, a gaussian blurred image, a uniformly noised image, an inverted colour image, as well as averaging attributions over several baseline choices. However, the correct baseline to represent a lack of salient features depends heavily on the specific classification task, on the trained model and on the input image. Indeed, the optimal baseline varies for each task–model–image combination (Akhtar & Jalwana, 2023); the baseline problem remains largely unsolved. Turning to the second problem of high-curvature output manifold, because trained neural networks exhibit approximately piece-wise linear decision boundaries (Goodfellow et al., 2014), inputs near function transitions are vulnerable to perturbative attacks. By simply adding norm-bounded, imperceptible adversarial noise to the input image, attackers can dramatically alter the attribution map without changing the model’s class prediction (Ghorbani et al., 2019a; Dombrowski et al., 2019). Methods of mitigation include explicit smoothing via averaging over multiple noised gradient attributions (Smilkov et al., 2017); adaptively optimising the integration path of attribution (Kapishnikov et al., 2021); imposing an attribution prior during training and optimising it at each step (Erion et al., 2021). However, all of these proposals starkly increase the complexity of attribution, requiring computationally costly forward and backward propagation steps.

To tackle the problematic triad of 1. *post-hoc attribution biases*, 2. *poor baseline definition*, 3. *high-curvature output manifold*, we propose UNI to discover debiased baselines by locally *unlearning* inputs, *i.e.* perturbing them in the unlearning direction of steepest ascent, as visualised in Figure 1. Towards better baselines, our unlearned reference is by definition explicitly optimised to lower output class confidence and can empirically erase or occlude salient features. We also say that the unlearned baseline is specific and featureless w.r.t. each task–model–input combination. Unlike the practice of using a black image baseline—which creates a post-hoc colour bias that darker pixels are less likely to be salient, UNI does not impose additional, pixel-wise colour, scale or geometric assumptions that are not already present in the model itself. Finally, we address the high-curvature decision boundaries problem by realising that this is a product of the training process—targeted unlearning smooths the decision boundary of the model within the vicinity of the input. For a more detailed overview on the principle of machine unlearning, we refer the reader to Section 6 of the supplement. We empirically verify this local smoothing effect by measuring the normal curvature of the model function before and after unlearning; we also demonstrate that unlearning makes attributions resistant to perturbative attacks. Our contributions can be summarised as follows:

1. *Post-hoc attribution can impose new biases.* We approach the baseline challenge from the fresh lens of post-hoc biases. We show that static baselines (*e.g.* black, blurred, random noise) inject additional colour, texture and frequency assumptions that are not present in the original model’s decision rule, which leads to explanation infidelity and inconsistency.

- 108 2. *A well-chosen baseline is specific and featureless.* We establish theoretically grounded principles  
 109 for sound baseline definitions, by formalising the idea of an “absence of signal” through an  
 110 unlearning direction of steepest ascent in model loss. By unlearning predictive features in the  
 111 model space and matching this reference model’s activations with a perturbation in the input  
 112 space, we introduce a new definition of “feature absence” and a novel attribution algorithm.
- 113 3. *Unlearning reduces the curvature of decision boundaries and increases robustness.* Targeted  
 114 unlearning simulates function statistics of unseen data, and smooths the curvature of the output  
 115 manifold around the sample. This is characterised by low geodesic path curvature and bounded  
 116 principal curvature of the output surface. This points to reduced variability of the gradient vector  
 117 under small-norm input perturbations, leading to better attribution robustness and faithfulness.

## 119 2 PRELIMINARIES

120  
 121 We consider feature attribution for trained deep neural networks within image classification.  
 122 Informally, we seek to assign scores to each pixel of an image for quantifying the pixel’s influence  
 123 (sign and magnitude) on the predicted output class confidence. It is noteworthy that attributions can  
 124 be signed: a negative value indicates that removing the pixel increases the target class probability.

### 126 2.1 NOTATION

127  
 128 The input (feature) space is denoted as  $\mathcal{X} \subset \mathbb{R}^{d_X}$ , where  $d_X$  is the number of pixels in an image.  
 129 The output (label) space is  $\mathcal{Y} \subset \mathbb{R}^{d_Y}$ ;  $\mathcal{Y}$  is the set of all probability distributions on the set of  
 130 classes. The model space is denoted as  $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$ . A trained model  $F : x \mapsto (F_1(x), \dots, F_{d_Y}(x))$   
 131 returns the probability score  $F_c(x)$  of each class  $c$ . Attribution methods are thus functions  $\mathcal{A} :$   
 132  $\{1, \dots, d_X\} \times \mathcal{F} \times \{1, \dots, d_Y\} \times \mathcal{X} \rightarrow \mathbb{R}$ , where  $\mathcal{A}(i, F, c, x)$  is the importance score of pixel  $i$  of  
 133 image  $x$  for the prediction made by  $F_c$ . For convenience, we use the shorthand  $\mathcal{A}_i(x)$  to refer to  
 134 the attributed saliency score of a pixel  $i$  for a specific class prediction  $c \in \{1, \dots, d_Y\}$ . We express  
 135 a linear path feature as  $\gamma(x', x, \alpha) : \mathbb{R}^{d_X} \times \mathbb{R}^{d_X} \times [0, 1] \rightarrow \mathbb{R}^{d_X}$ , where  $\gamma = (1 - \alpha)x' + \alpha x$  and  
 136 employ shorthands  $\gamma(0) = x', \gamma(1) = x$ .

## 137 3 GRADIENT-BASED ATTRIBUTIONS IN A NUTSHELL

### 139 3.1 LIMITATIONS

140  
 141 Taking the local gradients of a model’s output confidence map  $F_c(x)$  – for target class  $c$  – is a tried  
 142 and tested method for generating explanations. Commonly termed Simple Gradients (Erhan et al.,  
 143 2009; Baehrens et al., 2010; Simonyan et al., 2013),  $\mathcal{A}_i^{\text{SG}}(x) = \nabla_{x_i} F_c(x)$  can be efficiently computed  
 144 for most model architectures. However, it encounters output saturation when activation functions like  
 145 ReLU and Sigmoid are used, leading to zero gradients (hence null attribution) even for important  
 146 features (Sundararajan et al., 2017; 2016). DeepLIFT (Shrikumar et al., 2016) reduces saturation by  
 147 introducing a “reference state”. A feature’s saliency score is decomposed into positive and negative  
 148 contributions by backpropagating and comparing each neuron’s activations to that of the baseline.  
 149 Integrated Gradients (IG) (Sundararajan et al., 2017) similarly utilises a reference, black image and  
 150 computes the integral of gradients interpolated on a straight line between the image and the baseline.

$$151 \mathcal{A}_i^{\text{IG}}(x) = (x_i - x'_i) \int_{\alpha=0}^1 \nabla_{x_i} F_c(x' + \alpha(x - x')) d\alpha \quad (1)$$

152  
 153 Practically, the integral is approximated by a Riemann sum. Of existing methods, IG promises  
 154 desirable, game-theoretic properties of “sensitivity”, “implementation invariance”, “completeness”  
 155 and “linearity”. We consequently focus on analysing and developing the IG framework, though the  
 156 proposal to unlearn baselines can be applied to most mainstream gradient-based saliency methods.  
 157 Despite the advantages of IG, its soundness depends on a good *baseline definition*—an input which  
 158 represents the “absence” of predictive features; also on having stable *path features*—a straight-line of  
 159 increasing output confidence along the path integral from baseline to target image. In the conventional  
 160 setting where a black image is used, Akhtar & Jalwana (2023) prove that IG assumptions are violated  
 161 due to ambiguous path features, where extrema of model confidences lie along the integration  
 path instead of at the endpoints of the baseline (supposed minimum) and input image (supposed

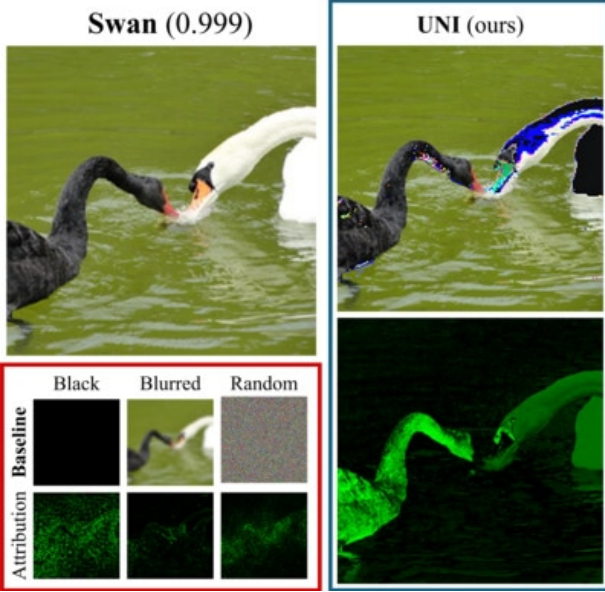


Figure 2: We visualise post-hoc biases imposed by static baselines—black baseline (colour), blurred (texture), random (frequency). UNI learns to mask out predictive features used by the model, generating reliable attributions.

**Algorithm 1** UNI: unlearning direction, baseline matching and path-attribution

- 1: **Given** model  $F(\cdot, \theta)$ ; inputs  $(x, y)$
- 2: **Choose** unlearning step-size  $\eta$ ; PGD steps  $T$ , budget  $\varepsilon$ , step-size  $\mu$ ; Riemann approximation steps  $B$
- 3: **Initialise** perturbation  $\delta^0$
- 4: **Unlearning direction.**  

$$\hat{\theta} = \theta + \eta \frac{\nabla_{\theta} \mathcal{L}(F_c(x; \theta), y)}{\|\nabla_{\theta} \mathcal{L}(F_c(x; \theta), y)\|}$$
- 5: **for**  $t = 0, \dots, T - 1$  **do**  

$$C = D_{KL}(F(x; \hat{\theta}) \| F(x + \delta^t; \theta))$$

$$\delta^{t+1} = \delta^t - \mu \nabla_{\delta} C$$

$$\delta^{t+1} = \varepsilon \frac{\delta^{t+1}}{\|\delta^{t+1}\|}$$
- 6: **end for**
- 7: **Baseline** definition  $x' = x + \delta^T$
- 8: **Attributions** computation:  $\mathcal{A}_i^{\text{UNI}}(x)$

$$= \frac{(x_i - x'_i)}{B} \sum_{k=1}^B \nabla_{x_i} F_c \left( x' + \frac{k}{B} (x - x'); \theta \right)$$

maximum). Sturmfels et al. (2020) enumerate problems with other baselines obtained via gaussian blurring, maximum distance projection, uniform noise. Despite the diversity of baseline alternatives, no candidate is optimal for each and every attributions setting. For instance, models trained with image augmentations (e.g. colour jittering, rescaling, gaussian blur) yield equivalent or even higher confidences for blurred and lightly-noised baselines—we need baselines that are well-optimised for each task–model–input combination. Without principled baselines, problems of non-conformant intermediate paths and counter-intuitive attribution scores will doubtlessly persist.

### 3.2 POST-HOC BIASES ARE IMPOSED

Since the baseline represents an absence of or reduction in salient features, static baseline functions (e.g. black, blurred, noised) implicitly assume that similar features (e.g. dark, smooth, high-frequency) are irrelevant for model prediction. To illustrate this intuition, we can consider IG with a black baseline, wherein it becomes more difficult to attribute dark but salient pixels. Due to the colour bias that “near-black features are unimportant”, the term  $(x_i - x'_i)$  is small and requires a disproportionately large gradient  $\nabla_{x_i} F_c(\cdot)$  to yield non-negligible attribution scores. Indeed, this is what we observe in Figures 2, 3, 11, where darker features belonging to the object-of-interest cannot be reliably identified. We further empirically verify that each static baseline imposes its own post-hoc bias by experimenting on ImageNet-C (Hendrycks & Dietterich, 2019). Corresponding to the 3 popular baseline choices for IG (all-black, gaussian blurred, gaussian noised), we focus on the families of *digital* (brightening and saturation), *blur* (gaussian and defocus blur) and *noise* (gaussian and shot noise) common-corruptions. Figures 4, 12 demonstrate that IG with a blurred baseline fails to attribute blurred inputs due to saturation and overly smoothed image textures; Figures 5, 13 visualise how a noised IG baseline encounters high-frequency noise and outputs irregular, high-variance attribution scores, even for adjacent pixels belonging to the same object. We crucially emphasise that such colour, texture and frequency biases are not present naturally in the pre-trained model but rather injected implicitly by a suboptimal choice of static baseline. The observation that poor baseline choices create *attribution bias* has so far been overlooked. As such, we depart entirely from the line of work on alternative static baseline towards adaptively (un)learning baselines with gradient-based optimisation. *UNI eliminates all external assumptions except for the model’s own predictive bias.*

216  
217  
218  
219  
220  
221  
222  
223  
224  
225

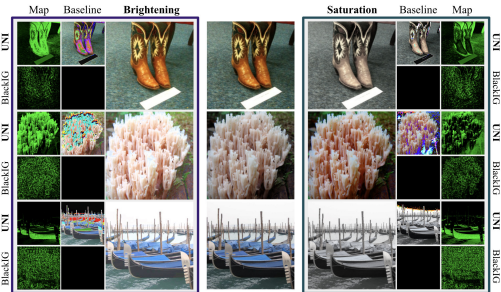


Figure 3: When the brightness or saturation is altered, IG with a black baseline fails to identify dark features, such as the boat’s hull (R3) or the top of the boot (R1).

226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238

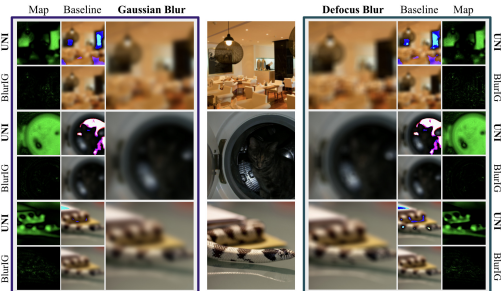


Figure 4: Under gaussian or defocus blur, IG with a blurred baseline suffers from saturation; has overly smooth texture; does not yield meaningful features.

239  
240  
241  
242  
243  
244  
245  
246  
247

## 4 UNI: UNLEARNING-BASED NEURAL INTERPRETATIONS

### 4.1 BASELINE DESIDERATA

A desirable baseline should preserve the game-theoretic properties of path-attribution (Section 3.1) and refrain from imposing post-hoc attribution biases (Section 3.2). For every given task-model-image triad, a well-chosen baseline should be 1. *image-specific*—be connected via a path feature of low curvature to the original image; 2. *reflect only the model’s predictive biases*—salient image features should be excluded from the baseline; be 3. *less task-informative than the original image*—interpolating from the baseline towards the input image should yield a path of increasing predictive confidence. We now introduce the UNI pipeline: first, unlearn predictive information in the model space; then, use activation-matching between unlearned and trained models to mine a featureless baseline in the image space; finally, interpolate along the low-curvature, conformant and consistent path from baseline to image to compute reliable explanations in the attributions space. Figure 8 visuals and Table 7 results attest to UNI’s ability to compute specific, unlearned baselines for attribution.

258  
259

### 4.2 DESIRABLE PATH FEATURES

260  
261  
262  
263  
264  
265  
266  
267  
268  
269

**Proximity** The meaningfulness of the attributions highly depends on the meaningfulness of the path. We aim for a smooth transition between absence and presence of features; and this intuitively cannot be achieved if the baseline and input are too far apart. Srinivas & Fleuret (2019) formalises this intuition through the concept of *weak dependence*, and proves that this property can only be compatible with completeness in the case where the baseline and the input lie in the same connected component (in the case of piecewise-linear models). An obvious implementation of this proximity condition in the general case is to bound the distance  $\|x - x'\|$  to a certain value  $\epsilon$ . This is strictly enforced in Algorithm 1 by normalising the perturbation at each step  $t$ .

**Low Curvature.** The curvature of the model prediction along the integrated path has been identified Dombrowski et al. (2019) as one of the key factors influencing both the sensitivity and faithfulness of

Table 1: Path monotonicity scores with Spearman correlation coefficient (higher = better). Integrating from a “featureless” baseline to the sample should give a path of monotonically increasing prediction confidence.

	UNI	IG	BlurIG	GIG
ResNet-18	.97 ±.222	.69 ±.460	.57 ±.576	.45 ±.476
Eff-v2-s	.95 ±.258	.28 ±.615	.34 ±.613	.38 ±.437
ConvNeXt-T	.99 ±.121	.76 ±.379	.77 ±.486	.46 ±.485
VGG-16-bn	.94 ±.286	.69 ±.474	.60 ±.544	.46 ±.479
ViT-B-16	.89 ±.396	.71 ±.399	.27 ±.648	.44 ±.468
SwinT	.97 ±.189	.88 ±.326	.88 ±.482	.45 ±.474



Figure 5: Gaussian and shot noise create visual artifacts prominent in noised-baseline IG. Frequency bias leads to disparate scores for adjacent pixels.

the computed attributions. We substantiate the intuition that a smooth and regular path is preferred by analysing the Riemannian sum calculation. Assuming that the function  $g : \alpha \in [0, 1] \mapsto \nabla F_c(x' + \alpha(x - x'))$  is derivable with a continuous derivative (i.e.  $C^1$ ) on the segment  $[x', x]$ , elementary calculations and the application of the Taylor-Lagrange inequality give the following error in the Riemann approximation of the attribution,

$$\left| (x_i - x'_i) \int_{\alpha=0}^1 g(\alpha) d\alpha - \frac{(x_i - x'_i)}{B} \sum_{k=1}^B g\left(\frac{k}{B}\right) \right| \leq \frac{M \|x - x'\|^2}{2B} \quad (2)$$

where  $M = \max_{\alpha \in [0,1]} \frac{dg}{d\alpha} = \max_{\alpha \in [0,1]} \frac{\partial^2 F_c(x' + \alpha(x - x'))}{\partial \alpha^2}$  exists by continuity of  $g'$  on  $[0, 1]$ . Thus, lower curvature along the path implies a lower value of the constant  $M$ , which in turn implies a lower error in the integration calculation. A smaller value  $B$  of Riemann steps is needed to achieve the same precision. More generally, a low curvature (i.e. eigenvalues of the hessian) on and in a neighbourhood of the baseline and path reduces the variability of the calculated gradients under small-norm perturbations, increasing the sensitivity and consistency of the method. Empirically, we observe a much lower curvature of paths computed by UNI, as per Table 1 and Appendix Figures 20, 21, 22, 23, 24, 25. Figure 10 also confirms the increased robustness to Riemann sum error induced.

**Monotonic.** Intuitively, the path  $\gamma$  defined by interpolating from the “featureless” baseline  $x'$  to the input image  $x$  should be *monotonically increasing* in output class confidence. At the image level, for all  $j, k$  such that  $j \leq k$ , since  $\|\gamma(j) - x\| \geq \|\gamma(k) - x\|$ , therefore the predictive confidence should be non-decreasing and order-preserving:  $F_c(\gamma(j)) \leq F_c(\gamma(k))$ . Constraining  $\gamma$  to be monotonically increasing suffices to satisfy a weak version of the criteria for *valid path features* (Akhtar & Jalwana, 2023):  $\text{sgn}(\nabla_x F_c(x)) \cdot \text{sgn}(\nabla_{\bar{x}} F_c(x')) = 1$  is naturally met.

### 4.3 EFFECTS OF UNLEARNING AND MATCHING

We explain the success of UNI with the illustrative example of a three gaussians mixture model. Figure 6 computes unlearning and activation matching for a model learned on three data points with gradient descent.  $F$  is chosen to be the output of the three gaussian components ( $G_1, G_2, G_3$ ). Note that the perturbation is not  $\varepsilon$ -normalised for clearer visualisation. We highlight two observations:

- The UNI path is monotonous, of low-curvature and proximal. Conversely, the path to the random baseline is long, non-monotonous, and goes through several zones of high second derivative.
- Optimizing KL divergence on  $(G_1, G_2, G_3)$  produces a better baseline. Figure 6b visualises the unlearning objective (i.e. the target probability after unlearning), which gives four points of intersection with the base model ( $a, b, c$  and  $d$ ). By constraining proximity of the baseline with the  $\varepsilon$  parameter, we restrain the optima found by gradient descent (on the global probability) to the closest two points  $a$  and  $b$ . UNI is then able yield the more optimal of the two, by optimising on each gaussian output. In fact, the idea of activation matching is to satisfy the crucial weak dependence property for conformal path attribution (Akhtar & Jalwana, 2023). Since modern ReLU networks have decision boundaries representable as piecewise linear functions (Xiong et al., 2020), activation matching supervises the baseline to use the same (piecewise linear) weights. In our case, we want to find a baseline for which  $G_1$  and  $G_2$  do not play a role, which is not the case for  $a$ . This is why Algorithm 1 optimises on  $F$  and not on  $F_c$ .

Finally,  $\varepsilon$  normalisation serves to regularise baseline GD learning and account for pathological cases where the locally shortest path would lead to further intersections than the closest one.

## 5 EXPERIMENTS

We experiment on ImageNet-1K (Deng et al., 2009), ImageNet-C (Hendrycks & Dietterich, 2019) and compare against various path-based and gradient-based attribution methods. This includes IG (Sundararajan et al., 2017), BlurIG (Xu et al., 2020), GIG (Kapishnikov et al., 2021), AGI (Pan et al., 2021), GBP (Springenberg et al., 2014) and DeepLIFT (Shrikumar et al., 2016). We consider a diverse set of pre-trained computer vision backbone models (Paszke et al., 2019), including ResNet-18 (He et al., 2016), EfficientNet-v2-small (Tan & Le, 2021), ConvNeXt-Tiny (Liu et al., 2022), VGG-16-bn (Simonyan & Zisserman, 2015), ViT-B\_16 (Dosovitskiy et al., 2020) and Swin-Transformer-Tiny

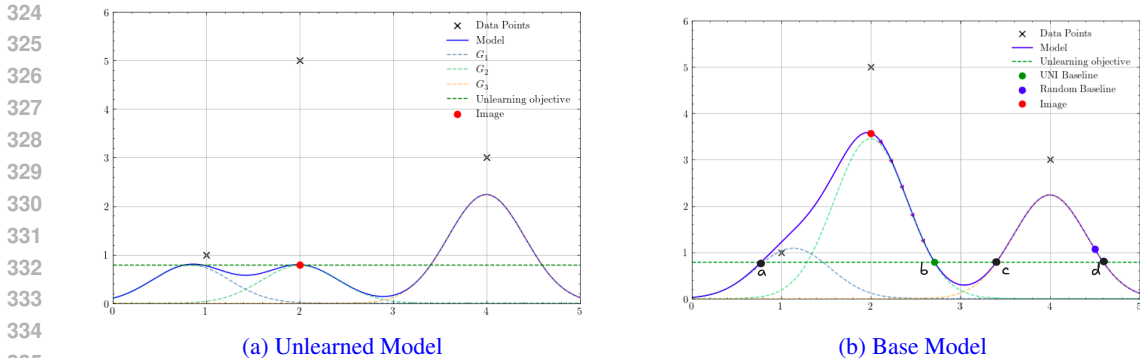


Figure 6: UNI baseline on a Gaussian mixture model of three gaussians  $G_1, G_2, G_3$ , each of fixed variance, parametrised by their mean and a scaling factor. (b) shows the model trained on the three datapoints (1,1), (2, 5) and (4, 3), while (a) shows the model after one gradient ascent step on the datapoint (2, 5). The path between UNI Baseline and the image is highlighted by arrows in (b).

(Liu et al., 2021). Unless otherwise specified, we use the following hyperparameters: unlearning step size  $\eta = 1$ ;  $l_2$  PGD with  $T = 10$  steps, a budget of  $\varepsilon = 0.25$ , step size  $\mu = 0.1$ ; Riemann approximation with  $B = 15$  steps. We further extend UNI to the NLP domain, to interpret generative language models using activation patching (Heimersheim & Nanda, 2024). UNI complements activation matching by computing a stable baseline without trading off attribution scalability, as observed in Appendix Table 8 and Figure 9. Our results verify UNI’s high faithfulness, stability and robustness.

Table 2: *MuFidelity scores* measure the correlation between a subset of pixels’ impact on the output (*i.e.* change in predictive confidence) and assigned saliency scores. Since attribution methods can yield strong positive or negative correlations, we report the absolute scores.

	UNI	IG	BlurIG	GIG	AGI	GBP	DeepLIFT
ResNet-18	<b>.12</b> $\pm$ .124	.06 $\pm$ .068	.07 $\pm$ .076	.07 $\pm$ .080	.10 $\pm$ .110	.09 $\pm$ .094	.08 $\pm$ .082
EfficientNetv2s	<b>.06</b> $\pm$ .046	.05 $\pm$ .043	.05 $\pm$ .044	.05 $\pm$ .044	.06 $\pm$ .045	.05 $\pm$ .043	.05 $\pm$ .043
ConvNeXt-Tiny	.16 $\pm$ .115	.11 $\pm$ .086	.15 $\pm$ .121	<b>.18</b> $\pm$ .149	.17 $\pm$ .131	.09 $\pm$ .072	.11 $\pm$ .084
VGG-16-bn	<b>.18</b> $\pm$ .141	.08 $\pm$ .066	.09 $\pm$ .076	.13 $\pm$ .108	.14 $\pm$ .104	.13 $\pm$ .108	.10 $\pm$ .082
ViT-B_16	<b>.15</b> $\pm$ .114	.10 $\pm$ .074	.10 $\pm$ .077	.11 $\pm$ .079	.14 $\pm$ .104	.09 $\pm$ .070	.10 $\pm$ .072
Swin-T-Tiny	<b>.13</b> $\pm$ .100	.09 $\pm$ .071	.12 $\pm$ .102	.12 $\pm$ .104	.13 $\pm$ .102	.09 $\pm$ .069	.10 $\pm$ .076

### 5.1 FAITHFULNESS

We report *MuFidelity scores* (Bhatt et al., 2021), *i.e.* the faithfulness of an attribution function  $\mathcal{A}$ , to a model  $F$ , at a sample  $x$ , for a subset of features of size  $|S|$ , given by  $\mu_f(F, \mathcal{A}; x) = \text{corr}_{S \in \binom{[d]}{|S|}} (\sum_{i \in S} \mathcal{A}(i, F, c, x), F_c(x) - F_c(x_{[x_s=\bar{x}_s]}))$ . We record the (absolute) correlation coefficient between a randomly sampled subset of pixels and their attribution scores. In line with open source exemplars (Fel et al., 2022a), we set  $|S|$  to be 25% of the total pixel count (slightly higher than the referenced 20%) as is required to adjust for ImageNet’s complexity and for obtaining less noisy measurements across all baseline methods. As from Table 2, UNI outperforms other methods across all settings but one, indicating high faithfulness. We supplement these numbers with visual comparisons in Appendix Figures 14, 15, 16, 17, 18, 19 against IG (black and noised baselines), BlurIG, GIG, AGI, GBP, DeepLift. Furthermore, we report deletion and insertion scores (Petsiuk et al., 2018)—a causally-motivated evaluation metric for interpretability methods—which measures the decrease (deletion) or increase (insertion) of a model’s output confidence as salient pixels are removed (from the original image) or inserted (into a featureless baseline). A steep drop in model confidence under pixel deletion results in a desirable and small area under the curve (AUC) score; a sharp rise under pixel insertion results in a desirably large AUC. Salient pixels are removed in descending order of importance, as identified by the tested interpretability method. We evaluate with a step size of 10% and average over 10,000 random image samples, where at each step, the next-10% most salient pixels are removed or inserted for inference. UNI reliably identifies pixels which are crucial for sample classification, achieving marked improvements especially in insertion AUC scores.

Table 3: *Deletion AUC* ↓ measures how confidence drops as pixels are removed (*lower = better*).

	UNI	IG	BlurIG	GIG	AGI	GBP	DeepLIFT
ResNet-18	<b>.06</b> ±.128	.10 ±.174	.27 ±.252	.11 ±.150	.13 ±.147	.08 ±.160	.13 ±.165
EfficientNetv2s	.19 ±.212	.26 ±.217	.50 ±.158	.19 ±.216	<b>.18</b> ±.207	.23 ±.163	.27 ±.215
ConvNeXt-Tiny	<b>.11</b> ±.139	.16 ±.164	.46 ±.172	.21 ±.160	.17 ±.123	.16 ±.099	.21 ±.162
VGG-16-bn	<b>.08</b> ±.143	.12 ±.181	.18 ±.241	.10 ±.163	.14 ±.178	.14 ±.194	.12 ±.186
ViT-B_16	.14 ±.185	.22 ±.207	.60 ±.166	.17 ±.190	<b>.13</b> ±.152	.23 ±.141	.17 ±.189
Swin-T-Tiny	<b>.13</b> ±.181	.22 ±.217	.47 ±.174	.22 ±.207	.21 ±.172	.21 ±.123	.23 ±.207

Table 4: *Insertion AUC* ↑ measures how confidence rises as pixels are inserted (*higher = better*).

	UNI	IG	BlurIG	GIG	AGI	GBP	DeepLIFT
ResNet-18	<b>.64</b> ±.138	.26 ±.045	.34 ±.131	.36 ±.048	.56 ±.068	.11 ±.066	.18 ±.042
EfficientNetv2s	<b>.64</b> ±.227	.38 ±.127	.51 ±.283	.37 ±.138	.38 ±.204	.23 ±.192	.37 ±.137
ConvNeXt-Tiny	<b>.63</b> ±.231	.21 ±.114	.40 ±.252	.56 ±.122	.52 ±.088	.22 ±.160	.17 ±.162
VGG-16-bn	<b>.56</b> ±.335	.37 ±.061	.31 ±.274	.38 ±.071	.47 ±.078	.26 ±.057	.17 ±.056
ViT-B_16	<b>.71</b> ±.237	.32 ±.107	.59 ±.292	.28 ±.125	.43 ±.089	.35 ±.172	.28 ±.123
Swin-T-Tiny	<b>.68</b> ±.245	.28 ±.145	.63 ±.282	.26 ±.153	.25 ±.156	.31 ±.202	.26 ±.152

### 5.2 ROBUSTNESS

Next, we evaluate UNI’s robustness to fragility adversarial attacks on model interpretations. Following Ghorbani et al. (2019a), we design norm-bounded attacks to maximise the disagreement in attributions whilst constraining that the prediction label remains unchanged. We consider a standard  $l_\infty$  attack designed with FGSM (Goodfellow et al., 2014), with perturbation budget  $\epsilon_f = 8/255$ .

$$\delta_f^* = \arg \max_{\|\delta_f\|_p \leq \epsilon_f} \frac{1}{d_X} \sum_{i=1}^{d_X} d(\mathcal{A}(i, F, c, x), \mathcal{A}(i, F, c, x + \delta_f)) \tag{3}$$

subject to  $\arg \max_{c'} F_{c'}(x) = \arg \max_c F_{c'}(x + \delta_f) = c$

We report robustness results using 2 distance measures—Spearman correlation coefficient in Table 5 and top- $k$  pixel intersection score in Table 6—pre and post attack. While other methods like DeepLIFT (DL), BlurIG, Integrated Gradients (IG) are misled to output irrelevant feature saliencies, UNI robustly maintains attribution consistency and achieves the lowest attack attribution disagreement scores (before and after FGSM attacks) for both metrics.

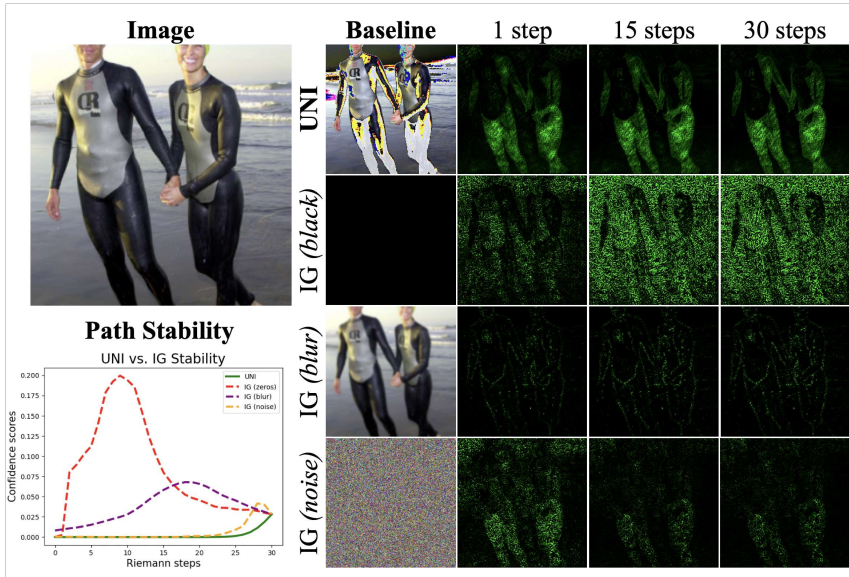


Figure 7: UNI path features monotonically increase in output confidence when interpolating from baseline to input. This eliminates instability and inconsistency problems caused by extrema and turning points along the Riemann approximation path, which is present in other methods.



### 5.3 STABILITY

We compare UNI and other methods’ sensitivity to Riemann approximation noise, which manifests in visual artefacts and misattribution of salient features. As seen from Figures 7, 10, UNI reliably finds unlearned, “featureless” baselines for consistent attribution, regardless of the number of approximation steps  $B \in \{1, 15, 30\}$ . This is due to the low geodesic curvature of  $\gamma^{\text{UNI}}$ , which approximately minimises the local distance between points used in Riemann approximation.

Table 5: *Robustness*: Spearman’s correlation coefficient. Higher scores indicate better path consistency pre/post FGSM attacks.

	UNI	IG	BlurIG	SG	DeepL
ResNet-18	<b>.271</b>	.088	.084	.014	.139
Eff-v2-s	<b>.302</b>	.009	.076	.008	.018
ConvNeXt-T	<b>.292</b>	.010	.127	.011	.012
VGG-16-bn	<b>.290</b>	.143	.098	.014	.108
ViT-B-16	<b>.319</b>	.018	.066	.023	.023
SwinT	<b>.271</b>	.088	.084	.014	.139

Table 6: *Robustness*: Top-1000 pixel intersection. Higher percentages indicate better attribution reliability pre/post FGSM attacks.

	UNI	IG	BlurIG	SG	DeepL
ResNet-18	<b>37.3</b>	20.0	25.3	18.2	24.8
Eff-v2-s	<b>39.4</b>	17.4	23.3	18.6	18.0
ConvNeXt-T	<b>34.8</b>	15.0	26.2	16.7	15.1
VGG-16-bn	<b>35.7</b>	25.5	25.3	18.8	25.2
ViT-B-16	<b>40.7</b>	17.1	21.7	19.6	17.2
SwinT	<b>37.3</b>	20.0	25.3	18.2	24.8

## 6 RELATED WORK

**Machine unlearning.** We draw inspiration from the high-level principle of unlearning, which concerns the targeted “forgetting” of a data-point for a trained model, by localising relevant information stored in network weights and introducing updates or perturbations (Bourtole et al., 2021). Formally, machine unlearning can be divided into exact and approximate unlearning (Nguyen et al., 2022). Exact unlearning seeks indistinguishability guarantees for output and weight distributions, between a model not trained on a sample and one that has unlearned said sample (Ginart et al., 2019; Thudi et al., 2022; Brophy & Lowd, 2021). However, provable exact unlearning is only achieved under full re-training, which can be computationally infeasible. Hence, approximate unlearning was proposed stemming from  $\epsilon$ -differential privacy (Dwork, 2011) and certified removal mechanisms (Guo et al., 2020; Golatkar et al., 2020). The former guarantees unlearning for  $\epsilon = 0$ , *i.e.* the sample has null influence on the decision function; the latter unlearns with first/second order gradient updates, achieving max-divergence bounds for single unlearning samples. Unlearning naturally lends itself to path-based attribution, to localise then delete information in the weight space, for the purposes of defining an “unlearned” activation. This “unlearned” activation can be used to match the corresponding, “featureless” input, where salient features have been deleted during the unlearning process. While the connection to interpretability is new, a few recent works intriguingly connect machine unlearning to the task of debiasing classification models during training and evaluation (Chen et al., 2024; Kim et al., 2019; Bevan & Atapour-Abarghouei, 2022).

**Perturbative methods.** Perturbative methods perturb inputs to change and explain outputs (Sculley et al., 2015), including LIME (Ribeiro et al., 2016), SHAP, KernelSHAP and GradientSHAP (Lundberg et al.), **RKHS-SHAP** (Chau et al., 2022), **ConceptSHAP** (Yeh et al., 2020), InterSHAP (Janzing et al., 2020), and DiCE (Kommiya Mothilal et al., 2021). LIME variants optimise a simulator of minimal functional complexity able to match the black-box model’s local behaviour for a given input-label pair. SHAP (Lundberg et al.) consolidates LIME, DeepLIFT (Shrikumar et al., 2016), Layerwise Relevance Propagation (LRP) (Montavon et al., 2019) under the general, game-theoretic framework of additive feature attribution methods. For this framework, they outline the desired properties of local accuracy, missingness, consistency; they propose SHAP values as a feature importance measure which satisfies these properties under mild assumptions to generate model-agnostic explanations. However, such methods fail to give a global insight of the model’s decision function and are highly unstable due to the reliance on local perturbations (Fel et al., 2022b). Bordt et al. (2022) show that this leads to variability, inconsistency and unreliability in generated explanations, where different methods give incongruent explanations which cannot be acted on. [Recent works have made considerable progress, including RISE \(Petsiuk et al., 2018\), which strives to causally explain model predictions by approximating the necessary saliency of pixels through random masking; Sobol \(Fel et al., 2021\), which adapts Sobol indices for perturbation masks towards](#)

variance-based sensitivity analysis; and FORGrad (Muzellec et al., 2023), which filters out high-frequency gradient noise induced by white-box methods (and network pooling or striding operations) and which can be complementarily applied to further UNI’s explanation faithfulness and efficiency.

**Backpropagative methods.** Beginning with simple gradients (Erhan et al., 2009; Simonyan et al., 2013), this family of methods—also, LRP (Montavon et al., 2019), DeepLIFT (Shrikumar et al., 2016), DeConvNet (Zeiler & Fergus, 2014), Guided Backpropagation (Springenberg et al., 2014) and GradCAM (Selvaraju et al., 2017)—leverages gradients of the output *w.r.t.* the input to proportionally project predictions back to the input space, for some given neuron activity of interest. Gradients of neural networks are, however, highly noisy and locally sensitive – they can only crudely localise salient feature regions. While this issue is partially remedied by SmoothGrad (Smilkov et al., 2017), we still observe that gradient-based saliency methods have higher sample complexity for generalisation than normal supervised training (Choi & Farnia, 2024) and often yield inconsistent attributions for unseen images at test time.

**Path-based attribution.** This family of post-hoc attributions is attractive due to its grounding in cooperative game-theory (Friedman, 2004). It comprises Integrated Gradients (Sundararajan et al., 2017), Adversarial Gradient Integration (Pan et al., 2021), Expected Gradients (EG) (Erion et al., 2021), Guided Integrated Gradients (GIG) (Kapishnikov et al., 2021) and BlurIG (Xu et al., 2020). Path attribution typically relies on a baseline – a “vanilla” image devoid of features; a path—an often linear path from the featureless baseline to the target image—along which the path integral is computed for every pixel. Granular control over the attribution process comes with difficulties of defining an unambiguously featureless baseline (for each (model, image) pair) (Sturmfels et al., 2020) and then defining a reliable path of increasing label confidence without intermediate inflection points (Akhtar & Jalwana, 2023). To measure the discriminativeness of features identified by attribution methods and the extent to which model predictions depend on them, experimental benchmarks and metrics such as ROAR (Hooker et al., 2019), DiffRAOR (Shah et al., 2021), deletion/insertion score (Petsiuk et al., 2018), the Hilbert-Schmidt independence criterion (HSIC) (Novello et al., 2022) and the Pointing Game (Zhang et al., 2018a) have been proposed.

## 7 CONCLUSION

In this work, we formally discuss the limitations of current path-attribution frameworks, outline a new principle for optimising baseline and path features, as well as introduce the UNI algorithm for unlearning-based neural interpretations. We empirically show that present reliance on static baselines imposes undesirable post-hoc biases which are alien to the model’s decision function. We account for and mitigate various infidelity, inconsistency and instability issues in path-attribution by defining principled baselines and conformant path features. UNI leverages insights from unlearning to eliminate task-salient features and mimic baseline activations in the “absence of signal”. It discovers low-curvature, stable paths with monotonically increasing output confidence, which preserves the completeness axiom necessary for path attribution. We visually, numerically and formally establish the utility of UNI as a means to compute robust, meaningful and debiased image attributions.

The contributions of UNI extend beyond the presented method and analyses, towards investigating machine unlearning as a tool for white-box interpretability. Unlearning at different granularities allows us to audit the various levels of a model’s learned feature hierarchy. In this work, we illustrate how first-order, sample-wise unlearning can identify salient input features important for a single prediction. A promising future direction involves interpreting higher-level, semantically complex concepts required to learn a task or fit a data distribution, by instead unlearning a set of concept-clustered exemplars. It is also of interest to delve into how interpretability methods impose additional assumptions onto trained models, prompting questions such as how to best design and align the correct interpretability method for a given model; how to use attribution methods to compare and contrast the inductive biases of different network architectures, of models trained with robust versus non-robust objectives, of models trained using different equivariant data augmentation strategies. Further technical extensions to UNI include going beyond first-order approximate unlearning towards certified, second-order machine unlearning; as well as granular investigations of how the baseline definition, model’s robustness and model’s inductive biases exert influence on path attribution results.

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

## 8 REPRODUCIBILITY STATEMENT

We report experimental details on the considered datasets, models and baselines, and hyperparameter values for UNI in Section 5. Evaluation metrics of MuFidelity scores (measuring faithfulness), Spearman correlation coefficient (measuring path monotonicity and attributions robustness), pixel intersection score (measuring attributions robustness) are detailed and consistent with existing literature. Theoretical arguments are substantiated by quantitative and qualitative results.

## REFERENCES

- 594  
595  
596 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
597 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.  
598 *arXiv preprint arXiv:2303.08774*, 2023.
- 599  
600 Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim.  
601 Sanity checks for saliency maps. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,  
602 N. Cesa-Bianchi, and R. Garnett (eds.), *NeurIPS*, volume 31. Curran Associates, Inc.,  
603 2018. URL [https://proceedings.neurips.cc/paper\\_files/paper/2018/  
604 file/294a8ed24blad22ec2e7efea049b8737-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/294a8ed24blad22ec2e7efea049b8737-Paper.pdf).
- 605  
606 Naveed Akhtar and Mohammad A. A. K. Jalwana. Towards credible visual model interpretation with  
607 path attribution. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan  
608 Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine  
609 Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 439–457. PMLR, 23–29  
610 Jul 2023. URL <https://proceedings.mlr.press/v202/akhtar23a.html>.
- 611  
612 David Alvarez Melis and Tommi Jaakkola. Towards robust interpretability with self-explaining neural  
613 networks. *NeurIPS*, 31, 2018.
- 614  
615 Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Towards better understanding  
616 of gradient-based attribution methods for deep neural networks. In *ICLR*, 2018. URL <https://openreview.net/forum?id=Sy21R9JAW>.
- 617  
618 Anthropic. Introducing the next generation of claude. Mar 2024. URL [https://www.  
619 anthropic.com/news/claude-3-family](https://www.anthropic.com/news/claude-3-family).
- 620  
621 Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. A diagnostic  
622 study of explainability techniques for text classification. In *Proceedings of the 2020 Conference on  
623 Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3256–3274, 2020.
- 624  
625 David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and  
626 Klaus-Robert Müller. How to explain individual classification decisions. *Journal of Machine  
627 Learning Research*, 11(61):1803–1831, 2010. URL [http://jmlr.org/papers/v11/  
628 baehrens10a.html](http://jmlr.org/papers/v11/baehrens10a.html).
- 629  
630 David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection:  
631 Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference  
632 on computer vision and pattern recognition*, pp. 6541–6549, 2017.
- 633  
634 Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Yuval Noah Harari,  
635 Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, Gillian Hadfield, et al. Managing ai risks in an era  
636 of rapid progress. *arXiv preprint arXiv:2310.17688*, 2023.
- 637  
638 Peter Bevan and Amir Atapour-Abarghouei. Skin deep unlearning: Artefact and instrument debiasing  
639 in the context of melanoma classification. In *International Conference on Machine Learning*, pp.  
640 1874–1892. PMLR, 2022.
- 641  
642 Umang Bhatt, Adrian Weller, and José MF Moura. Evaluating and aggregating feature-based model  
643 explanations. In *Proceedings of the Twenty-Ninth International Conference on International Joint  
644 Conferences on Artificial Intelligence*, pp. 3016–3022, 2021.
- 645  
646 Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan,  
647 Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron,  
Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models  
across training and scaling, 2023. URL <https://arxiv.org/abs/2304.01373>.
- Moritz Bohle, Mario Fritz, and Bernt Schiele. Convolutional dynamic alignment networks for  
interpretable classifications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
Pattern Recognition (CVPR)*, pp. 10029–10038, June 2021.

- 648 Moritz Böhle, Mario Fritz, and Bernt Schiele. B-cos networks: Alignment is all we need for  
649 interpretability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
650 Recognition*, pp. 10329–10338, 2022.
- 651 Sebastian Bordt, Michèle Finck, Eric Raidl, and Ulrike von Luxburg. Post-hoc explanations fail to  
652 achieve their purpose in adversarial contexts. In *Proceedings of the 2022 ACM Conference on  
653 Fairness, Accountability, and Transparency*, FAccT '22, pp. 891–905, New York, NY, USA, 2022.  
654 Association for Computing Machinery. ISBN 9781450393522. doi: 10.1145/3531146.3533153.  
655 URL <https://doi.org/10.1145/3531146.3533153>.
- 656 Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers,  
657 Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE Symposium  
658 on Security and Privacy (SP)*, pp. 141–159. IEEE Computer Society, 2021.
- 659 Wieland Brendel and Matthias Bethge. Approximating cnns with bag-of-local-features models works  
660 surprisingly well on imagenet. In *International Conference on Learning Representations*, 2018.
- 661 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe  
662 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video  
663 generation models as world simulators. 2024. URL [https://openai.com/research/  
664 video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 665 Jonathan Brophy and Daniel Lowd. Machine unlearning for random forests. In *International  
666 Conference on Machine Learning*, pp. 1092–1104. PMLR, 2021.
- 667 Siu Lun Chau, Robert Hu, Javier Gonzalez, and Dino Sejdinovic. RKHS-SHAP: Shapley values for  
668 kernel methods. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.),  
669 *Advances in Neural Information Processing Systems*, 2022. URL [https://openreview.  
670 net/forum?id=gnc2VJHXmSG](https://openreview.net/forum?id=gnc2VJHXmSG).
- 671 Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K  
672 Su. This looks like that: Deep learning for interpretable image recognition. In  
673 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.),  
674 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.,  
675 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/2019/  
676 file/adf7ee2dcf142b0e11888e72b43fcb75-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/adf7ee2dcf142b0e11888e72b43fcb75-Paper.pdf).
- 677 Ruizhe Chen, Jianfei Yang, Huimin Xiong, Jianhong Bai, Tianxiang Hu, Jin Hao, Yang Feng,  
678 Joey Tianyi Zhou, Jian Wu, and Zuozhu Liu. Fast model debias with machine unlearning. *Advances  
679 in Neural Information Processing Systems*, 36, 2024.
- 680 Ching Lam Choi and Farzan Farnia. On the generalization of gradient-based neural network  
681 interpretations, 2024. URL <https://openreview.net/forum?id=EwAGztBkJ6>.
- 682 Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-  
683 Alonso. Towards automated circuit discovery for mechanistic interpretability, 2023. URL [https://  
684 arxiv.org/abs/2304.14997](https://arxiv.org/abs/2304.14997).
- 685 Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse  
686 autoencoders find highly interpretable features in language models, 2023. URL [https://  
687 arxiv.org/abs/2309.08600](https://arxiv.org/abs/2309.08600).
- 688 Google Deepmind. Veo: Our most capable generative video model. May 2024. URL [https://  
689 deepmind.google/technologies/veo/](https://deepmind.google/technologies/veo/).
- 690 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
691 hierarchical image database. In *CVPR*, pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 692 Ann-Kathrin Dombrowski, Maximillian Alber, Christopher Anders, Marcel Ackermann, Klaus-  
693 Robert Müller, and Pan Kessel. Explanations can be manipulated and geometry is to blame.  
694 *Advances in neural information processing systems*, 32, 2019.
- 695 Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning.  
696 *arXiv preprint arXiv:1702.08608*, 2017.
- 697  
698  
699  
700  
701

- 702 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
703 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image  
704 is worth 16x16 words: Transformers for image recognition at scale. In *International Conference  
705 on Learning Representations*, 2020.
- 706 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
707 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
708 *arXiv preprint arXiv:2407.21783*, 2024.
- 709 Cynthia Dwork. *Differential Privacy*, pp. 338–340. Springer US, Boston, MA, 2011. ISBN 978-1-  
710 4419-5906-5. doi: 10.1007/978-1-4419-5906-5\_752. URL [https://doi.org/10.1007/  
711 978-1-4419-5906-5\\_752](https://doi.org/10.1007/978-1-4419-5906-5_752).
- 712 Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer  
713 features of a deep network. 2009.
- 714 Gabriel Erion, Joseph D Janizek, Pascal Sturmfels, Scott M Lundberg, and Su-In Lee. Improving  
715 performance of deep learning models with axiomatic attribution priors and expected gradients.  
716 *Nature machine intelligence*, 3(7):620–631, 2021.
- 717 European Commission. Laying down harmonised rules on artificial intelligence  
718 (artificial intelligence act) and amending certain union legislative acts. *CNECT*, Apr  
719 2021. URL [https://digital-strategy.ec.europa.eu/en/library/  
720 proposal-regulation-laying-down-harmonised-rules-artificial-intelligence](https://digital-strategy.ec.europa.eu/en/library/proposal-regulation-laying-down-harmonised-rules-artificial-intelligence).
- 721 Thomas Fel, Rémi Cadène, Mathieu Chalvidal, Matthieu Cord, David Vigouroux, and Thomas  
722 Serre. Look at the variance! efficient black-box explanations with sobol-based sensitivity analysis.  
723 *Advances in neural information processing systems*, 34:26005–26014, 2021.
- 724 Thomas Fel, Lucas Hervier, David Vigouroux, Antonin Poche, Justin Plakoo, Remi Cadene, Mathieu  
725 Chalvidal, Julien Colin, Thibaut Boissin, Louis Bethune, Agustin Picard, Claire Nicodeme, Laurent  
726 Gardes, Gregory Flandin, and Thomas Serre. Xplique: A deep learning explainability toolbox.  
727 *Workshop on Explainable Artificial Intelligence for Computer Vision (CVPR)*, 2022a.
- 728 Thomas Fel, David Vigouroux, Rémi Cadène, and Thomas Serre. How good is your explanation?  
729 algorithmic stability measures to assess the quality of explanations for deep neural networks. In  
730 *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 720–730,  
731 2022b.
- 732 Jiahai Feng and Jacob Steinhardt. How do language models bind entities in context?, 2024. URL  
733 <https://arxiv.org/abs/2310.17191>.
- 734 Eric J. Friedman. Paths and consistency in additive cost sharing. *Int. J. Game Theory*, 32(4):501–518,  
735 aug 2004. ISSN 0020-7276. doi: 10.1007/s001820400173. URL [https://doi.org/10.  
736 1007/s001820400173](https://doi.org/10.1007/s001820400173).
- 737 Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. In  
738 *AAAI*, volume 33, pp. 3681–3688, 2019a.
- 739 Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based  
740 explanations. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett  
741 (eds.), *NeurIPS*, volume 32, 2019b. URL [https://proceedings.neurips.cc/paper\\_  
742 files/paper/2019/file/77d2afcb31f6493e350fca61764efb9a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/77d2afcb31f6493e350fca61764efb9a-Paper.pdf).
- 743 Antonio Ginart, Melody Guan, Gregory Valiant, and James Y Zou. Making ai forget you: Data  
744 deletion in machine learning. *Advances in neural information processing systems*, 32, 2019.
- 745 Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net:  
746 Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer  
747 Vision and Pattern Recognition*, pp. 9304–9312, 2020.
- 748 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial  
749 examples. *arXiv preprint arXiv:1412.6572*, 2014.

- 756 Google. Introducing paligemma, gemma 2, and an upgraded responsible ai  
757 toolkit. May 2024. URL [https://developers.googleblog.com/en/  
758 gemma-family-and-toolkit-expansion-io-2024/](https://developers.googleblog.com/en/gemma-family-and-toolkit-expansion-io-2024/).  
759
- 760 Chuan Guo, Tom Goldstein, Awni Hannun, and Laurens Van Der Maaten. Certified data removal  
761 from machine learning models. In *Proceedings of the 37th International Conference on Machine  
762 Learning*, pp. 3832–3842, 2020.
- 763 Michael Hanna, Ollie Liu, and Alexandre Variengien. How does gpt-2 compute greater-than?:  
764 Interpreting mathematical abilities in a pre-trained language model, 2023. URL [https://  
765 arxiv.org/abs/2305.00586](https://arxiv.org/abs/2305.00586).  
766
- 767 Michael Hanna, Sandro Pezzelle, and Yonatan Belinkov. Have faith in faithfulness: Going beyond  
768 circuit overlap when finding model mechanisms, 2024. URL [https://arxiv.org/abs/  
769 2403.17806](https://arxiv.org/abs/2403.17806).
- 770 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
771 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
772 pp. 770–778, 2016.  
773
- 774 Stefan Heimersheim and Neel Nanda. How to use and interpret activation patching, 2024. URL  
775 <https://arxiv.org/abs/2404.15255>.
- 776 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common  
777 corruptions and perturbations. In *International Conference on Learning Representations*, 2019.  
778
- 779 Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability  
780 methods in deep neural networks. *NeurIPS*, 32, 2019.
- 781 Dominik Janzing, Lenon Minorics, and Patrick Bloebaum. Feature relevance quantification in  
782 explainable ai: A causal problem. In *Proceedings of the Twenty Third International Conference on  
783 Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pp.  
784 2907–2916. PMLR, 26–28 Aug 2020. URL [https://proceedings.mlr.press/v108/  
785 janzing20a.html](https://proceedings.mlr.press/v108/janzing20a.html).  
786
- 787 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
788 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
789 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 790 Andrei Kapishnikov, Subhashini Venugopalan, Besim Avci, Ben Wedin, Michael Terry, and Tolga  
791 Bolukbasi. Guided integrated gradients: An adaptive path method for removing noise. In *CVPR*,  
792 pp. 5050–5058, 2021.  
793
- 794 Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al.  
795 Interpretability beyond feature attribution: Quantitative testing with concept activation vectors  
796 (tcav). In *International conference on machine learning*, pp. 2668–2677. PMLR, 2018.
- 797 Byungju Kim, Hyunwoo Kim, Kyungsu Kim, Sungjin Kim, and Junmo Kim. Learning not to learn:  
798 Training deep neural networks with biased data. In *Proceedings of the IEEE/CVF conference on  
799 computer vision and pattern recognition*, pp. 9012–9020, 2019.  
800
- 801 Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and  
802 Percy Liang. Concept bottleneck models. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of  
803 the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine  
804 Learning Research*, pp. 5338–5348. PMLR, 13–18 Jul 2020. URL [https://proceedings.  
805 mlr.press/v119/koh20a.html](https://proceedings.mlr.press/v119/koh20a.html).
- 806 Ramaravind Kommiya Mothilal, Divyat Mahajan, Chenhao Tan, and Amit Sharma. Towards  
807 unifying feature attribution and counterfactual explanations: Different means to the same  
808 end. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES  
809 ’21, pp. 652–663. Association for Computing Machinery, 2021. ISBN 9781450384735. doi:  
10.1145/3461702.3462597. URL <https://doi.org/10.1145/3461702.3462597>.

- 810 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*,  
811 2023. URL <https://openreview.net/forum?id=w0H2xGH1kw>.
- 812
- 813 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.  
814 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*  
815 *IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 816 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie.  
817 A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and*  
818 *pattern recognition*, pp. 11976–11986, 2022.
- 819
- 820 Scott M. Lundberg, Gabriel Erion, Hugh Chen, Alex DeGrave, Jordan M. Prutkin, Bala Nair, Ronit  
821 Katz, Jonathan Himmelfarb, Nisha Bansal, and Su-In Lee. From local explanations to global  
822 understanding with explainable ai for trees. *Nature machine intelligence*, 2. doi: 10.1038/  
823 s42256-019-0138-9. URL <https://par.nsf.gov/biblio/10167481>.
- 824 Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller.  
825 Sparse feature circuits: Discovering and editing interpretable causal graphs in language models,  
826 2024. URL <https://arxiv.org/abs/2403.19647>.
- 827
- 828 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual  
829 associations in gpt, 2023. URL <https://arxiv.org/abs/2202.05262>.
- 830 Meta. Introducing meta llama 3: The most capable openly available llm to date. Apr 2024. URL  
831 <https://ai.meta.com/blog/meta-llama-3/>.
- 832
- 833 Grégoire Montavon, Alexander Binder, Sebastian Lapuschkin, Wojciech Samek, and Klaus-Robert  
834 Müller. *Layer-Wise Relevance Propagation: An Overview*, pp. 193–209. Springer International  
835 Publishing, Cham, 2019. ISBN 978-3-030-28954-6. doi: 10.1007/978-3-030-28954-6\_10. URL  
836 [https://doi.org/10.1007/978-3-030-28954-6\\_10](https://doi.org/10.1007/978-3-030-28954-6_10).
- 837 Sabine Muzellec, Léo Andéol, Thomas Fel, Rufin VanRullen, and Thomas Serre. Gradient strikes  
838 back: How filtering out high frequencies improves explanations. *arXiv preprint arXiv:2307.09591*,  
839 2023.
- 840
- 841 Neel Nanda. Attribution patching: Activation patching at industrial scale, 2023.
- 842
- 843 Thanh Tam Nguyen, Thanh Trung Huynh, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin,  
844 and Quoc Viet Hung Nguyen. A survey of machine unlearning. *arXiv preprint arXiv:2209.02299*,  
845 2022.
- 846 Paul Novello, Thomas Fel, and David Vigouroux. Making sense of dependence: Efficient black-box  
847 explanations using dependence measure. *Advances in Neural Information Processing Systems*, 35:  
848 4344–4357, 2022.
- 849 Deng Pan, Xin Li, and Dongxiao Zhu. Explaining deep neural network models with adversarial  
850 gradient integration. In *IJCAI*, pp. 2876–2883. International Joint Conferences on Artificial  
851 Intelligence Organization, 8 2021. doi: 10.24963/ijcai.2021/396. URL [https://doi.org/](https://doi.org/10.24963/ijcai.2021/396)  
852 [10.24963/ijcai.2021/396](https://doi.org/10.24963/ijcai.2021/396). Main Track.
- 853
- 854 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
855 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas  
856 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,  
857 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-  
858 performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp.  
859 8024–8035. Curran Associates, Inc., 2019. URL [http://papers.neurips.cc/paper/](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)  
860 [9015-pytorch-an-imperative-style-high-performance-deep-learning-library.](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)  
861 [pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).
- 862 Vitali Petsiuk, Abir Das, and Kate Saenko. RISE: randomized input sampling for explanation of  
863 black-box models. In *BMVC*, pp. 151. BMVA Press, 2018. URL [http://bmvc2018.org/](http://bmvc2018.org/contents/papers/1064.pdf)  
[contents/papers/1064.pdf](http://bmvc2018.org/contents/papers/1064.pdf).



- 864 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
865 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.  
866
- 867 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the  
868 predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference*  
869 *on Knowledge Discovery and Data Mining*, KDD '16, pp. 1135–1144, New York, NY, USA, 2016.  
870 Association for Computing Machinery. ISBN 9781450342322. doi: 10.1145/2939672.2939778.  
871 URL <https://doi.org/10.1145/2939672.2939778>.
- 872 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
873 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF*  
874 *conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.  
875
- 876 Andrew Slavin Ross, Michael C. Hughes, and Finale Doshi-Velez. Right for the right reasons:  
877 Training differentiable models by constraining their explanations. In *Proceedings of the Twenty-*  
878 *Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pp. 2662–2670, 2017.  
879 doi: 10.24963/ijcai.2017/371. URL <https://doi.org/10.24963/ijcai.2017/371>.
- 880 D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay  
881 Chaudhary, Michael Young, Jean-François Crespo, and Dan Dennison. Hidden technical  
882 debt in machine learning systems. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and  
883 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran  
884 Associates, Inc., 2015. URL [https://proceedings.neurips.cc/paper\\_files/](https://proceedings.neurips.cc/paper_files/paper/2015/file/86df7dcfd896fcaf2674f757a2463eba-Paper.pdf)  
885 [paper/2015/file/86df7dcfd896fcaf2674f757a2463eba-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/86df7dcfd896fcaf2674f757a2463eba-Paper.pdf).
- 886 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh,  
887 and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based  
888 localization. In *CVPR*, pp. 618–626, 2017.  
889
- 890 Harshay Shah, Prateek Jain, and Praneeth Netrapalli. Do input gradients highlight discriminative  
891 features? *Advances in Neural Information Processing Systems*, 34:2046–2059, 2021.  
892
- 893 Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black  
894 box: Learning important features through propagating activation differences. *arXiv preprint*  
895 *arXiv:1605.01713*, 2016.
- 896 K Simonyan and A Zisserman. Very deep convolutional networks for large-scale image recognition.  
897 In *3rd International Conference on Learning Representations (ICLR 2015)*, 2015.  
898
- 899 Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:  
900 Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- 901 Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad:  
902 removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.  
903
- 904 Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for  
905 simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- 906 Suraj Srinivas and Francois Fleuret. Full-gradient representation for neural network visualization,  
907 2019.  
908
- 909 Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya Sachan. A mechanistic interpretation of  
910 arithmetic reasoning in language models using causal mediation analysis, 2023. URL <https://arxiv.org/abs/2305.15054>.  
911
- 912 Pascal Sturmfels, Scott Lundberg, and Su-In Lee. Visualizing the impact of feature attribution  
913 baselines. *Distill*, 2020. doi: 10.23915/distill.00022. <https://distill.pub/2020/attribution-baselines>.  
914
- 915 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Gradients of counterfactuals, 2016.  
916
- 917 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In  
*International conference on machine learning*, pp. 3319–3328. PMLR, 2017.

- 918 Aaquib Syed, Can Rager, and Arthur Conmy. Attribution patching outperforms automated circuit  
919 discovery. *arXiv preprint arXiv:2310.10348*, 2023.
- 920
- 921 Mingxing Tan and Quoc Le. Efficientnetv2: Smaller models and faster training. In *International  
922 conference on machine learning*, pp. 10096–10106. PMLR, 2021.
- 923
- 924 Anvith Thudi, Gabriel Deza, Varun Chandrasekaran, and Nicolas Papernot. Unrolling sgd:  
925 Understanding factors influencing machine unlearning. In *2022 IEEE 7th European Symposium on  
926 Security and Privacy (EuroS&P)*, pp. 303–319. IEEE, 2022.
- 927
- 928 Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Simas Sakenis, Jason  
929 Huang, Yaron Singer, and Stuart Shieber. Causal mediation analysis for interpreting neural nlp:  
The case of gender bias, 2020. URL <https://arxiv.org/abs/2004.12265>.
- 930
- 931 Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.  
932 Interpretability in the wild: a circuit for indirect object identification in gpt-2 small, 2022. URL  
933 <https://arxiv.org/abs/2211.00593>.
- 934
- 935 White House OSTP. Blueprint for an ai bill of rights: Making automated systems work for  
936 the american people. Oct 2022. URL [https://www.whitehouse.gov/wp-content/  
uploads/2022/10/Blueprint-for-an-AI-Bill-of-Rights.pdf](https://www.whitehouse.gov/wp-content/uploads/2022/10/Blueprint-for-an-AI-Bill-of-Rights.pdf).
- 937
- 938 H. Xiong, L. Huang, M. Yu, L. Liu, F. Zhu, and L. Shao. On the number of linear regions of  
939 convolutional neural networks, 2020. URL <https://arxiv.org/abs/2006.00978>.
- 940
- 941 Shawn Xu, Subhashini Venugopalan, and Mukund Sundararajan. Attribution in scale and space. In  
942 *CVPR*, June 2020.
- 943
- 944 Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Suggala, David I Inouye, and Pradeep K Ravikumar. On the  
945 (in) fidelity and sensitivity of explanations. *NeurIPS*, 32, 2019.
- 946
- 947 Chih-Kuan Yeh, Been Kim, Sercan Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar.  
948 On completeness-aware concept-based explanations in deep neural networks. *Advances in neural  
949 information processing systems*, 33:20554–20565, 2020.
- 950
- 951 Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In  
952 *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12,  
953 2014, Proceedings, Part I 13*, pp. 818–833. Springer, 2014.
- 954
- 955 Jianming Zhang, Sarah Adel Bargal, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff.  
956 Top-down neural attention by excitation backprop. *IJCV*, 126(10):1084–1102, 2018a.
- 957
- 958 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
959 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on  
960 computer vision and pattern recognition*, pp. 586–595, 2018b.
- 961
- 962 Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Interpretable basis decomposition for  
963 visual explanation. In *Proceedings of the European Conference on Computer Vision (ECCV)*,  
964 September 2018.
- 965
- 966
- 967
- 968
- 969
- 970
- 971

972 A APPENDIX

973  
974  
975  
976 A.1 VERIFYING ATTRIBUTION SPECIFICITY

977  
978 To verify that UNI computes explanations that are specific to each task–model–input triplet, we  
979 compare its saliency attributions across models for the same image input. Visually, we observe in  
980 Figure 8 that attributions differ significantly and even reflect the inductive biases of respective models  
981 (e.g. grid-like artefacts are present in ViT attributions whereas smoother attributions are computed  
982 for convolutional architectures). We further present numerical results in Table 7—LPIPS (Zhang  
983 et al., 2018b) scores reflect the dissimilarity/distance between the original image and the unlearned  
984 baseline; the percentage change in confidence scores reflect how the unlearned baseline effectively  
985 reduces predictive confidence (relative to the original input).  
986

987 Table 7: UNI computes different baselines for  
988 network architectures with different inductive  
989 biases on the same input, as seen from the drop  
990 in model confidence ( $\Delta\%$  Confidence) and image-  
991 baseline similarity scores (LPIPS<sub>vgg</sub>, LPIP<sub>alex</sub>).

	$\Delta\%$ Confidence	LPIPS <sub>vgg</sub>	LPIP <sub>alex</sub>
ResNet-18	-82.3%	.021 ± .025	.003 ± .005
Eff-v2-s	-76.9%	.025 ± .024	.004 ± .011
ConvNeXt-T	-95.1%	.018 ± .016	.002 ± .003
VGG-16-bn	-71.6%	.017 ± .020	.001 ± .002
ViT-B-16	-69.7%	.014 ± .015	.004 ± .007
SwinT	-84.6%	.014 ± .015	.002 ± .002

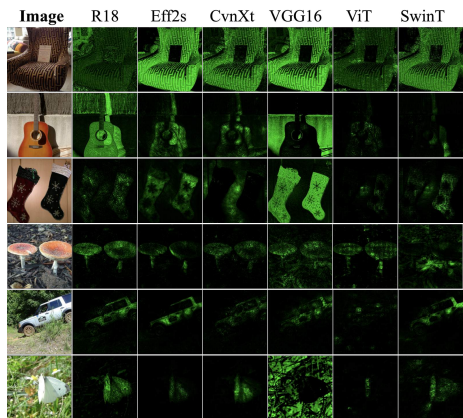


Figure 8: UNI computes different attributions to explain the predictions of each model.

1002  
1003  
1004  
1005 A.2 PRELIMINARY RESULTS ON NLP

1006  
1007 Table 8: Faithfulness: L<sub>2</sub>-Distance from activation patching to attribution patching results on the  
1008 residual stream (averaged over 100 samples).  
1009

	UNI	Random
Pythia-1b-v0	<b>3.12</b>	6.64
GPT2-medium	<b>15.26</b>	35.00
Llama-3.2-1B	<b>5.25</b>	10.17

1010 We extend the testing of our method to the case of Natural Language Processing (NLP). We choose  
1011 to test the application of UNI in the general framework of generative models (which includes  
1012 classification models), and attribution of not only inputs but more generally activations. Activation  
1013 patching (Heimersheim & Nanda, 2024) is one of the most widely used technique in Mechanistic  
1014 Interpretability, and more generally to study the properties of LLMs’ internals (Vig et al., 2020; Meng  
1015 et al., 2023; Wang et al., 2022; Feng & Steinhardt, 2024; Cunningham et al., 2023; Stolfo et al., 2023;  
1016 Hanna et al., 2023). This attribution method consists in analyzing a model’s output variation after  
1017 replacing its internal activations, following the equation:  
1018  
1019  
1020  
1021  
1022

$$\mathcal{A}_e^{\text{ACT}}(x) = F(x|e = e(x')) - F(x) \tag{4}$$

1023 where  $e$  denotes one activation in the model,  $x'$  a chosen baseline, and  $F$  a function of the model’s  
1024 output (usually the logit value of the maximum probability token of the model run on  $x$ ).  
1025

1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079

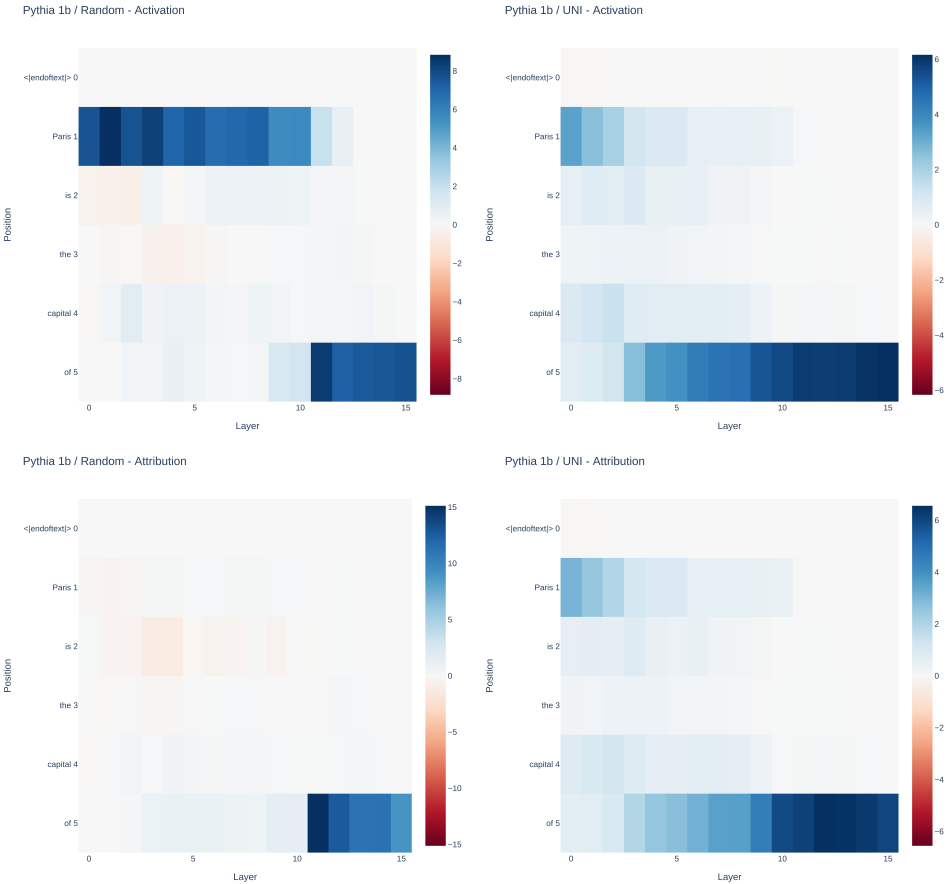


Figure 9: Visual comparison of attribution results for Activation vs. Attribution patching, with UNI versus Random baselines, on Pythia-1b-v0. Each cell shows the logit variation obtained by patching at a specific token and layer the residual stream of our baseline with the original activation.

Unfortunately activation patching is computationally costly, especially for purposes such as circuit discovery (Conmy et al., 2023). One of the main alternatives that solves the scalability problem is attribution patching (Nanda, 2023; Syed et al., 2023), which computes a first order Taylor approximation of Equation 4:

$$\mathcal{A}_e^{\text{ATTR}}(x) = (e(x) - e(x'))^T \nabla_e F(x') \tag{5}$$

for which the attributions for all of the activations can be computed at the same time (no patching of one single activation is performed). Despite its scalability, attribution patching suffers from a lack of faithfulness for causal interventions, mainly due to saturation and lack of linearity of the studied dependencies.

The analogy with integrated gradients seems quite striking, and indeed two recent works (to our knowledge) have tried to investigate the use of IG for more faithful attribution patching. While Marks et al. (2024) applies a very computationally complex version of such a method to small models, Hanna et al. (2024) proves the potential of IG-based attribution patching, while showing it still gets outperformed by activation patching.

We here provide a new UNI-based attribution method algorithm outputting faithful attributions while maintaining the scalability advantage of attribution patching. Mainly, we apply Algorithm 1 to compute a baseline  $x'$  that is then used to compute Equation 1:

$$\mathcal{A}_e^{\text{UNI-ATTR}}(x) = (e(x) - e(x'))^T \nabla_e F(\text{UNI}(x)) \tag{6}$$

where we take  $x$  to be the embedding of the input, to allow for continuous operations on it. Considering the known high faithfulness of activation patching (Hanna et al., 2024), we approximate faithfulness

of attributions computed from a baseline, as the  $L_2$ -distance between these attributions and the activation patching ones. The dataset used is a subset of 100 counterfactual prompts taken from Meng et al. (2023), and three different models are tested: Pythia-1b-v0 (Biderman et al., 2023), GPT2-medium (Radford et al., 2019) and Llama-3.2-1B (Dubey et al., 2024). The results can be seen in Table 8, and visuals in Figure 9. Note that no fine-tuning of UNI hyperparameters has been done, so that we expect even better results when adapting for each models. Unsurprisingly, decoding the baselines by shortest distance to the rows of the embedding matrix yields the same input, and a direct decoding of the perturbation  $\delta$  doesn't provide any interesting information.

### A.3 ADDITIONAL VISUALISATIONS

We supplement the main text with visualisations of the UNI baseline, attributions and path features (properties, stability and robustness). We additionally include figures elucidating the colour, texture and frequency biases post-hoc imposed by path attribution methods. From Figure 10, we observe the stability of UNI path features: our attributions can be reliably and efficiently computed with Riemann approximation. In Figures 11, 12 and 13, we present visualisations on ImageNet-C, highlighting how static choices of baselines may bias the path-attribution procedure, leading to null or noisy explanations. UNI does not impose additional post-hoc assumptions that are alien to the model's decision function. Furthermore, we present qualitative comparisons of attribution results of pre-trained models on the ImageNet-1K test set, in Figures 14, 15, 16, 17, 18 and 19. UNI attributions are visibly better localised and more semantically meaningful. Finally, we visualise the consistent, geodesic paths of monotonically increasing output confidence, discovered by UNI. As seen from Figures 20, 21, 22, 23, 24 and 25, while other path attribution methods might encounter extrema and turning points along the interpolation path from baseline to input, UNI's path features are monotonic and preserve the crucial completeness property on which the path attribution framework depends.

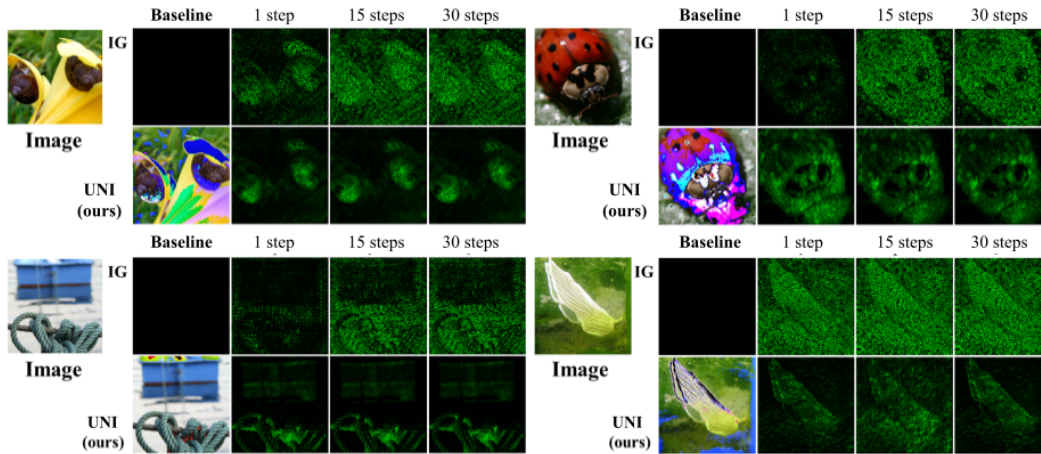
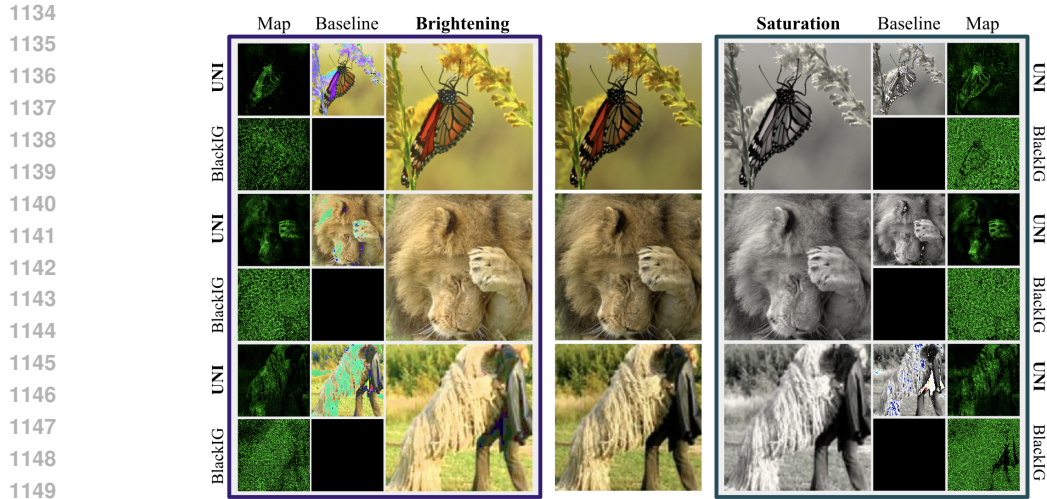


Figure 10: Comparison of attribution maps computed by Integrated Gradients and UNI, for a pre-trained ResNet-18 on the ImageNet-1K test set. UNI occludes and unlearns predictive input features; reliably localises predictive image regions; can be efficiently computed with only 1 Riemann step.



1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

Figure 11: *Colour bias*: When an image’s brightness or saturation is altered, IG with a black baseline fails to identify dark features, such as the wings of the butterfly (R1) or black jacket (R3).

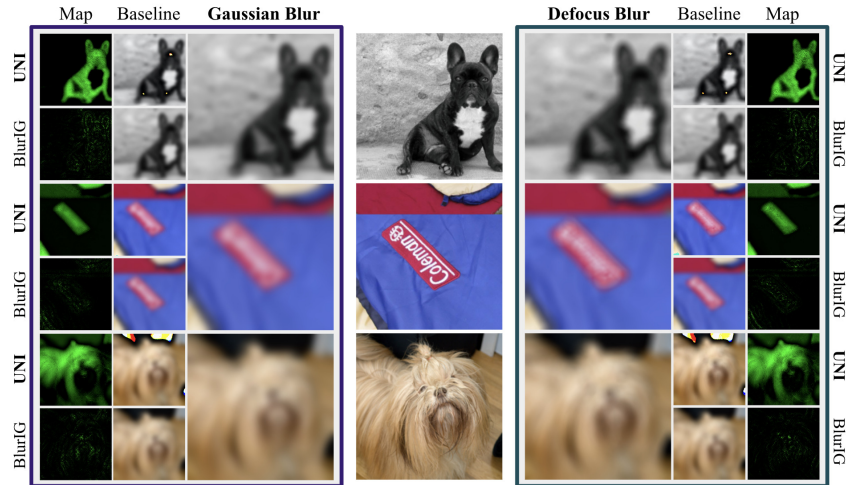


Figure 12: *Texture bias*: Using a blurred baseline for IG leads to a smoothness assumption in image texture, which leads to missingness in attribution when the input is also gaussian or defocus blurred.

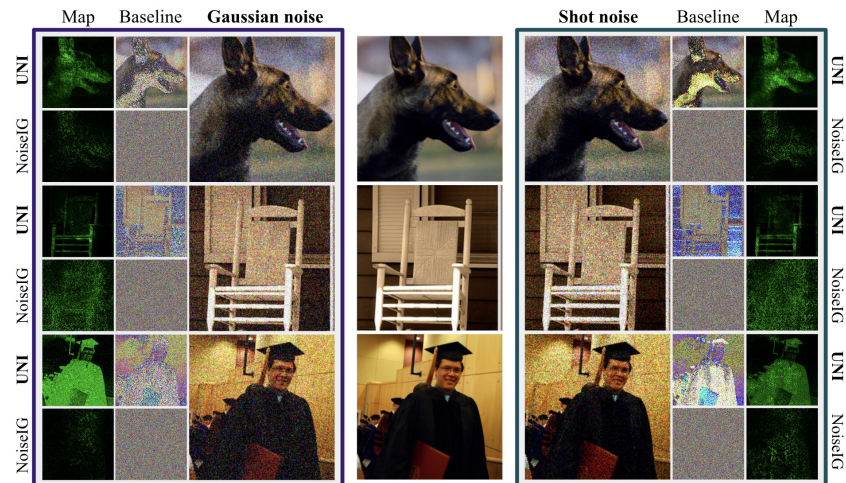


Figure 13: *Frequency bias*: A gaussian noised baseline for IG renders it vulnerable to high-frequency corruptions. Adding gaussian or shot noise to the image yields unmeaningful, noisy attributions.

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

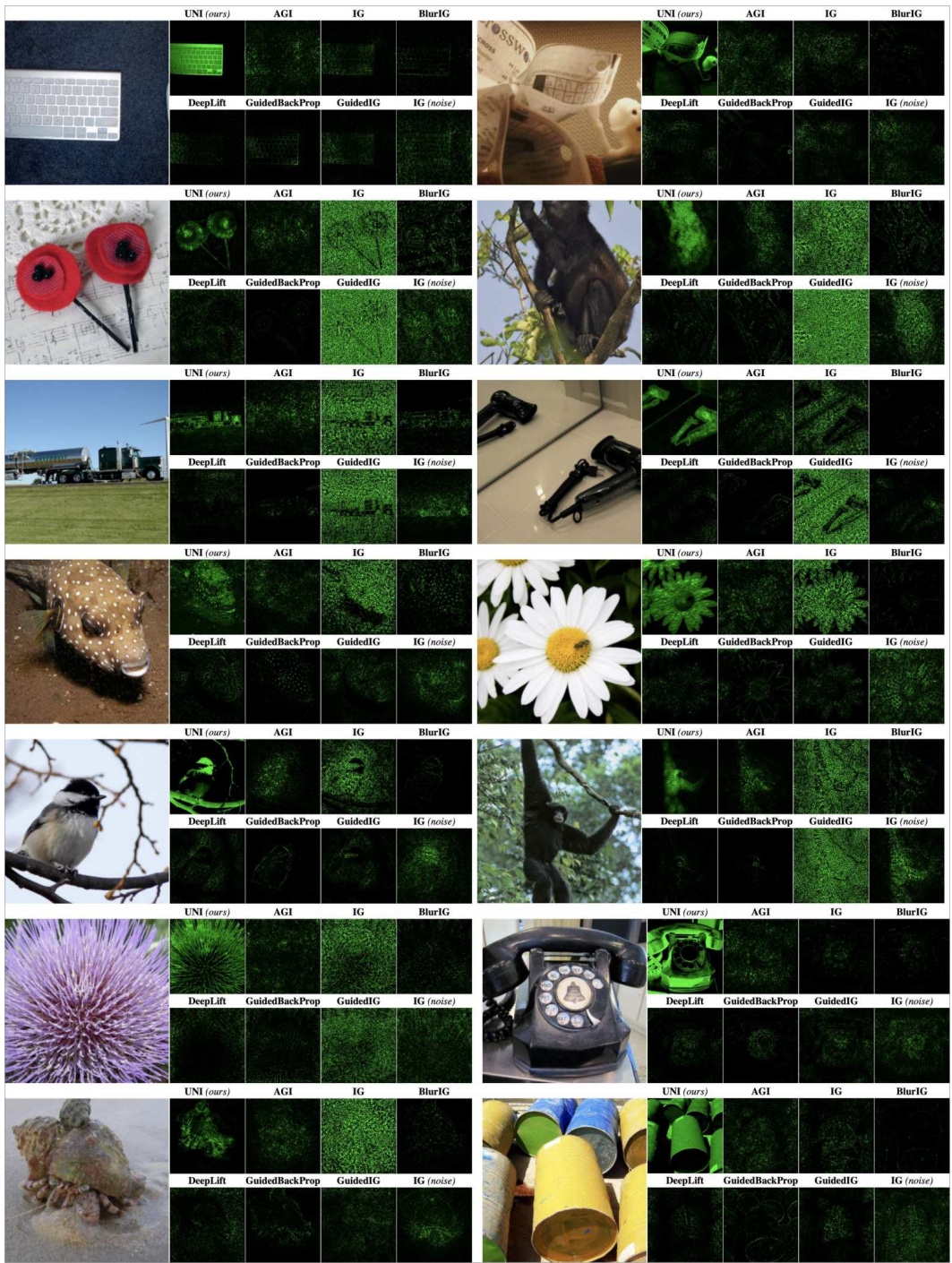


Figure 14: Comparing attributions (*ResNet-18*): UNI attributions demonstrate higher saliency, fidelity and faithfulness relative to conventional baselines on the ImageNet test set.

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

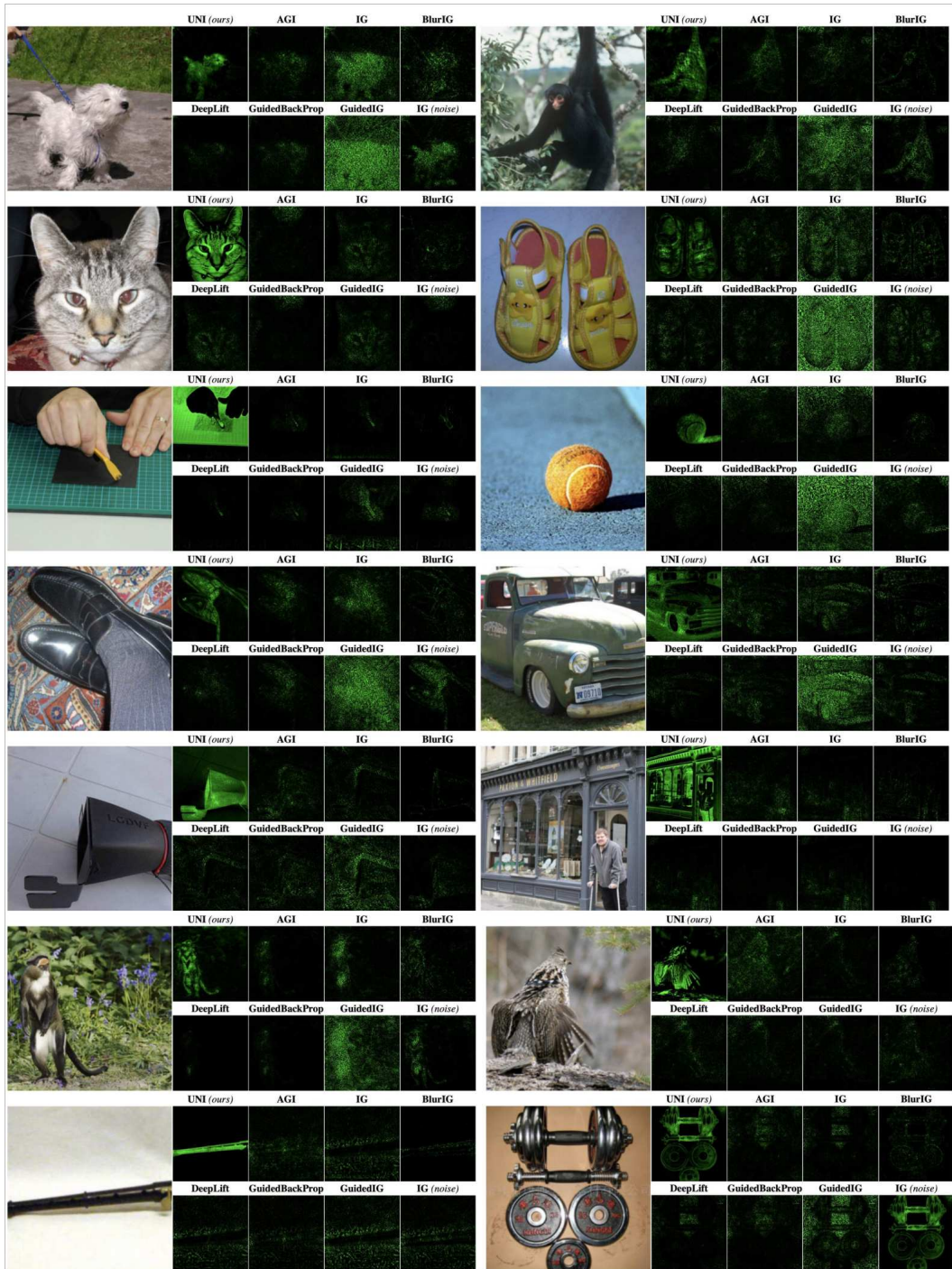


Figure 15: Comparing attributions (EfficientNet-v2-small): UNI attributions demonstrate higher saliency, fidelity and faithfulness relative to conventional baselines on the ImageNet test set.



1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

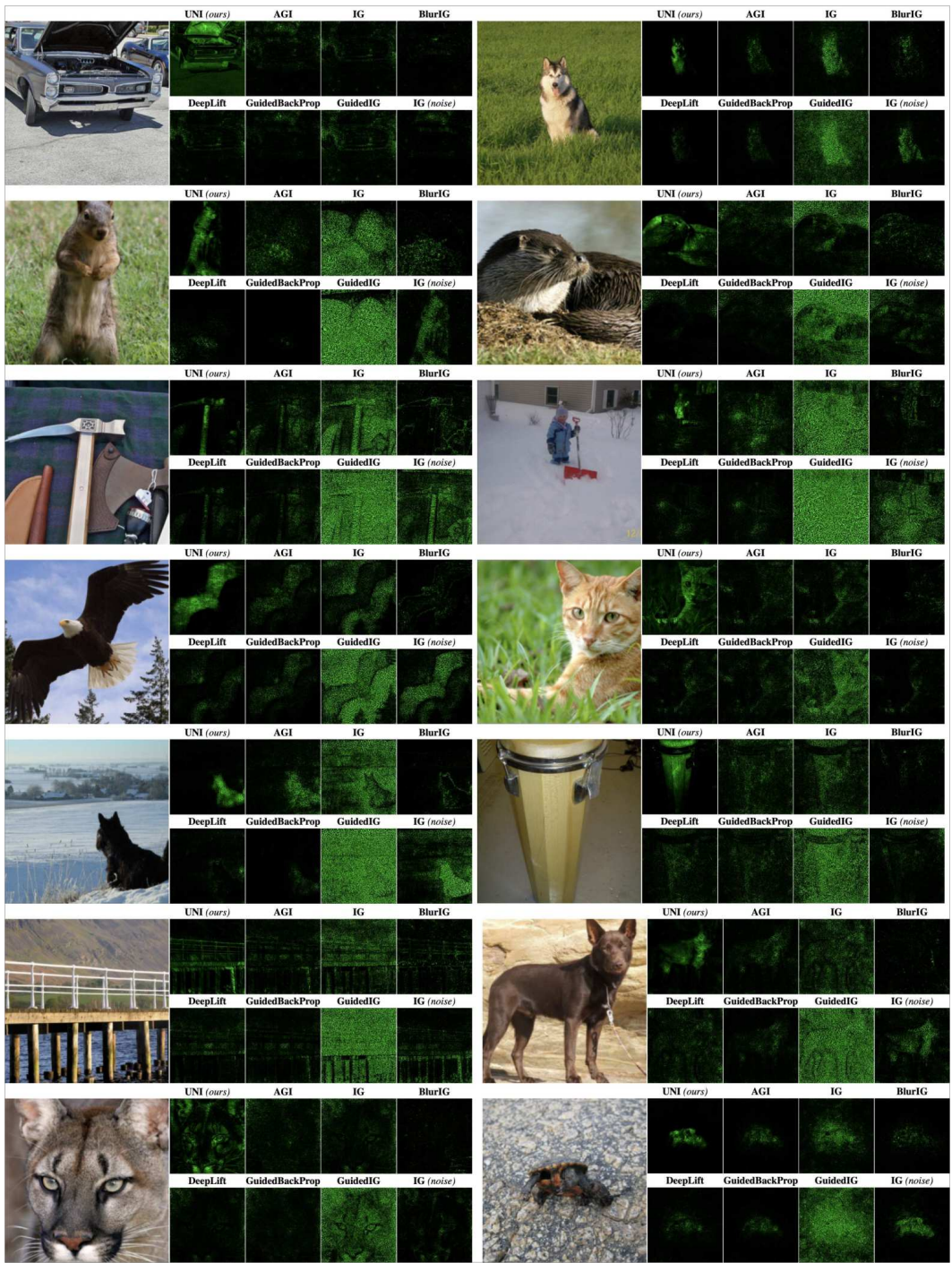


Figure 16: *Comparing attributions (ConvNext-Tiny):* UNI attributions demonstrate higher saliency, fidelity and faithfulness relative to conventional baselines on the ImageNet test set.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

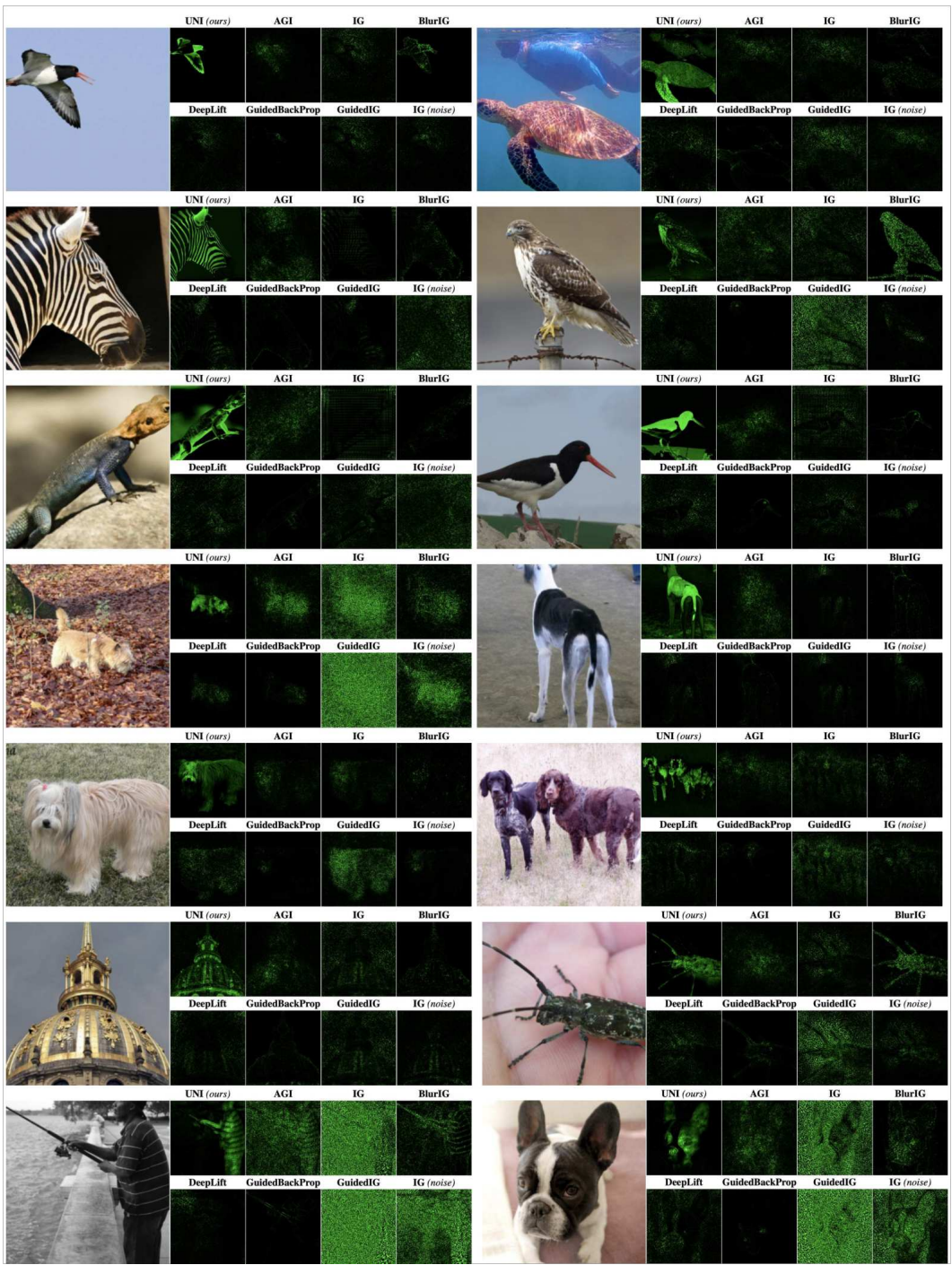


Figure 17: *Comparing attributions (VGG-16-bn):* UNI attributions demonstrate higher saliency, fidelity and faithfulness relative to conventional baselines on the ImageNet test set.

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

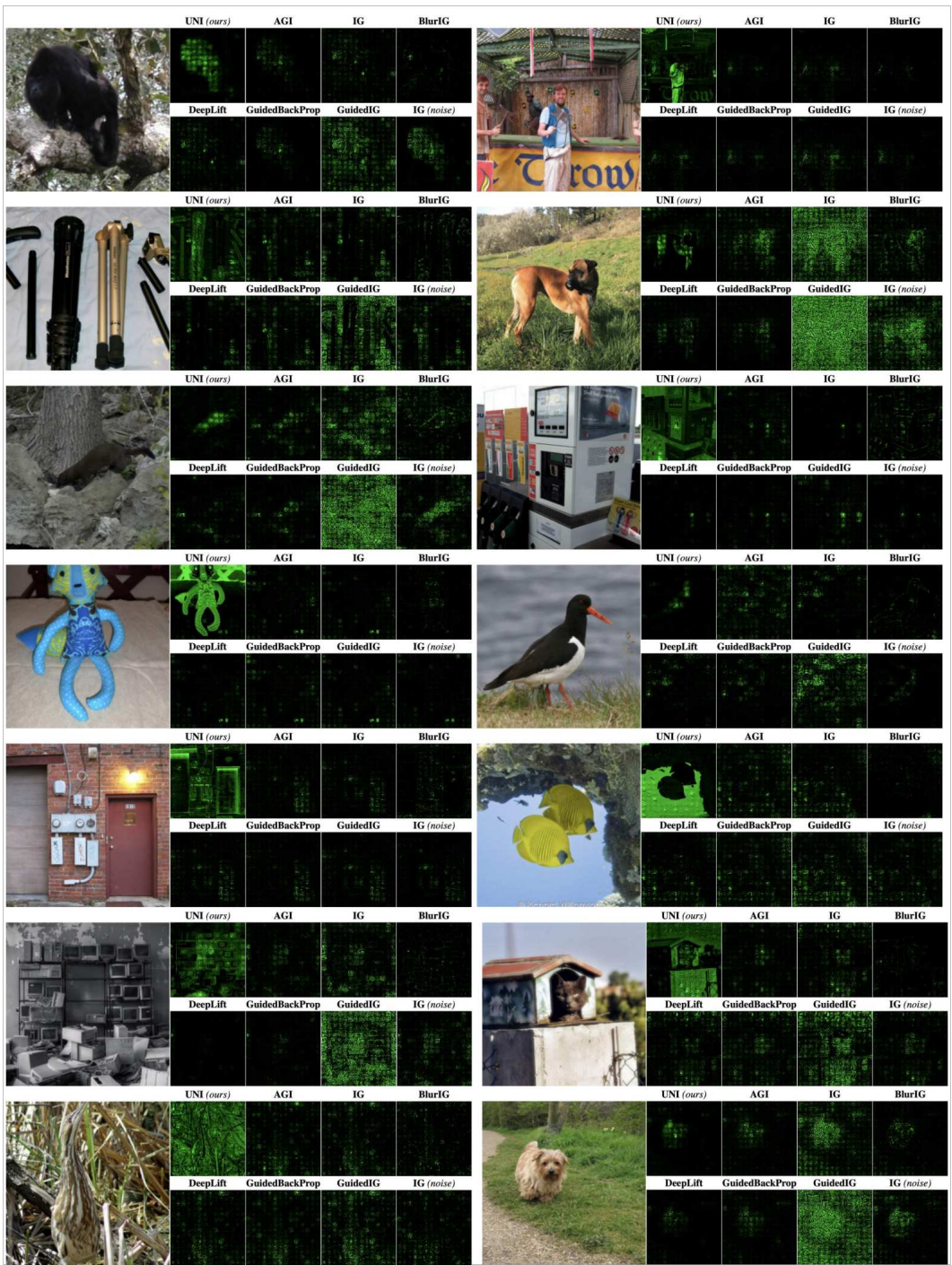


Figure 18: Comparing attributions (ViT-B\_16): UNI attributions demonstrate higher saliency, fidelity and faithfulness relative to conventional baselines on the ImageNet test set.

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

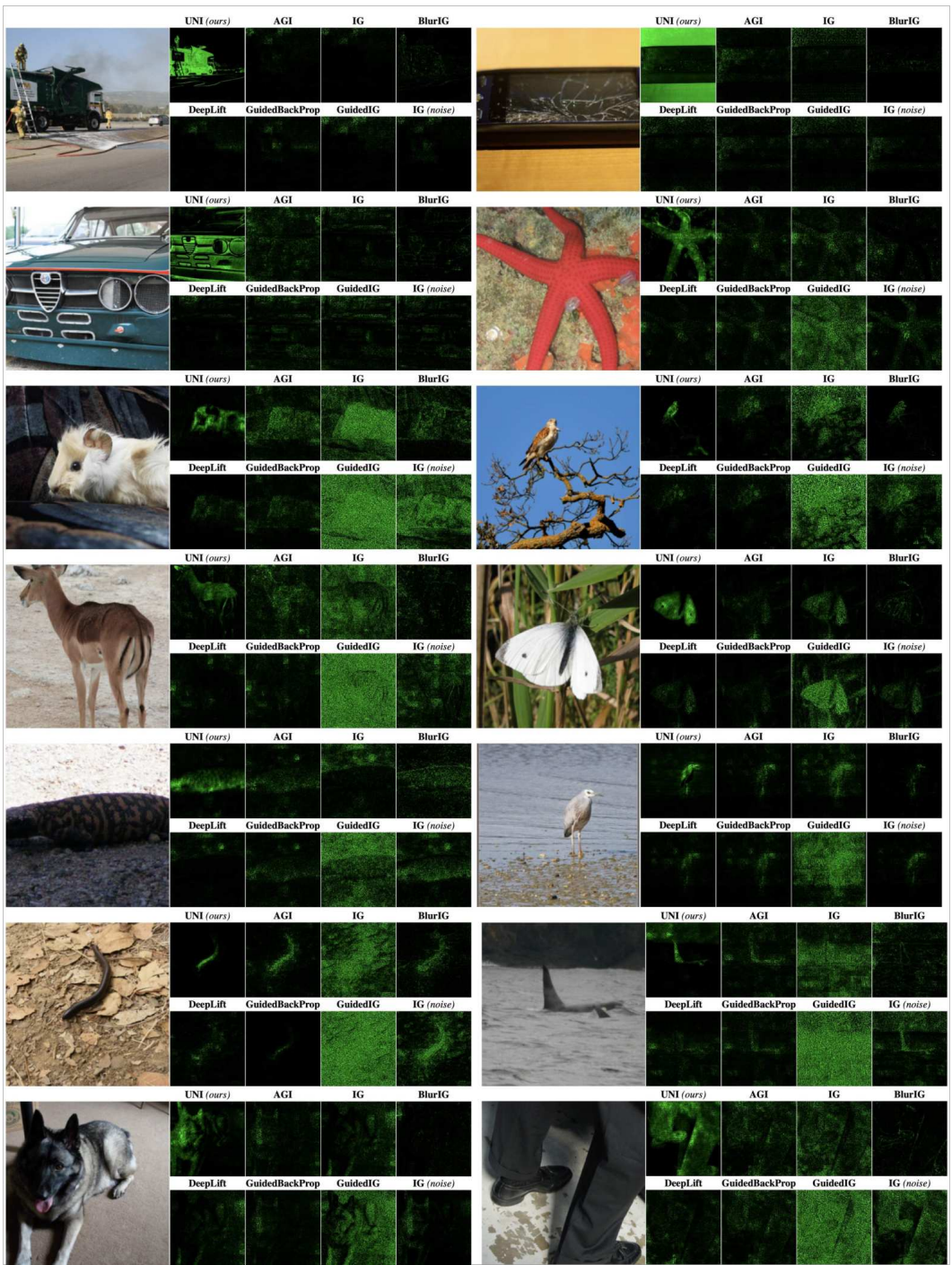


Figure 19: Comparing attributions (Swin-Transformer-Tiny): UNI attributions demonstrate higher saliency, fidelity and faithfulness relative to conventional baselines on the ImageNet test set.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

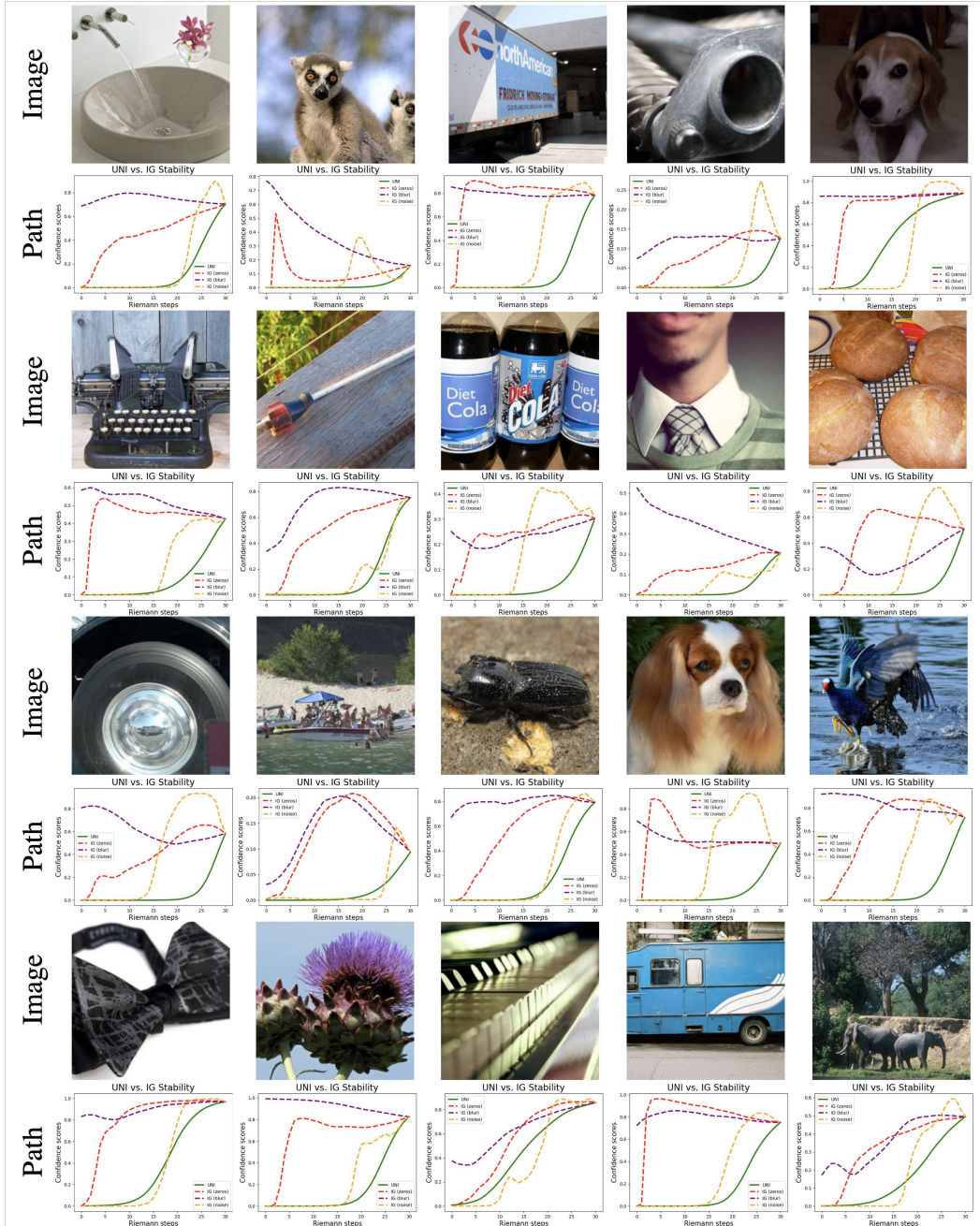


Figure 20: Comparing paths (ResNet-18): UNI discovers geodesic paths of monotonically increasing output confidence, preserving the completeness property required for robust attributions.

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

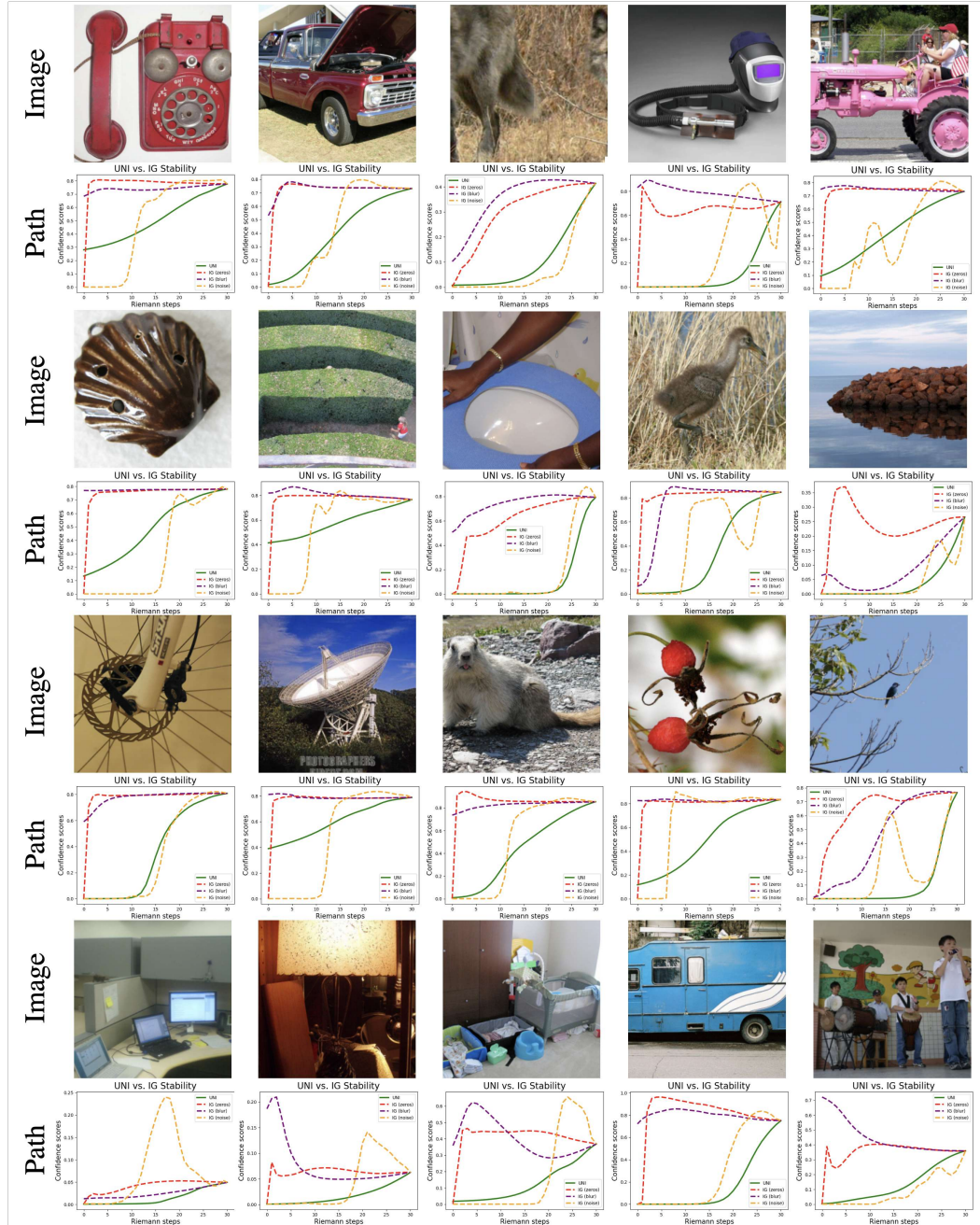


Figure 21: Comparing paths (EfficientNet-v2-small): UNI discovers geodesic paths of monotonically increasing output confidence, preserving the completeness property required for robust attributions.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

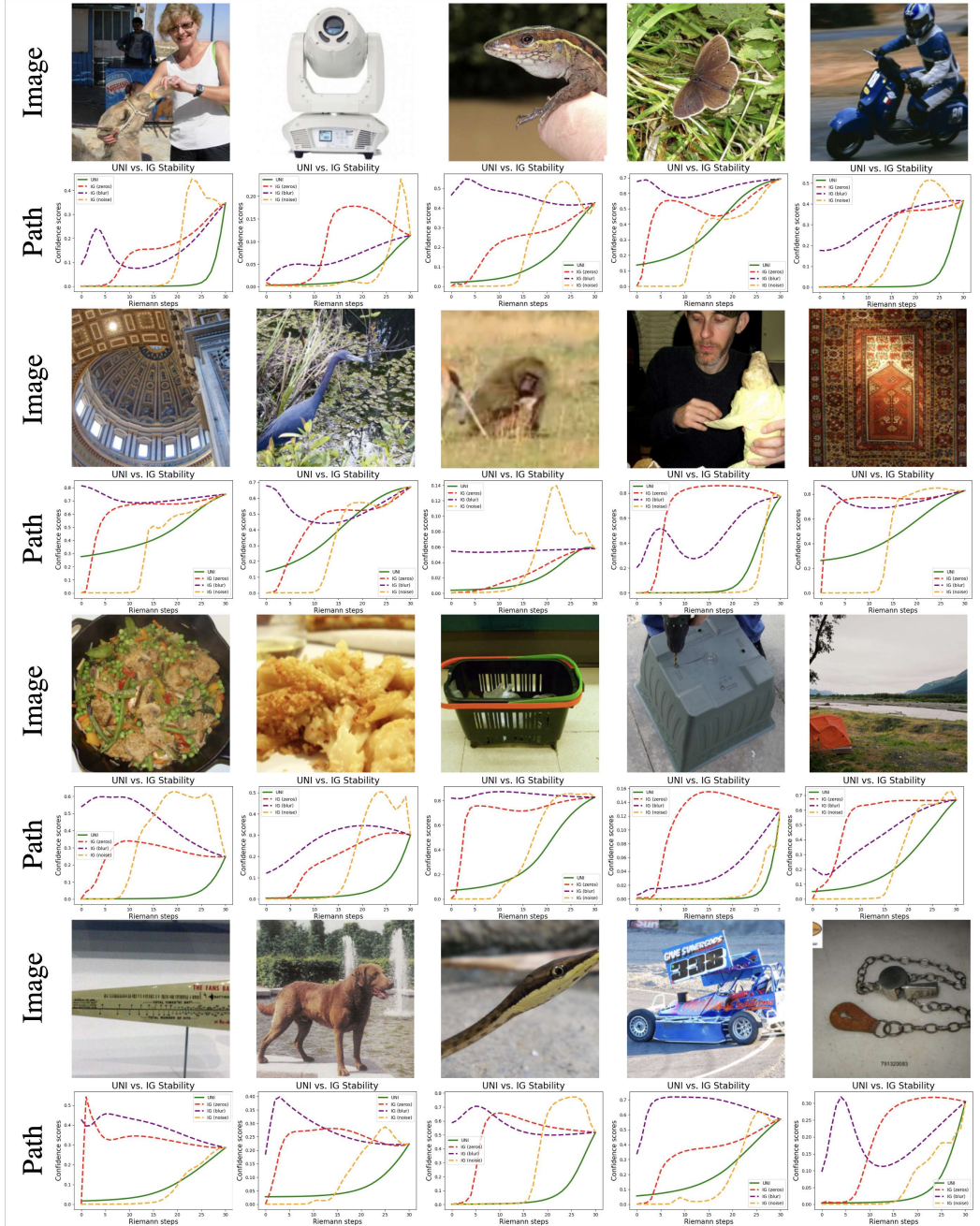


Figure 22: Comparing paths (ConvNeXt-Tiny): UNI discovers geodesic paths of monotonically increasing output confidence, preserving the completeness property required for robust attributions.

1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727

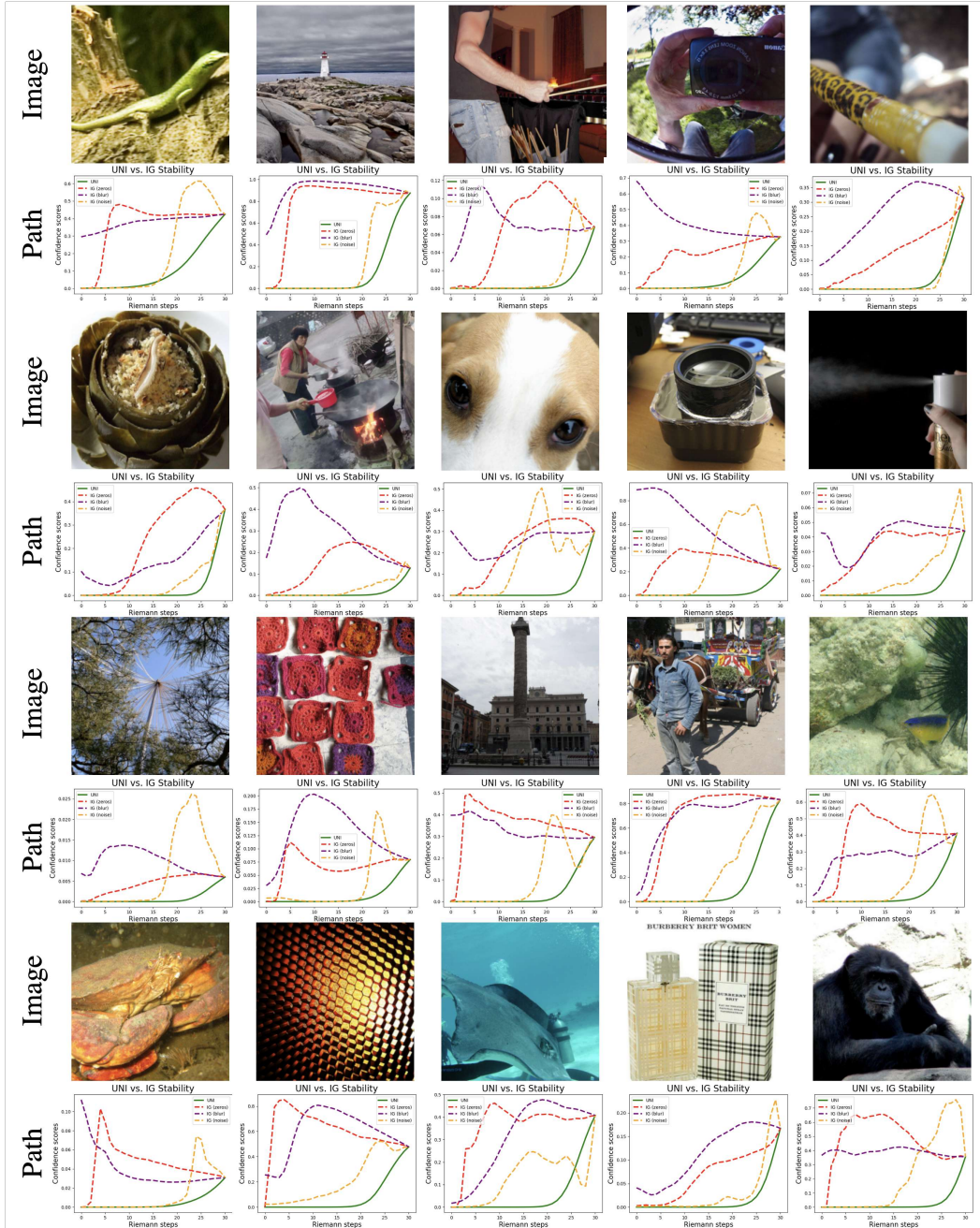


Figure 23: Comparing paths (VGG-16-bn): UNI discovers geodesic paths of monotonically increasing output confidence, preserving the completeness property required for robust attributions.



1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

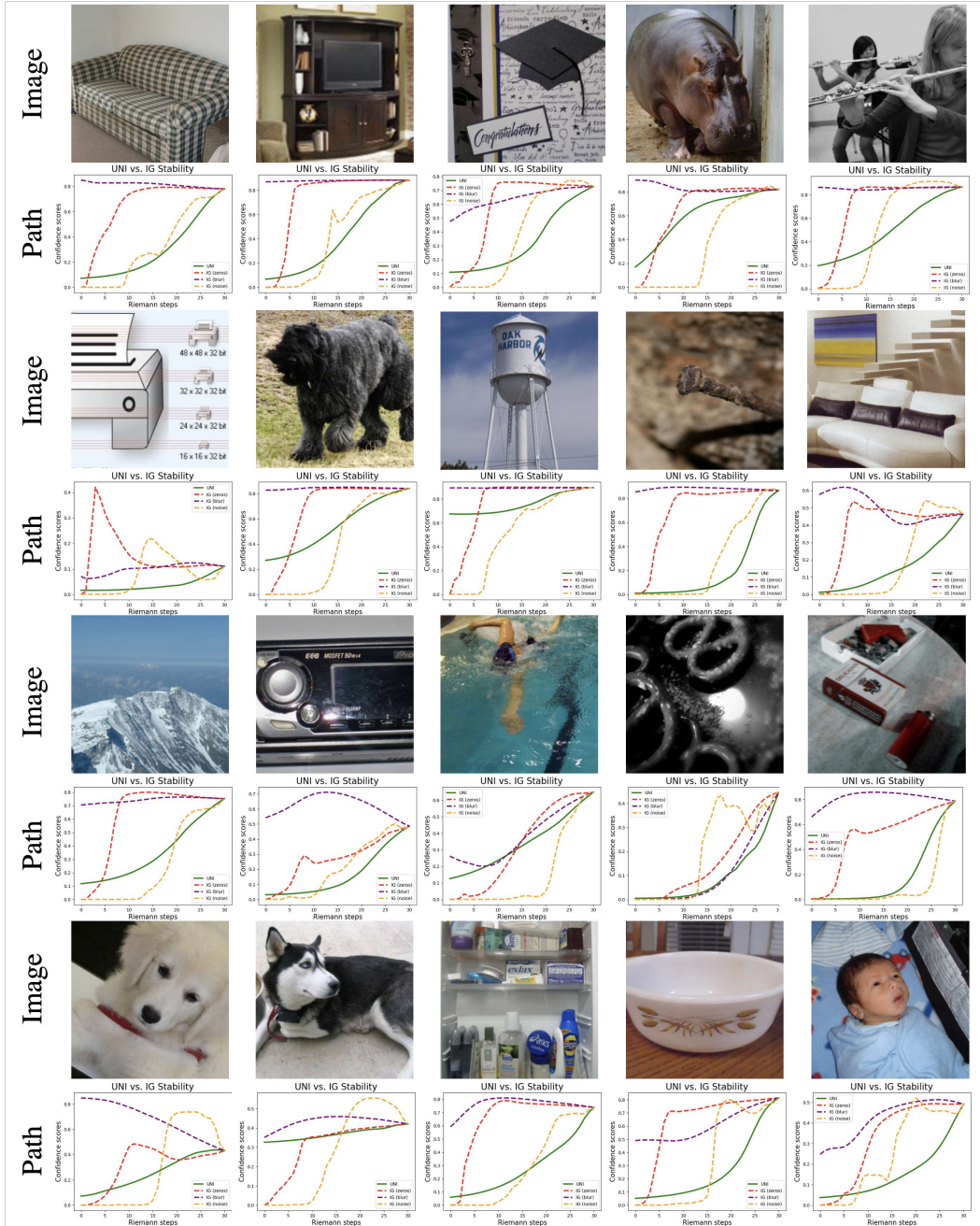


Figure 24: Comparing paths (ViT-B\_16): UNI discovers geodesic paths of monotonically increasing output confidence, preserving the completeness property required for robust attributions.

1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835

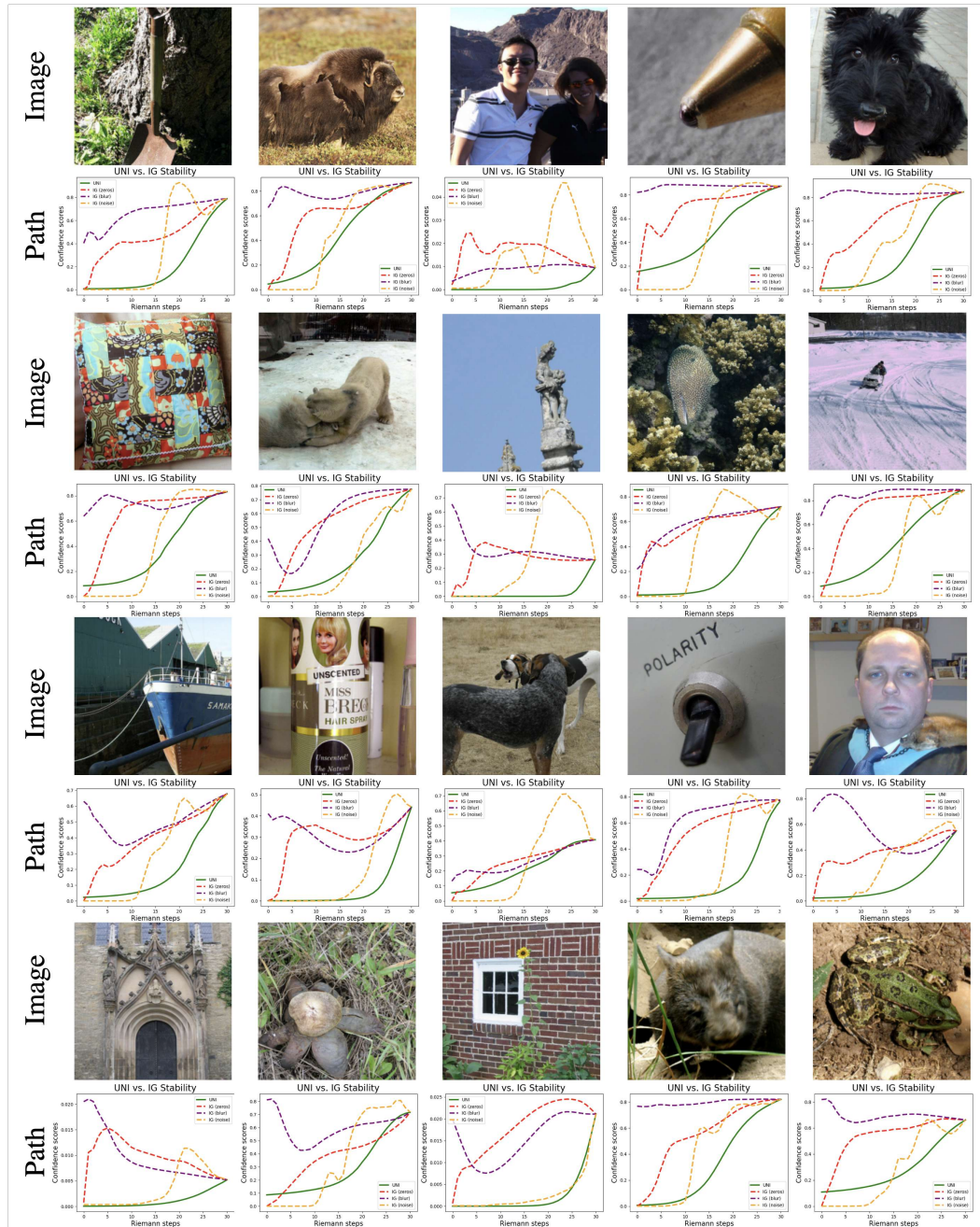


Figure 25: Comparing paths (Swin-Transformer-Tiny): UNI discovers geodesic paths of monotonically increasing output confidence, preserving the completeness property required for robust attributions.