

Evaluating Interpretable Methods via Geometric Alignment of Functional Distortions

Anna Hedström^{1,6,8,†}
Philine Bommer^{1,8}
Thomas F Burns^{3,4,5}
Sebastian Lapuschkin⁶
Wojciech Samek^{1,2,6}
Marina M.-C. Höhne^{2,7,8,†}

anna.hedstroem@tu-berlin.de
philine.bommer@tu-berlin.de
t.f.burns@gmail.com
sebastian.lapuschkin@hhi.fraunhofer.de
wojciech.samek@hhi.fraunhofer.de
mhoehne@atb-potsdam.de

¹ Department of Electrical Engineering and Computer Science, Technical University of Berlin

² BIFOLD – Berlin Institute for the Foundations of Learning and Data

³ Institute for Computational and Experimental Research in Mathematics, Brown University

⁴ Scientific Artificial Intelligence Center, Cornell University, USA

⁵ Neural Coding and Brain Computing Unit, Okinawa Institute of Science and Technology

⁶ Department of Artificial Intelligence, Fraunhofer Heinrich Hertz Institute

⁷ Department of Computer Science, University of Potsdam


⁸ UMI Lab, Leibniz Institute of Agricultural Engineering and Bioeconomy e.V. (ATB)

[†] corresponding authors

Reviewed on OpenReview: <https://openreview.net/forum?id=5ceyt8qT4e>

Abstract

Interpretability researchers face a universal question: without access to ground truth labels, how can the faithfulness of an explanation to its model be determined? Despite immense efforts to develop new evaluation methods, current approaches remain in a pre-paradigmatic state: fragmented, difficult to calibrate, and lacking cohesive theoretical grounding. Observing the lack of a unifying theory, we propose a novel evaluative criterion entitled *Generalised Explanation Faithfulness (GEF)* which is centered on explanation-to-model alignment, and integrates existing perturbation-based evaluations to eliminate the need for singular, task-specific evaluations. Complementing this unifying perspective, from a geometric point of view, we reveal a prevalent yet critical oversight in current evaluation practice: the failure to account for the learned geometry, and non-linear mapping present in the model, and explanation spaces. To solve this, we propose a general-purpose, threshold-free faithfulness evaluator *GEF* that incorporates principles from differential geometry, and facilitates evaluation agnostically across tasks, and interpretability approaches. Through extensive cross-domain benchmarks on natural language processing, vision, and tabular tasks, we provide first-of-its-kind insights into the comparative performance of various interpretable methods. This includes local linear approximators, global feature visualisation methods, large language models as post-hoc explainers, and sparse autoencoders. Our contributions are important to the interpretability and AI safety communities, offering a principled, unified approach for evaluation.

 <https://github.com/annahedstroem/GEF>

1 Introduction

Explaining the general behaviour, and predictions of machine learning (ML) models, particularly those functioning as black boxes, is critical, especially in domains such as healthcare, finance, and law. Driven by the urgency to comply with regulations like the EU AI Act, and GDPR, the interpretability (or eXplainable

AI (XAI) research community has produced a plethora of interpretable methods in recent years (Baehrens et al., 2010; Zeiler & Fergus, 2014a; Lundberg & Lee, 2017; Bykov et al., 2022; Fel et al., 2024; Lieberum et al., 2024). Simultaneously, the rise of large-scale, multi-tasking large language models (LLMs) (or “foundation models”) (OpenAI, 2023; Mesnard et al., 2024) has spurred a significant shift in the interpretability landscape, with the mechanistic interpretability community producing a new generation of methods specifically designed to decompose, and reverse-engineer these increasingly black-box models (Elhage et al., 2022; Conmy et al., 2023; Bykov et al., 2023; Bills et al., 2023; Templeton et al., 2024). Despite this immense activity, consensus is lacking whether existing methods are of sufficient quality or trustworthy (Adebayo et al., 2018; Ghassemi et al., 2021; Bordt & von Luxburg, 2024; Bhattacharjee & von Luxburg, 2024). Since black-box models lack ground truth explanation labels (Bellido & Fiesler, 1993; Benitez et al., 1997), the universal question: “how faithful is the explanation to the model it seeks to explain?” remains difficult to answer. The prevalence of disagreements within the interpretability community about which methods¹ work, and under what conditions (Neely et al., 2021; Watson et al., 2022; Krishna et al., 2022; Koenen & Wright, 2024) signals that the challenge of evaluation is still unsolved.

To *approximate* explanation quality (Agarwal et al., 2022b; Hedström et al., 2023b), researchers commonly use perturbation-based evaluations, where *robustness* (Montavon et al., 2018; Alvarez-Melis & Jaakkola, 2018b; Yeh et al., 2019; Nguyen & Martinez, 2020; Dasgupta et al., 2022), *sensitivity* (Adebayo et al., 2018; Hedström et al., 2024), and *faithfulness* (Bach et al., 2015; Samek et al., 2017; Ancona et al., 2018; Rieger & Hansen, 2020; Dasgupta et al., 2022; Bhatt et al., 2020; Rong et al., 2022) are well-embraced criteria to examine the relationship between explanation, and model outputs under perturbation, albeit with different emphases. Here, robustness, and sensitivity criteria refer to making small or large perturbations (*e.g.*, adding noise to the input or randomising model parameters), and then measuring corresponding changes in the explanation output. Faithfulness criterion generally measures how much the model’s performance degrades when inputs, such as tokens or pixels, are cumulatively perturbed according to the explanation values. Significant changes in model behaviour are interpreted as indicators of explanation faithfulness.

Lack of Cohesive, Unified Theory. Despite repeated attempts to define, and measure faithfulness (Montavon et al., 2018; Jacovi & Goldberg, 2020; Bhatt et al., 2020; Turpin et al., 2023; Lanham et al., 2023; Agarwal et al., 2024), fragmented mathematical terminology makes it an ongoing, and unresolved matter (Bordt & von Luxburg, 2024). What exactly is explanation faithfulness, and how do robustness, and sensitivity evaluations differ from it? From a conceptual standpoint, although these evaluations share common steps—such as perturbing the inputs or the model parameters, measuring the effects, and interpreting the functional outcomes—the overwhelming number of evaluation methods under these distinct criteria (Lakkaraju et al., 2022), and the absence of a cohesive, unified theory makes it difficult to answer such seemingly straightforward questions. To better understand these evaluations’ shared attributes, assumptions, and outcomes, a mathematical discussion is required. In Section 2, we propose a *unifying* perspective that formalises robustness, sensitivity, and faithfulness evaluations, providing a principled *Generalised Explanation Faithfulness* (GEF) definition in Section 3 to substitute singular evaluations.

Ignoring the Impact of Geometry. Alongside the lack of a cohesive, unified theory, most perturbation-based evaluations (Section 2.1.2)—while well-intended, and intuitive—often rely on overly simplistic assumptions about the underlying geometry of both model, and explanation spaces. When perturbations are introduced, the functional outcomes of models, and explanations are frequently compared using direct distance measures or correlation coefficients (Alvarez-Melis & Jaakkola, 2018b; Yeh et al., 2019; Ancona et al., 2018; Bhatt et al., 2020; Nguyen & Martinez, 2020; Agarwal et al., 2022a). From a *geometric* perspective, this overlooks a simple yet critical fact: that a uniform perturbation such as input noise or parameter shifts can affect non-linear systems in highly non-uniform ways. Only in a linear system, the perturbation effects would be uniform. By neglecting the geometric differences (*e.g.*, differences in curvatures) between the model, and explanation spaces, current evaluations risk misjudging how faithful the explanation is *w.r.t.* its underlying model. For fair measurements across non-linear systems, perturbation effects must be measured in the context of the distinct geometric structures of the respective manifolds (Lee, 2012). In Section 4, we

¹Throughout this work, we use the terms “interpretable methods” and “explanation methods” interchangeably, without implying a difference in their scope or function.

examine these geometric factors, and introduce a solution that accounts for the intrinsic geometry of each space, thereby improving current evaluation practice.

To address the research gaps in unified theory (Section 3), and the neglected impact of geometry (Section 4), our work offers a fourfold contribution.

- (C1) In the absence of cohesive theory, we systematise common steps in numerous perturbation-based evaluation algorithms (Section 2), and provide a unifying evaluative criterion for robustness, sensitivity, and faithfulness evaluations (Section 3).
- (C2) To account for geometric discrepancies in many evaluation methods, we propose a solution based on differential geometry that ensures fair measurements across non-linear mappings (Section 4).
- (C3) Recognising the need for a general-purpose, threshold-free, task-agnostic faithfulness evaluators, we provide **GEF** and **Fast-GEF** methods, serving distinct compute budgets (Section 5).
- (C4) Observing the lack of cross-domain insights on the faithfulness across distinct explanation approaches such as local, global, LLM as an explainer, and sparse autoencoders (SAEs), we perform extensive experiments across vision, tabular, and natural language processing (NLP) tasks (Section 6).

Our contributions carry substantial importance to the interpretability (and related) communities. The reliability of individual explanation methods, and XAI as a field is already under hot debate, thus it is not only timely but relevant to provide clarity on the matter of explanation faithfulness. As we enter a new era of interpretability, it is of utmost importance to revisit, and revise existing evaluation approaches. We hope this work will clarify how best to approach and perform faithfulness evaluation, ultimately empowering researchers to confidently select and develop new interpretable methods.

2 Interpretability Evaluation: Where Are We Now?

In this section, we present the scope of this work. We begin by outlining preliminaries to estimate explanation quality, followed by a description of the general workflow of perturbation-based evaluation. Finally, we mathematically formalise robustness, sensitivity, and faithfulness evaluation, revealing critical assumptions that are essential for their validity. Complete notation tables are provided in Appendix A.9.

2.1 Preliminaries

Let $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ be a differentiable neural network (NN) that maps inputs $\mathbf{x} \in \mathbb{R}^D$ to predictions $\mathbf{y} \in \mathbb{R}^C$ of C classes. By functionally mapping $\mathbf{x} \in \mathcal{X}$ to $\mathbf{y} \in \mathcal{Y}$ with parameters θ such that $\mathbf{y} = f(\mathbf{x}; \theta)$, a trained model f_θ is obtained, which we refer to as f . Here, θ includes weights, and biases, and exists in parameter space $\Theta \in \mathbb{R}^W$ for a fixed architecture in function space $f_\theta \in \mathcal{F}$. The model f may represent NN architectures ranging from simple feedforward MLPs, CNNs to highly parameterised transformer-based models.

Local Explanations. To interpret a specific model prediction (*i.e.*, logit) $y := y_c$ of a class $c \in [1, 2, \dots, C]$, we may employ a *local* method. Let $\phi_L : \mathcal{F} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^V$ be a local explanation function that takes an input, and logit pair, and assigns importance scores to a subset (or all) of its input features such that

$$\mathbf{e} = \phi_L(f, \mathbf{x}, y; \lambda), \quad (1)$$

where $\mathbf{e} \in \mathbb{R}^V$ is the explanation output, parameterised by λ .

A broad variety of local explanation approaches fall within the scope of our work, *e.g.*, gradient-based (Simonyan & Zisserman, 2015; Smilkov et al., 2017; Sundararajan et al., 2017; Bykov et al., 2022; Krishna et al., 2023; Selvaraju et al., 2020), back-propagation-based (Bach et al., 2015; Shrikumar et al., 2017), model-agnostic (Zeiler & Fergus, 2014a; Lundberg & Lee, 2017), local surrogate (Ribeiro et al., 2016a), attention-based (Chefer et al., 2021; Covert et al., 2022), or prototypical explanation methods (Simonyan & Zisserman, 2015). More recent approaches (Krishna et al., 2023; Kroeger et al., 2023) that leverage separate LLMs as the explanation function ϕ to interpret local predictions in a post-hoc manner, are also within the scope of this work.

Global Explanations. To study the model f from a *global* point of view, an explanation is produced independent of a specific instance, *i.e.*, \mathbf{x} . Here, a global explanation method $\phi_G : \mathcal{F} \times \mathcal{Y} \rightarrow \mathbb{R}^V$ takes a trained model f , and generates an explanation $\mathbf{e} \in \mathbb{R}^V$ for specific neural activation associated with a target class c . Here, c is represented by logit y such that

$$\mathbf{e} = \phi_G(f, y; \kappa), \quad (2)$$

where ϕ_G is parameterised by κ . Here ϕ_G may be variants of activation-maximisation (or “feature visualisation”) which provide either natural, or synthetic data points of maximal activation (Berkes & Wiskott, 2006; Erhan et al., 2009; Olah et al., 2017; Nguyen, 2020; Fel et al., 2024) or concept-based explanations Bykov et al. (2023). Recently, trained SAEs Bricken et al. (2023); Lieberum et al. (2024); Huben et al. (2024) have emerged as an alternative formulation for ϕ_G , aiming to produce interpretable “monosemantic” feature encodings at the layer level, providing insight into a model’s intermediate representations.

For convenience, we let $\phi \in \mathcal{E}$ denote ϕ_L , and ϕ_G although they formally reside in different spaces. To avoid label leakage (Jethani et al., 2023), we use the predicted class (and not the true class) to generate the explanation \mathbf{e} .

2.1.1 Estimate Explanation Quality

Without ground truth explanation labels, the task of estimating the quality of an explanation ϕ is non-trivial. To approximate explanation quality, researchers rely on metric-based heuristics (or “metrics”). Following Hedström et al. (2023a), we define a general evaluation function $\Psi_\tau : \mathcal{E} \times \mathcal{X} \times \mathcal{F} \times \mathcal{Y} \rightarrow \mathbb{R}$

$$q = \Psi(\phi, \mathbf{x}, f, y; \tau) \quad (3)$$

which returns a quality estimate $q \in \mathbb{R}$, indicating the quality of a given explanation, parameterised by τ . When global explanations ϕ_G are evaluated, \mathbf{x} is omitted from Equation 3. Unless required, we omit hyperparameters $\tau, \lambda, \kappa, \zeta$ for notational convenience.

2.1.2 Related Works

Within approaches that evaluate explanation quality by approximation, we concentrate on those that examine the *functional relationship* between the explanation, and the model through means of perturbation, *i.e.*, assessing qualities such as robustness, sensitivity, and faithfulness. These are briefly introduced below, and mathematically formalised in Section 2.3.

Robustness. Robustness (also referred to as “continuity”, and “stability”) methods evaluate the explanation function’s resilience to infinitesimal input noise, and is a widely used evaluation technique (Yeh et al., 2019; Montavon et al., 2018; Alvarez-Melis & Jaakkola, 2018b; Nguyen & Martinez, 2020; Agarwal et al., 2022a; Dasgupta et al., 2022). Most commonly, robustness is evaluated by first perturbing an input sample, then generating the explanation for the perturbed input, and finally comparing this explanation to the original explanation. Higher similarity between the original, and perturbed explanation indicates higher quality. Existing robustness measures differ in how noise is applied to the input (*e.g.*, using a Gaussian (Alvarez-Melis & Jaakkola, 2018b; Yeh et al., 2019) or a uniform distribution (Agarwal et al., 2022a)), and how explanation similarity is measured (*e.g.*, Yeh et al. (2019) computes difference with *Monte-Carlo* sampling, and Alvarez-Melis & Jaakkola (2018b); Agarwal et al. (2022a) rely on variants of a *Lipschitz* constant).

Sensitivity. Sensitivity (or “randomisation”) methods (Adebayo et al., 2018; Hedström et al., 2024) act complementary to robustness, and assesses a critical, and perhaps indisputable evaluative quality: that the explanation function ϕ should be sensitive to a randomisation of model parameters. Existing sensitivity measures differ in how the change in the explanation outputs is measured (*e.g.*, Adebayo et al. (2018) relies on *Structural Similarity Index* (SSIM), and Hedström et al. (2024) uses discrete entropy calculations), and how perturbation is applied (*e.g.*, Adebayo et al. (2018) randomises model parameters layer-by-layer in a *top-down* fashion, and Hedström et al. (2024) uses *bottom-up* or *full* parameter randomisation). The sensitivity criterion asks that the explanation should change significantly when the model parameters are randomised, whether layer-by-layer (Adebayo et al., 2018) or entirely (Hedström et al., 2024).

Faithfulness. Faithfulness (or “fidelity”) methods (Bach et al., 2015; Samek et al., 2017; Montavon et al., 2018; Ancona et al., 2018; Rieger & Hansen, 2020; Dasgupta et al., 2022; Bhatt et al., 2020; Rong et al., 2022; Atanasova et al., 2023; Blücher et al., 2024; Chuang et al., 2024) evaluate explanations by gradually perturbing the input based on the importance of pixels or tokens indicated by the explanation, and observing the resulting degradation in model performance. Methods differ in how model responses are reported (with logits (Alvarez-Melis & Jaakkola, 2018a; Yeh et al., 2019; Bhatt et al., 2020) or softmax probabilities (Montavon et al., 2018; Ancona et al., 2018; Rieger & Hansen, 2020; Nguyen & Martinez, 2020; Dasgupta et al., 2022; Rong et al., 2022)), how perturbations are ordered (ascending (Arya et al., 2019; Nguyen & Martinez, 2020) or descending (Bach et al., 2015; Samek et al., 2017; Rong et al., 2022)), and in the general approach to perturbation (whether using single-pixel changes (Bach et al., 2015), patch-based masking with a constant value (Samek et al., 2017), or linear interpolation (Rong et al., 2022)). Faithfulness methods typically aggregate model responses into a single quality estimate, such as AUC (Bach et al., 2015; Samek et al., 2017; Rong et al., 2022). For faithfulness to be considered fulfilled, the model’s performance should rapidly decrease as perturbations are applied—the steeper the degradation, the higher the explanation quality.

Beyond approximation techniques, interpretability researchers have explored alternative ways to evaluate explanation quality, such as using human judgment (Zeiler & Fergus, 2014b; Ribeiro et al., 2016b), and restricting tasks to synthetic or toy environments (Guidotti, 2021; Carmichael & Scheirer, 2023). While such approaches complement evaluation methods that approximate explanation quality, they lack scalability, and generalisability to real-world scenarios, and are not covered in this work.

2.2 Perturbation-based Evaluation

A key observation is that robustness, sensitivity, and faithfulness evaluations generally rely on three common steps. First, a perturbation is applied to either the input (*e.g.*, by adding infinitesimal noise) or the model parameters (*e.g.*, by randomisation). Second, the effect of the perturbation is measured on the output of either the explanation function ϕ or the model f . Third, an interpretation is made to assess whether this change in functional outputs is acceptable given a criterion, such as requiring the distance in explanation outputs to be small when the perturbation is small. We refer to Figure 1 for an illustration.

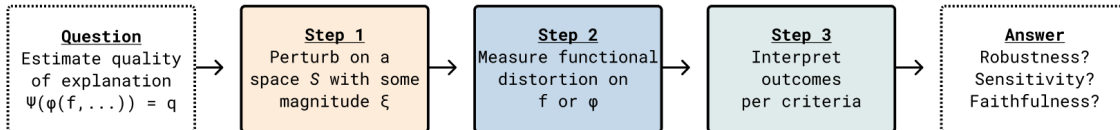


Figure 1: An overview of the “perturb, measure, and interpret” evaluation methodology (Section 2.3).

To facilitate mathematical unification (Section 4), and further insights (Section 5), we next formalise the three steps of perturbation-based evaluation. Therefore, some general notation for perturbation (Eqs 4-5), and measurement (Definition 1) is introduced. By systematising evaluation, we can advance our conceptual understanding, especially in clarifying how existing methods differ, and what attributes are shared.

2.2.1 Step 1. Perturbation

First, a perturbation is initiated. This is typically done either on the model parameter space in large magnitudes, *e.g.*, by randomising weights, or on the input in small magnitudes, *e.g.*, by adding Gaussian noise. Alternatively, perturbations can be applied cumulatively, such as by masking pixels or regions of pixels, or by replacing tokens in textual inputs. To accommodate diverse evaluation methods across different data modalities, we follow Hedström et al. (2023a), and define a general perturbation function that can be applied on any real-valued space $\mathcal{S} \subseteq \{\mathcal{X}, \Theta, \mathcal{Y}\}$. Let $\mathcal{P}_{\mathcal{S}} : \mathcal{S} \rightarrow \mathcal{S}$ be a perturbation function of $\mathbf{s} \in \mathcal{S}$ with parameters $\omega \in \mathbb{R}$ such that

$$\mathcal{P}_{\mathcal{S}}(\mathbf{s}; \omega) = \hat{\mathbf{s}}, \tag{4}$$

where $\forall \hat{\mathbf{s}}, \mathbf{s} \in \mathcal{S}$, and $\hat{\mathbf{s}} \neq \mathbf{s}$. For brevity, we may omit ω such that $\mathcal{P}_{\mathcal{S}}(\mathbf{s}) := \mathcal{P}_{\mathcal{S}}(\mathbf{s}; \omega)$. With Equation 4, we may, *e.g.*, generate a perturbed instance $\hat{\mathbf{s}}$ with input perturbation, *i.e.*, $\hat{\mathbf{x}} = \mathcal{P}_{\mathcal{X}}(\mathbf{x})$ or model parameter

randomisation, *i.e.*, $\hat{\theta} = \mathcal{P}_\Theta(\theta)$. Since robustness, sensitivity, and faithfulness evaluations require distinct perturbation magnitudes, we let ξ denote the difference between \mathbf{s} , and $\hat{\mathbf{s}}$ as follows

$$\delta(\mathbf{s}, \hat{\mathbf{s}}) = \xi, \quad (5)$$

where $\delta : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ is a general discrepancy function, *e.g.*, an ℓ_p -norm, cosine distance or Pearson correlation.

2.2.2 Step 2. Measurement

As a second step, the impact of the perturbation is measured on relevant functions. Common approaches include measuring the distance between explanation outputs or recording the change in model responses under random or cumulative masking guided by the explanation output. We define a general approach to measure the perturbation impact on a separate function (*e.g.*, the impact of input perturbation on the model function) below.

Definition 1 (Functional Distortion) *Let $\mathbf{s}, \hat{\mathbf{s}} \in \mathcal{S}$ denote instances in space $\mathcal{S} \subseteq \{\mathcal{X}, \Theta, \mathcal{Y}\}$, before, and after perturbation, respectively. Let $k : \mathcal{S} \rightarrow \mathcal{H}$ denote a separate function that maps $\mathbf{s}, \hat{\mathbf{s}}$ to a distinct space $\mathcal{H} \subseteq \{\mathcal{F}, \mathcal{E}\}$ from \mathcal{S} . Then, perturbation impact in function k is measured by functional distortion $\mathbf{D}_k : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ as follows*

$$\mathbf{D}_k(\mathbf{s}, \hat{\mathbf{s}}) = \delta(k(\mathbf{s}), k(\hat{\mathbf{s}})), \quad (6)$$

where $k(\mathbf{s}) = h$ with $h \in \mathcal{H}$, and $\delta : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$.

Model, and Explanation Distortion. With Definition 1, we can flexibly apply perturbation in one space, and then evaluate the effect in a different space². For example, assume we have applied perturbation on the input space, *i.e.*, $\hat{\mathbf{x}} = \mathcal{P}_\mathcal{X}(\mathbf{x})$ (Equation 4), and therefore have two instances \mathbf{x} , and $\hat{\mathbf{x}}$. Then, to measure the perturbation impact on the model function f , we follow Definition 1, and set $k = f$ where $h = \mathbf{y}$. Evaluating $\mathbf{D}_f(\mathbf{x}, \hat{\mathbf{x}})$ from Equation 6 effectively means that we compare model evaluations on perturbed, and non-perturbed inputs, *i.e.*, $\delta(y, \hat{y})$ with $\hat{y} = f_c(\hat{\mathbf{x}}; \theta)$ for the same class c . Alternatively, to measure perturbation impacts on the explanation function ϕ , we set $k = \phi$ where $h = \mathbf{e}$. Evaluating $\mathbf{D}_\phi(\mathbf{x}, \hat{\mathbf{x}})$, practically means that we compute $\delta(\mathbf{e}, \hat{\mathbf{e}})$ where $\hat{\mathbf{e}} = \phi(\hat{\mathbf{x}}, \dots)$ is the explanation *w.r.t.* perturbed input $\hat{\mathbf{x}}$. For comparability, $\hat{\mathbf{e}}$ is generated *w.r.t.* the same class c as its non-perturbed counterpart \mathbf{e} . Similarly, to compute functional distortion after parameter perturbation, *i.e.*, $\hat{\theta} = \mathcal{P}_\Theta(\theta)$, we compute $\mathbf{D}_f(\theta, \hat{\theta})$, and $\mathbf{D}_\phi(\theta, \hat{\theta})$ using logit $\hat{y} = f(\mathbf{x}; \hat{\theta})$, and explanation $\hat{\mathbf{e}} = \phi(f_{\hat{\theta}}, \dots)$, respectively. To generalise the notation across different perturbation types, we let \mathbf{D}_f , and \mathbf{D}_ϕ denote the model, and explanation distortion quantities, respectively.

2.2.3 Step 3. Interpretation

In the final step of the evaluation workflow (Figure 1), the distortion quantities are examined separately according to their evaluative criteria. For example, if robustness is evaluated, generally low values for \mathbf{D}_ϕ are expected, assuming perturbation magnitude ξ is small. Conversely, if sensitivity is evaluated, high values for \mathbf{D}_ϕ are expected, assuming perturbation magnitude ξ is large. If faithfulness is evaluated, model distortion \mathbf{D}_f is anticipated to increase as perturbation is cumulatively applied according to the explanation function output. Notably, a key limitation of this step is the need for thresholds to be set by researchers in order to distinguish between low-, and high-quality evaluation outcomes, which has shown could be adversarially manipulated (Wickstrøm et al., 2024).

2.3 Formalising Robustness, Sensitivity, and Faithfulness

Equipped with a general perturbation function $\mathcal{P}_\mathcal{S}$ (Equation 4), and its magnitude ξ (Equation 5) as well as a measure to compute functional distortion of the explanation, and model functions (Definition 1), we can combine a wide variety of existing evaluation techniques into general formalisations of robustness, sensitivity,

²While both the perturbation magnitude ξ (Equation 5), and the distortion \mathbf{D}_k (Equation 6) use the discrepancy function $\delta(\cdot, \cdot)$, their outputs differ. Notably, ξ expresses the discrepancy between the original, and perturbed instance, and \mathbf{D}_k measures the discrepancy in a distinct space from the perturbation space.

and faithfulness evaluation methods (*cf.* Section 2.1.1). Based on these three main criteria (Definitions 2-4), we show that the validity of each explanation criterion critically depends on fulfilling a separate, implicit model assumption (Assumptions 1-3).

We proceed by presenting a definition of explanation robustness, absorbing the spirit of numerous existing robustness methods³ (Yeh et al., 2019; Montavon et al., 2018; Alvarez-Melis & Jaakkola, 2018b; Nguyen & Martinez, 2020; Agarwal et al., 2022a).

Definition 2 (Explanation Robustness) Let $\hat{\mathbf{x}} = \mathcal{P}_{\mathcal{X}}(\mathbf{x})$ be a perturbed input, and Ψ^{RO} be a quality estimator to yield robustness estimates $q^{RO} \in \mathbb{R}$ such that $q^{RO} = \mathbf{D}_{\phi}(\mathbf{x}, \hat{\mathbf{x}})$. Given thresholds $\alpha, \varepsilon_{\mathbf{D}_{\phi}}^{RO} \in \mathbb{R}^+$, an explanation function ϕ is robust if the perturbation magnitude $\xi^{RO} \leq \alpha$:

$$q^{RO} \leq \varepsilon_{\mathbf{D}_{\phi}}^{RO}. \quad (7)$$

For an explanation function ϕ to be considered robust, the estimator Ψ^{RO} should yield low values, *i.e.*, $q^{RO} \leq \varepsilon_{\mathbf{D}_{\phi}}^{RO}$, reflecting minor differences between the original explanation \mathbf{e} , and the perturbed explanation $\hat{\mathbf{e}}$. Since the stability expectations of the explanation function ϕ are dictated by the robustness of f (Yeh et al., 2019; Chalasani et al., 2020; Agarwal et al., 2022a; Tan & Tian, 2023), it would be false to expect ϕ to exhibit robustness if its underlying model is not robust. Consequently, the validity of Equation 7 depends on the fulfillment of model robustness (Assumption 1).

Assumption 1 (Model Robustness) Given an input perturbation $\mathcal{P}_{\mathcal{X}}$ of magnitude ξ^{RO} , and thresholds $\alpha, \varepsilon_{\mathbf{D}_f}^{RO} \in \mathbb{R}^+$, $\xi^{RO} \leq \alpha$, the model distortion (Equation 6) is bounded by $\mathbf{D}_f(\mathbf{x}, \hat{\mathbf{x}}) \leq \varepsilon_{\mathbf{D}_f}^{RO}$.

In line with works of Adebayo et al. (2018); Hedström et al. (2024), we define explanation sensitivity in the following.

Definition 3 (Explanation Sensitivity) Let $\hat{\theta} = \mathcal{P}_{\Theta}(\theta)$ create a model $f_{\hat{\theta}}$ with perturbed parameters, and Ψ^{SE} be a quality estimator that yields sensitivity estimates $q^{SE} \in \mathbb{R}$ such that $q^{SE} = \mathbf{D}_{\phi}(\theta, \hat{\theta})$. Given thresholds $\alpha, \varepsilon_{\mathbf{D}_{\phi}}^{SE} \in \mathbb{R}^+$, an explanation function ϕ is sensitive if the perturbation magnitude $\xi^{SE} > \alpha$:

$$q^{SE} > \varepsilon_{\mathbf{D}_{\phi}}^{SE}. \quad (8)$$

For ϕ to be considered sensitive to randomness, the differences between explanations should be substantial, meaning Ψ^{SE} yields high estimates, *i.e.*, $q^{SE} > \varepsilon_{\mathbf{D}_{\phi}}^{SE}$, reflecting significant discrepancies between \mathbf{e} , and $\hat{\mathbf{e}}$. This expectation that q^{SE} should be large is based on the assumption that the model responded strongly to the perturbation. Similar to how explanation robustness depends on the stability of f , the emphasis on a large q^{SE} assumes a different model response. Therefore, the validity of the sensitivity evaluation (Equation 8) depends on model sensitivity (Assumption 2).

Assumption 2 (Model Sensitivity) Given a parameter perturbation \mathcal{P}_{Θ} of magnitude ξ^{SE} , and thresholds $\alpha, \varepsilon_{\mathbf{D}_f}^{SE} \in \mathbb{R}^+$, $\xi^{SE} > \alpha$, the model distortion (Equation 6) is bounded by $\mathbf{D}_f(\theta, \hat{\theta}) > \varepsilon_{\mathbf{D}_f}^{SE}$.

With various existing interpretations of explanation faithfulness (Section 2.1.2), we focus on common criteria to combine these interpretations into a single definition below.

Definition 4 (Explanation Faithfulness) Let $\hat{\mathbf{x}}^z = \mathcal{P}_{\mathcal{X}}(\mathbf{x}; z)$ denote the input after the z^{th} perturbation for $z \in [1, Z]$, where $\mathcal{P}_{\mathcal{X}}$ progressively masks the top- z features according to the indices given by $\text{argmax}(\mathbf{e})$, with perturbation magnitudes ξ_z satisfying $\xi_1 \leq \xi_2 \leq \dots \leq \xi_Z$. A quality estimator Ψ^{FA} yields a vector of faithfulness estimates $\mathbf{q}^{FA} \in \mathbb{R}^Z$ with entries $q_z^{FA} = f(\hat{\mathbf{x}}^z, \theta)$, where the overall faithfulness score $q^{FA} \in \mathbb{R}$ is obtained by aggregating these estimates via a function $\nu: \mathbb{R}^Z \rightarrow \mathbb{R}$:

$$q^{FA} = \nu(\hat{\mathbf{q}}^{FA}). \quad (9)$$

³Some algorithmic details are omitted in the definition. For completeness, mathematical definitions are provided for each evaluation method in Appendix A.4.5.

When ν is defined using AUC, a faithful explanation is expected to produce low aggregated scores q^{FA} (Equation 9). The conventional expectation in faithfulness evaluation (Bach et al., 2015; Samek et al., 2017; Rong et al., 2022) is that significant distortions should occur early, as the “more important features” are removed first. To ensure that the faithfulness score is solely driven by the quality of ϕ , and not by other factors, such as out-of-distribution samples (OOD) (Hase et al., 2021; Hesse et al., 2024), non-linear feature effects or artifacts introduced by cumulative perturbations (Hooker et al., 2019; Brunke et al., 2020; Hase et al., 2021; Rong et al., 2022; Brocki & Chung, 2022), the model distortion to these perturbations should be monotonically non-decreasing, *i.e.*, satisfy model faithfulness (Assumption 3).

Assumption 3 (Model Faithfulness) *Given Z cumulative perturbations $\mathcal{P}_{\mathcal{X}}$ of magnitudes ξ_z with $\xi_1 \leq \xi_2 \leq \dots \leq \xi_Z$ the corresponding model distortions (Equation 6) are: $\mathbf{D}_f^1 \leq \mathbf{D}_f^2 \leq \dots \leq \mathbf{D}_f^Z$ with $\mathbf{D}_f^z = \mathbf{D}_f(\mathbf{x}, \hat{\mathbf{x}}^z)$.*

2.4 Model Assumptions in Practice

Evaluations under Definitions 2-4 typically assume that model distortions are proportional to perturbation magnitudes, *i.e.*, that larger perturbations lead to greater distortions, and smaller perturbations result in lesser distortions. This naturally raises the question: with commonly used perturbation techniques for evaluating robustness (*e.g.*, additive Gaussian noise), sensitivity (*e.g.*, layer-wise randomisation), and faithfulness (*e.g.*, cumulative input masking, does this assumption hold in practice? In Appendix A.5, we extensively analyse the extent to which Assumptions 1-3 hold versus fail across various explanation methods, and NN models.

Notably, we find that Assumptions 1-3 are systematically violated in practice. While this is expected due to the inherent non-linearity embedded in f , it has significant consequences for the validity of existing evaluations (Definitions 2-4). Evaluation outcomes may be misleading when explanation robustness is enforced for models that fundamentally lack it (Chalasan et al., 2020; Tan & Tian, 2023; Agarwal et al., 2022a), or when faithfulness scores are attributed to explanation quality without considering OOD scenarios (Hase et al., 2021; Hesse et al., 2024). In Section 3.3, we propose a mitigation strategy to address this issue.

3 A Unifying Perspective

With clear definitions of robustness, sensitivity, and faithfulness evaluations (Section 2.3), we may now explore their shared attributes, and outcomes. In the following, we discuss the unifying aspects of these evaluations, and introduce a novel definition to evaluate faithfulness, which integrates these distinct evaluations into a single criterion of explanation quality.

3.1 Unifying Attributes

Upon formalising the evaluation criteria (Definitions 2-4), a notable observation is that robustness, sensitivity, and faithfulness exhibit common attributes. Each of the evaluative criteria (1) introduces a *perturbation* of a specific magnitude ξ , (2) *measures* the functional effect, and (3) *interprets* the results, *i.e.*, the quality estimate q . Also, the evaluation is performed under distinct *model assumptions* about its distortion \mathbf{D}_f . We refer to Table 1 for a summary of these findings.

In Figure 2 (A), we illustrate these theoretical similarities on a graph, with axes corresponding to the shared attributes ξ , and q . Here, we can observe that robustness evaluation (*green*) involves minimal perturbation with a small difference in expected explanation output (or low q). Sensitivity (*red*) employs substantial perturbation, expecting a significant difference in explanation output (or high q). Faithfulness (*blue*) uses cumulative perturbation of Z steps, evaluating the corresponding variations in model output. By placing the different perspective of explanation quality onto Figure 2 (A), and thereafter examining the positions of the post-perturbed instances $\hat{\mathbf{s}} \in \mathcal{S}$, we can advance our understanding of how the criteria relate to one another: specifically, that diverse evaluation methods can be unified under a shared conceptual framework.

Table 1: A concise overview of the attributes of the *robustness*, *sensitivity*, and *faithfulness* evaluations. The last row presents our *unified* evaluation proposal (Def 5), whose theory and practical implementation are described in Section3 and Section5, respectively.

EVALUATION (Ψ) (DEFINITIONS 2-4)	Step 1. PERTURBATION; MAGNITUDE (EQUATIONS 4-5)	Step 2. MEASUREMENT (DEFINITION 6)	Step 3. INTERPRETATION (EQUATION 7-9)	MODEL ASSUMPTIONS (ASSUMPTIONS 1-3)
ROBUSTNESS (Ψ^{RO})	$\mathcal{P}_X(\mathbf{x}); \xi^{RO} \leq \alpha$	$D_\phi(\mathbf{x}, \hat{\mathbf{x}})$	$q^{RO} \leq \varepsilon_{D_\phi}^{RO}$	$D_f \leq \varepsilon_{D_f}^{RO}$
SENSITIVITY (Ψ^{SE})	$\mathcal{P}_\Theta(\theta); \xi^{SE} > \alpha$	$D_\phi(\theta, \hat{\theta})$	$q^{SE} > \varepsilon_{D_\phi}^{SE}$	$D_f > \varepsilon_{D_f}^{SE}$
FAITHFULNESS (Ψ^{FA})	$\mathcal{P}_X(\mathbf{x}, z); \xi_1^{FA} \leq \xi_2^{FA} \leq \dots \leq \xi_Z^{FA}$	$f(\hat{\mathbf{x}}^z, \theta)$	$q^{FA} = \nu(\hat{q}^{FA})$	$D_f^1 \leq D_f^2 \leq \dots \leq D_f^Z$
UNIFIED (Ψ^{GEF})	$\mathcal{P}_\Theta(\theta, z); \xi_1^{GEF} \leq \xi_2^{GEF} \leq \dots \leq \xi_Z^{GEF}$	$D_\phi(\theta, \hat{\theta})$, AND $D_f(\theta, \hat{\theta})$	$\rho(d_f, d_\phi) \approx 1$	NONE

3.2 Unifying Outcomes

Another point of unification emerges when considering the outcomes of these evaluation criteria, and how they interact in practice. In Figure 2 (B), similar to the traditional confusion matrix (in ML) or contingency table (in statistics), we provide a visual representation of the possible model, and explanation outcomes, post-perturbation. Although model, and explanation outcomes are typically continuous in reality—for conceptual clarity, we classify them into four distinct quadrants: true positive (TP), true negative (TN), false positive (FP), and false negative (FN). A key benefit of *discretising* evaluation outcomes in this way, is that we can distinguish between *aligned*, and *misaligned* explanation behaviour:

- *Aligned outcomes (TP + TN)*. The green quadrant represents outcomes where the explanation, and model agree, indicating explanation robustness, *i.e.*, $e = \hat{e}$, and $y = \hat{y}$, and satisfying Assumption 1. Conversely, the red quadrant contains outcomes where both explanation, and model outputs differ, reflecting explanation sensitivity, *i.e.*, $e \neq \hat{e}$, and $y \neq \hat{y}$, and satisfying Assumption 2. Explanation faithfulness is achieved when evaluation outcomes are aligned over Z steps (Assumption 3).
- *Misaligned outcomes (FP + FN)*. The orange quadrants highlight misalignment between ϕ , and f . The top-left quadrant shows explanation dissimilarity despite prediction stability (*i.e.*, $y = \hat{y}$, and $e \neq \hat{e}$), failing Assumption 2. The bottom-right quadrant shows explanation similarity despite a prediction change (*i.e.*, $y \neq \hat{y}$, and $e = \hat{e}$), failing Assumption 1.

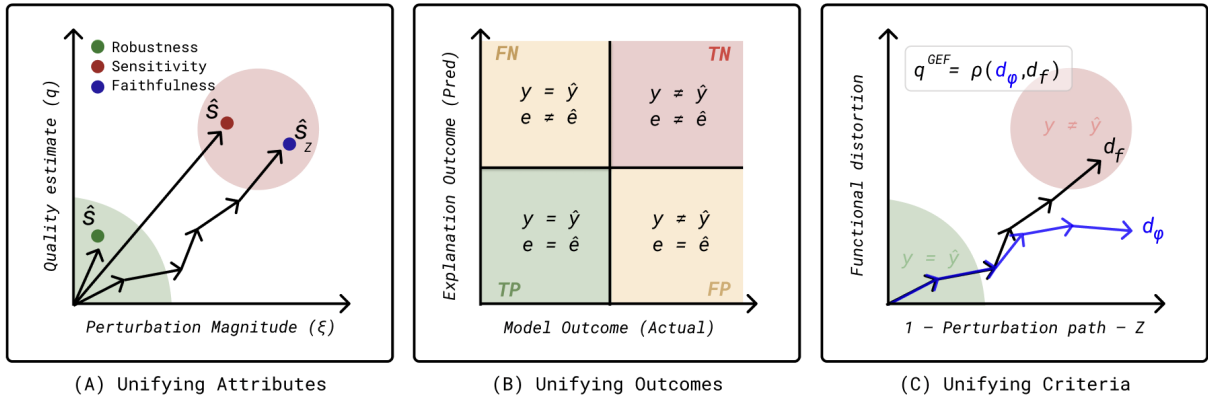


Figure 2: Intuition behind the relationship between robustness, sensitivity, and faithfulness evaluations. (A) illustrates the shared attributes, *i.e.*, perturbation magnitude ξ , and quality estimate q that unifies robustness (green), sensitivity (red), and faithfulness (blue) evaluations. (B) displays a confusion matrix of *discretised* model, and explanation outcomes, with green, and red quadrants indicating aligned behaviour, and orange quadrants showing misalignment. (C) shows our proposed GEF criteria (Definition 5) which measures explanation to model alignment over diverse evaluation perspectives.

Explanation Faithfulness is Alignment. Our analysis reveals that evaluation according to Definitions 2-4, fundamentally concerns the *alignment* between the explanation, and the model’s behaviour, whether across single (Definitions 2-3) or multiple (Definition 4) perturbation steps. A key observation is that

existing robustness, and sensitivity measures provide a limited view of isolated model conditions: robustness evaluates alignment when the model’s predictions remain stable (TP quadrant), while sensitivity evaluates alignment when predictions change (TN quadrant). Faithfulness (Definition 4) evaluates alignment over Z steps, assuming non-decreasing, monotonic model responses under cumulative perturbations (Assumption 3). These *singular* perspectives require strict adherence to specific model conditions, and consequently fail to evaluate the full behaviour of the explanation function. Next, we propose a unifying criterion for faithfulness evaluation. We refer to Figure 2 (C) for an illustration.

3.3 Unifying Criteria

We extend, and generalise the current faithfulness criterion (Definition 4) by integrating the robustness, and sensitivity evaluations into a combined criterion, that is free from restrictive model assumptions. Using a series of Z perturbations, we measure explanation alignment across a spectrum of model outcomes—from cases where model predictions remain consistent, *i.e.*, $y = \hat{y}$, to those where predictions diverge, *i.e.*, $y \neq \hat{y}$. In this way, a generalised definition of explanation faithfulness is obtained.

Definition 5 (Generalised Explanation Faithfulness) Let $\mathbf{d}_f = [\mathbf{D}_f^1, \mathbf{D}_f^2, \dots, \mathbf{D}_f^Z]$ and $\phi = [\mathbf{D}_\phi^1, \mathbf{D}_\phi^2, \dots, \mathbf{D}_\phi^Z]$ be the model, and explanation distortion vectors, where \mathbf{D}_f^z , and \mathbf{D}_ϕ^z are distortion quantities of the z^{th} step along a perturbation path $z \in [1, Z]$, from robustness at $z = 1$ to sensitivity at $z = Z$ such that $\forall y, \hat{y} \in Y$:

$$(z = 1 : y = \hat{y}) \quad \text{and} \quad (z = Z : y \neq \hat{y}),$$

where \hat{y} , and y are perturbed versus unperturbed model outputs, respectively. Let Ψ^{GEF} be a quality estimator that yields estimates $q^{GEF} \in \mathbb{R}$ via the correlation coefficient $\rho : \mathbb{R}^Z \times \mathbb{R}^Z \rightarrow \mathbb{R}$ such that $q^{GEF} = \rho(\mathbf{d}_f, \phi)$. An explanation function $\phi \in \mathcal{E}$ is faithful to $f \in \mathcal{F}$ if:

$$q^{GEF} \approx 1. \tag{10}$$

With Equation 10, we define a quality estimator Ψ^{GEF} that yields values ranging between $[-1, 1]$, with a value of 1 implying perfect generalised faithfulness, 0 suggesting an absence of it, and -1 an inverse relationship. GEF estimation is, therefore, *threshold-free* in the sense that the correlation coefficient directly indicates the quality of the explanation, eliminating the need for arbitrary cut-offs. Note that in Definition 5, we implicitly rely on predicted class c to generate the perturbed logit \hat{y} as the target for the explanation, and model distortion. In Appendix A.1.2, we discuss a broader application of GEF where the targets \hat{y} , and y are replaced by any c^{th} neuron within a layer $l \in [1, L]$ of a feed-forward model. Moreover, Definition 5 applies to a wide range of explanation functions, as discussed in Section 2.1. The choice of ρ , and perturbation applied to construct the distortion vectors depends on the practical implementation (Section 5.2).

Remarks. Our definition shares similarities with faithfulness estimation (Definition 4) in that it assesses explanation quality along a perturbation path. However, it fundamentally differs by focusing on *general alignment* rather than a specific scenario of measuring the *magnitude* of model response to *cumulative input* perturbation. A key benefit of our proposal is that we use the model distortion to *anchor* the expectations of the explanation distortion, and as such, eliminate the need to rely on arbitrary thresholds. In this way, the evaluation will be grounded in the exact functional response of the model, and thus resilient to OOD scenarios: expecting small explanation distortions only when model distortions are small, and vice-versa.

Theoretical Benefits. A good faithfulness measure should assign low scores to unfaithful explanations, and high scores to faithful explanations. In Appendix A.1.3, we prove that a linear model $f = \theta \mathbf{x} + c$ where θ acts as the explanation, attains a perfect faithfulness score, *i.e.*, $q^{GEF} = 1$ with GEF. Conversely, unfaithful explanations are penalised by GEF. For instance, constant explanations that generate no distortion, *i.e.*, $\mathbf{D}_\phi(\mathbf{e}, \hat{\mathbf{e}}) = 0$, pass the conventional robustness test $q^{RO} \leq \varepsilon_{\mathbf{D}_\phi}^{RO}$ (Definition 2), but correctly fails in the GEF formulation. Similarly, random explanations (*e.g.*, generated by uniform sampling, *i.e.*, $\hat{\mathbf{e}}_i \sim \mathcal{U}(0, 1)$) produce maximal distortion (Binder et al., 2022), and thus generally pass the sensitivity test, *i.e.*, $q^{SE} > \varepsilon_{\mathbf{D}_\phi}^{SE}$ (Definition 3) but fails in our definition. We provide proofs for both cases in Appendix A.1.1, and outline empirical evidence in (Section 4.2).

4 A Geometric Perspective

With an advanced understanding of explanation faithfulness (Section 3), we can more systematically study the behaviour of the explanation function. Without assuming the restricted model conditions (Assumptions 1-3) are met, which are often violated in practice (Appendix A.5), we can more objectively measure the *true* explanation function behaviour. We can formalise questions such as: do explanation functions in real-world evaluation scenarios align or misalign with their model? How do different types of perturbation impact model and explanation functions? In the following, and in Appendix A.6, we empirically examine such questions across common explanation methods (Section 2.1) for different NN models. Guided by differential geometry, we provide theoretical considerations on the impact of geometry (Section 4.2).

4.1 Explanation Alignment Patterns

To empirically analyse whether explanation functions are aligned with their model, we study how the distortions of various local, and global explanation functions, and models change under perturbation. Here, we use additive Gaussian noise, *i.e.*, $\nu_i \sim \mathcal{N}(0, \sigma)$ to generate perturbed inputs $\hat{\mathbf{x}}_i = \mathbf{x} + \nu_i$, with a standard deviation σ increasing until the model behaves randomly (*i.e.*, with an accuracy equal to $1/C$) using $Z = 10$ perturbation steps. We refer to Table 2 in Section 6 for details regarding datasets, and models, and to Appendix A.6 for extended results.

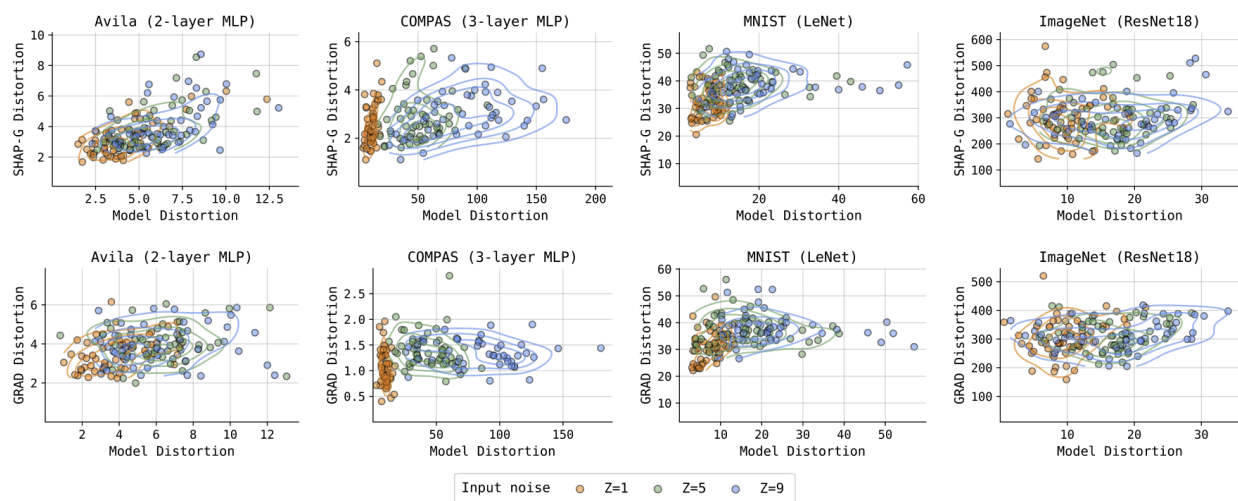


Figure 3: Model (x-axis), and explanation distortions (y-axis) under varying levels of additive Gaussian input noise for vision, and tabular tasks. The scatter points represent individual samples, coloured by perturbation magnitude ($z=1$, $z=5$, $z=9$), with overlapping contours highlighting the relative alignment patterns. The top, and bottom rows represent *GradientSHAP* (SHAP-G), and *Gradient* (GRAD) explanations, respectively.

Figure 3 (top, and bottom) presents the results, which can be interpreted as a continuous analogue of the confusion matrix presented in Section 3.2. The scatter points, coloured by perturbation magnitude, reveal that ϕ , and f rarely align fully. Instead, the relative alignment varies with both the model, and the explanation functions. Here, we include *GradientSHAP* (SHAP-G) (Lundberg & Lee, 2017), and *Gradient* (GRAD) (Morch et al., 1995; Baehrens et al., 2010)) on top, and bottom rows in Figure 3, respectively, with more results in Appendix A.6). The overlapping contours (*e.g.*, Avila results in Figure 3) underscore a simple but nonetheless systematically overlooked aspect of perturbation-based evaluation (Section 2.3): that a uniform perturbation of its inputs may affect highly non-linear systems in a non-uniform way. If the effects were uniform, the system would likely be linear.

4.2 The Impact of Geometry

By considering the geometric nature of the spaces these functions inhabit, we can understand the observed misalignment better. In differential geometry, each space—whether it is the model output space \mathcal{Y} or the

explanation output space \mathcal{E} —can be viewed as a manifold with its unique geometric characteristics (Lee, 2012). When a perturbation is applied, a new point on these manifolds may be accessed, and then, when functional distortion (Definition 1) is computed in each space, we are effectively computing a distance between two points on each manifold. For example, with model parameter perturbation, *i.e.*, $\hat{\theta} = \mathcal{P}_\Theta(\theta)$ (Equation 4) we obtain perturbed model outputs \hat{y} (or, a logit \hat{y}) given $\hat{y} = f_{\hat{\theta}}(\mathbf{x})$. From this, model distortion (Definition 1) is calculated using, *e.g.*, Euclidean distance between the original, and the perturbed instance.

A key observation is that, when distances in two different spaces are directly compared, we ignore the fact that manifolds have their own separate geometric characteristics which are distinct, where distances in one space not necessarily reflect equivalent distances in another. In direct comparisons such as correlation (Ancona et al., 2018; Bhatt et al., 2020) or Lipschitz calculations (Alvarez-Melis & Jaakkola, 2018a; Agarwal et al., 2022a), a global flat metric is assumed. We refer to Figure 4 (A), and (B) for an illustration of the problem of ignoring the impact of geometry. As a result, the quality estimation of the explanation may be misleading.

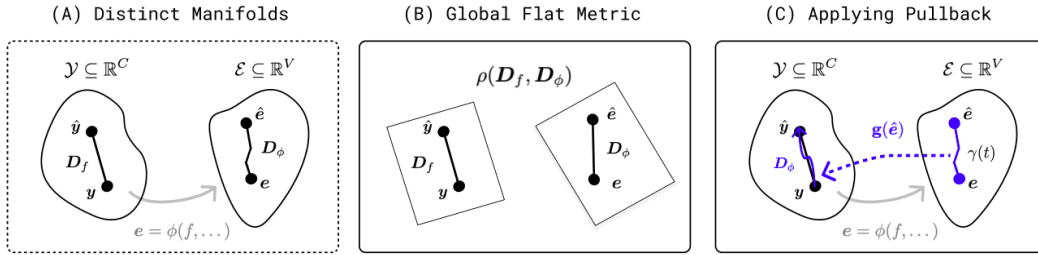


Figure 4: An illustration of the relationship between the manifolds of the model, and explanation. (A) shows how the explanation function maps between the model, and explanation spaces, \mathcal{Y} , and \mathcal{E} . (B) displays the problem with directly comparing distortions across spaces, assuming a global flat metric. (C) illustrates the pullback operation using metric tensor \mathbf{g} to adjust distortions in \mathcal{E} for comparison in \mathcal{Y} .

4.3 Reconciling Geometric Discrepancies

To enable a geometrically sound comparison between explanation, and model distortion, the aim is to recompute D_ϕ to incorporate the non-linear mappings used in generating explanations. This can be achieved by mapping the distortion from the explanation space \mathcal{E} to the model space \mathcal{Y} , effectively “pulling back” the measured distance into \mathcal{Y} (see Ch. 11 of Lee (2012) for further details). Guided by differential geometry, we create a metric tensor \mathbf{g} that serves as this pullback onto \mathcal{Y} . This process is illustrated as $\mathbf{g}(\hat{e})$ in Figure 4 (C).

To construct the metric tensor \mathbf{g} , we consider an infinitesimal neighbourhood around the parameter perturbation $\theta + du$, for a fixed \mathbf{x} , and y . By applying a first-order Taylor expansion in this neighbourhood, we obtain

$$\phi(f_{\theta+du}, \dots) \approx e + J_f du, \quad (11)$$

where $J_f \in \mathbb{R}^{V \times C}$ is the Jacobian for fixed input \mathbf{x} , with elements $J_{i,j} = \frac{\partial e_i}{\partial f_j}$. We use f_j as shorthand for $f_j(\mathbf{x})$. Effectively, $\theta + du$ yields a new perturbed model $f_{\hat{\theta}}$ which is computed with model parameter perturbation (Section 5.1). With Equation 11, we can compute the elements of the pullback tensor $\mathbf{g} \in \mathbb{R}^{V \times V}$ as the sum of the resulting changes in each explanation element e_v *w.r.t.* the changes in each model element f_j

$$g_{i,j}(e) = \sum_{v=1}^V \frac{\partial e_v}{\partial f_i} \frac{\partial e_v}{\partial f_j}. \quad (12)$$

Thus, Equation 12 captures the sensitivity of ϕ to model output changes, with \mathbf{g} corresponding to the squared Jacobian $\mathbf{g} = J_f^\top J_f$. In this way, we can obtain a more reliable measurement of distances in the *pseudo-Riemannian manifold* $(\mathcal{Y}, \mathbf{g})$ of space \mathcal{Y} .

With the pullback metric tensor \mathbf{g} in place, we can measure explanation distortion that is equivalent to computing the path length under the induced parameter changes in the “pulled-back” space

$$\mathbf{D}_\phi := L(\gamma) = \int_0^1 \frac{d\gamma(t)}{dt}^\top g_{\gamma(t)} \frac{d\gamma(t)}{dt} dt, \quad (13)$$

where $\gamma(t)$ is a path between endpoints $\mathbf{e}, \hat{\mathbf{e}} \in \mathcal{E}$ derived from the original, and perturbed models, respectively. Here, t denotes the step size. Now, we replace our original definition of explanation distortion $\mathbf{D}_\phi = \delta(\mathbf{e}, \hat{\mathbf{e}})$ (Definition 1) with the *total accumulated* distortion along the path, *i.e.*, Equation 13. Here, longer paths correspond to greater distortions. Upon taking this geometric perspective, we can study \mathcal{Y} using extrinsically-defined geometry, and contrast it with the simpler assumption of a flat, intrinsic Euclidean metric. As a result, \mathbf{D}_ϕ , and \mathbf{D}_f are more fairly compared in the same space.

5 Method: From Theory to Practice

While our unified theory (Section 3), and solution to reconcile geometric discrepancies in measurement (Section 4), provide first steps towards resolving issues in perturbation-based evaluation, many practical concerns remain regarding the choice of perturbation. In this section, describe how to reliably translate our theory (Definition 5) to practice—we propose a general-purpose, task-agnostic perturbation technique based on model parameter scaling (Section 5.1), and introduce the full evaluation algorithms, *i.e.*, **GEF** and **Fast-GEF**(Section 5.2).

5.1 Selecting Perturbation Strategy

While all perturbation-based evaluations inherently require parameterisation, input-based perturbation (Definitions 2 and 4), has proven particularly challenging to calibrate (Sturmfels et al., 2020; Haug et al., 2021). Without ground truth labels, selecting parameters such as patch size, pixel, or token replacement strategies is typically based on researchers’ judgment. Small changes to input parameters have been shown to significantly impact evaluation outcomes (Brunke et al., 2020; Brocki & Chung, 2022; Rong et al., 2022; Blücher et al., 2024), raising concerns about reliability.

Moreover, perturbing on the input space is not only impractical from a practitioner’s standpoint but also compromises impartiality—if parameters must be adjusted for each model, and dataset, how can task-specific confounds be controlled? In Appendix A.5.1, we provide empirical evidence for the existence of confounds in faithfulness evaluations (Definition 4).

Researchers need a general-purpose, dataset, and architecture-agnostic perturbation strategy that facilitates evaluation across distinct explanation approaches (*e.g.*, local, and global methods), and magnitudes, *i.e.*, ξ . Following Bykov et al. (2022), we propose a simple perturbation strategy in the following.

Model Parameter Scaling. *Introduce perturbations $\forall z \in [1, Z]$ by scaling parameters $\theta \in \mathbb{R}^W$ with Gaussian noise $\eta_i \sim \mathcal{N}(\mathbf{1}, \sigma_z^2 \mathbf{1})$, and $\sigma_z^2 \in \mathbb{R}^+$ such that $\hat{\theta}_z = \theta \cdot \eta_i$, yielding a perturbed model $f_{\hat{\theta}_z}$.*

By systematically perturbing model parameters instead of the input, from low to high magnitudes with incremental increases of σ_z^2 , ranging from robustness at $z = 1$ to sensitivity at $z = Z$, explanation behaviour is evaluated comprehensively, and agnostically across tasks. With $\xi := \delta(f(\mathbf{x}), f_{\hat{\theta}_z}(\mathbf{x}))$ (Equation 5), we can measure the perturbation impact at each z^{th} step so that robustness, *i.e.*, $y = \hat{y}$, and sensitivity, *i.e.*, $y \neq \hat{y}$ criteria are fulfilled (Definition 5). Our approach contrasts with the model parameter randomisation procedure of Adebayo et al. (2018), which proposes layer-wise *randomisation* in a top-down order, an approach that faces methodological concerns (Sundararajan & Taly, 2018; Binder et al., 2022; Kokhlikyan et al., 2021; Yona & Greenfeld, 2021). For an illustration of how model parameter scaling affects the classifier’s decision boundary, we refer to Fig. 1 of Bykov et al. (2022).

5.2 Introducing GEF Evaluator

From an algorithmic perspective, three steps are necessary to perform the evaluation. First, given a model, and a test set of input-output pairs, we generate perturbed models $f_{\hat{\theta}_1}, \dots, f_{\hat{\theta}_Z}$ given Z sets of parameters $\hat{\theta}_1, \dots, \hat{\theta}_Z$ along a perturbation path (see Algorithm 1, line 6). Then, for each model $f_{\hat{\theta}_z}$, we compute the model, and explanation distortion quantities, *i.e.*, \mathbf{D}_f^z , and \mathbf{D}_ϕ^z , using the pullback tensor \mathbf{g} (lines 7, 9, and 10). Finally, distortion vectors are constructed, and correlated using $\rho(\mathbf{d}_f, \mathbf{d}_\phi)$ (lines 14, and 15). Due to the stochastic nature of model perturbation, we repeat this process M times to average out the effects. We refer to Figure 5, and Algorithm 1 for an overview of the steps involved. An ablation study on hyperparameter choices is provided in Appendix A.7.

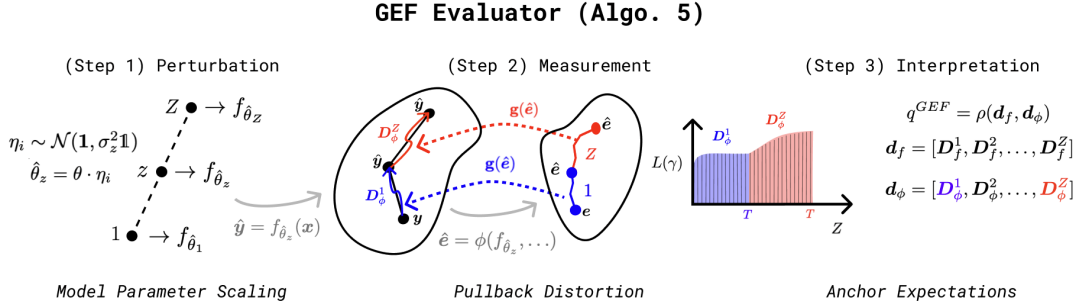


Figure 5: The three steps of GEF evaluation (Algorithm 1) to estimate generalised explanation faithfulness (Definition 5).

Practical Benefits. Our proposed evaluation (Algorithm 1) provides several practical benefits. First, *anchoring*, negates the need to rely on arbitrary thresholds in evaluation, *e.g.*, when determining a permissible value for the evaluations themselves (Equations 7, 8, and 9) or what perturbation magnitude leads to model alignment for a particular task. Second, perturbing via *model parameter scaling*, at varying intensities combines distinct criteria of explanation quality into a single unified evaluation metric (Section 3.3) that is agnostic to the data, model, and explanation approach. Third, the pullback metric calculation provides a geometrically grounded faithfulness measurement, capturing the true functional impacts of the explanation *w.r.t.* its model.

Algorithm 1 GEF Evaluator

```

1: Require: Model  $f$ , explanation function  $\phi$ , input-prediction pairs  $\mathbf{x}, y \in \mathbf{X}, \mathbf{Y}$  with  $\mathbf{X} \subseteq \mathcal{X}, \mathbf{Y} \subseteq \mathcal{Y}$ 
2: Parameters: Integers  $Z, M, T, K$ , correlation measure  $\rho$ 
3: for  $\mathbf{x}, y$  in range( $\mathbf{X}, \mathbf{Y}$ ) do
4:    $e \leftarrow \phi(f, y, \dots)$ 
5:   for  $z$  in range( $Z$ ) do
6:      $\hat{y} \leftarrow f_{\hat{\theta}_z}(\mathbf{x})$ 
7:      $\mathbf{D}_f^z \leftarrow \delta(y, \hat{y})$  // Equation (6)
8:     if Fast-GEF then
9:        $\mathbf{D}_\phi^z \leftarrow \delta(e, \hat{e})$  with  $\hat{e} \leftarrow \phi(f_{\hat{\theta}_z}, \dots)$  // Equation (6)
10:    else
11:       $\mathbf{D}_\phi^z \leftarrow \text{compute\_path\_length}(f_{\hat{\theta}_z}, \mathbf{x}, y, T, K)$  // Equation (13)
12:    end if
13:  end for
14: Construct:  $\mathbf{d}_f \leftarrow [D_f^1, D_f^2, \dots, D_f^Z]$ , and  $\mathbf{d}_\phi \leftarrow [D_\phi^1, D_\phi^2, \dots, D_\phi^Z]$ 
15: Calculate:  $q^{GEF} \leftarrow \rho(\mathbf{d}_f, \mathbf{d}_\phi)$ 
16: Return:  $q^{GEF}$ 
17: end for

```

Implementation Details. Unless stated otherwise, we use *Euclidean distance* for δ in the functional distortion calculations (Definition 1), and define ρ using *Spearman Rank Correlation*, assessing the degree of

monotonic relationship between the distortion quantities. For the experiments, we set $Z = 5$ (see discussion of the influence of Z in Appendix A.1.4) but it is a choice that can be flexibly updated in the open-source implementation. In Appendix A.2, we provide further details on the implementation, including how to generate the perturbation path (line 2), and how to tune parameters (line 2). This also includes information on how we compute the path length (line 11); where we follow an approximation procedure outlined in Equations 18, and 19.

5.3 Balancing Computational Constraints

While the pullback operation ensures a fair geometric comparison of distortion quantities, its use of high-dimensional Jacobian calculations, and integral steps (Equation 12) also increases computational demands. To accommodate evaluation contexts involving large model architectures or high-dimensional explanations, we offer an alternative method. For a faster yet *naive* approximation of explanation quality, we omit the pullback operation, and instead define D_ϕ according to Equation 6. This approach, entitled **Fast-GEF**, is less computationally demanding, and complements the *exact* approach with pullback, entitled **GEF**, providing a geometrically sound quality estimate.

Choosing between GEF or Fast-GEF. Users can choose between these methods based on their specific computational constraints and demands for accurate quality estimates. We recommend using **GEF** wherever possible due to its ability to account for manifold-specific distortions. However, **Fast-GEF** provides a computationally efficient alternative that is suitable for large-scale tasks or resource-constrained environments. Empirical results show that while the **GEF** or **Fast-GEF** may diverge in individual estimates (Appendix A.6), they often share categorical rankings of explanation methods (Appendix A.8.2).

6 Experiments

Our experiments aim to answer the following questions:

- (Q1) Are unified, **GEF** and **Fast-GEF** evaluations more empirically reliable than competitive singular approaches?
- (Q2) How does generalised faithfulness of local, and global explanation methods compare across distinct data domains?
- (Q3) How faithful are LLMs as a top- K token post-hoc explainer for NLP classifications?
- (Q4) Are SAEs generally faithful, and does more capacity in their width improve their faithfulness?

To answer these questions, we select a diverse set of datasets, model architectures on tabular, vision, and NLP classification tasks. See Table 2 for an overview. Our experiments evaluate the faithfulness of various explanation approaches, as detailed below.

Global, and Local Methods. For global methods, we include feature visualisation techniques with different regularisation, and optimization procedures: *Deep-Viz* (DV) (Yosinski et al., 2015), *Magnitude Constrained Optimization* (MACO) (Fel et al., 2024), and *Fourier preconditioning* (FO) (Olah et al., 2017). Optimization steps are set to 50, 100, and 250, otherwise, default values are used as provided in the respective publications ((Fel et al., 2024), and (Nguyen, 2020)). For local methods, two variants of *Layer-wise Relevance Propagation* (LRP), the ε -rule (LRP- ε) (Bach et al., 2015) with $\varepsilon = 1e^{-6}$, and the z^+ -rule (LRP- z^+) (Montavon et al., 2017) are employed. Also, we include several gradient-based approaches such as *Gradient* (GRAD) (Morch et al., 1995; Baehrens et al., 2010), *Saliency* (SAL) (Simonyan et al., 2014), *Input \times Gradient* (IXG) (Shrikumar et al., 2016), *GradCAM* (G-CAM) (Selvaraju et al., 2020), *Guided Back-propagation* (GBPG) (Springenberg et al., 2015), *SmoothGrad* (SMG) (Smilkov et al., 2017) with 10 noisy samples, and noise level $0.1/(x_{\max} - x_{\min})$, *Integrated Gradients* (INT-G) (Sundararajan et al., 2017) with 10 iterations, and zero baseline. For NLP tasks, we evaluate *LayerIntegratedGradients* (L-INTG) explanations *w.r.t.* the first embedding layer. Two Shapley-based algorithms (Lundberg & Lee, 2017) are included: *GradientSHAP* (SHAP-G) with 10 samples, and *PartitionShap* (SHAP-P) for NLP tasks.

LLM-x Methods. An emerging research area in explainability uses separate LLMs to generate post-hoc attributions for important features of a given model (Bills et al., 2023; Kroeger et al., 2023; Krishna et al., 2023; Amara et al., 2024). We create LLM-x explanations by prompting Gemma-2B-IT (Mesnard et al., 2024) to rank the top- K most important tokens given a textual input, which is then parsed, decoded, and mapped to input tokens, producing binary attribution vectors. LLM prompts describe the model’s classification task, and prediction certainty before, and after model perturbation (Section 5.1). The temperature is set to 0 for deterministic outputs. Varying synonyms, the order of tokens, and the number of top- K values to $\{5, 10\}$ contribute to the robustness of our findings. The full explanation methodology is described in Appendix A.4.4 with an illustration in Figure A.2.

Sparse Autoencoders. SAEs have lately come forth as an interpretability method for understanding the internal representations of LLMs (Templeton et al., 2024; Huben et al., 2024). In our work, we generate SAE explanations using **Gemma-Scope** (Lieberum et al., 2024), pretrained on the residual block representations of the Gemma-2-2B model. Explanations are saved for all 26 layers at both 16K, and 65K widths. Given the sparsity of the explanation vectors, we use cosine distance to compute explanation distortion, defined as $1 - \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$, as it effectively measures similarity regardless of magnitude. Appendix A.4.4 provides a detailed description of the generation process for each SAE explanation.

Control Variants We also evaluate the faithfulness of two control variants: a random explanation (RAN) sampled from a uniform distribution, $\hat{e}_i \sim \mathcal{U}(1, 0)$, and a top- K control variant (RAN- K) with K non-zero attributions, each equal to 1. Unless specified, all experiments evaluate 250 explanations for the logit of the predicted class. For comparability, global, and local explanations are normalised by dividing the attribution map by the square root of its average second-moment estimate (Equation 21) (Binder et al., 2022), with further explanation preprocessing details provided in Appendix A.4.4. For metric implementation, and meta-evaluation, we use the **Quantus** (Hedström et al., 2023b), and **MetaQuantus** (Hedström et al., 2023a) libraries, respectively. Further experimental details for Q1, Q2, and Q3 are provided in Appendix A.8.1, A.8.4, and A.8.5, respectively.

Table 2: An overview of datasets, and models, with references in Appendix A.4. A semicolon separates models used per dataset.

MODALITY	DATASET (N. CLASSES)	MODEL (SIZE)	ACC. %	SOURCE	EXPL. DIM	TASK
TEXT	SMS SPAM (2)	BERT-TINY FT (4.4M)	98.0	HF	128	SPAM
	IMDB (2)	PYTHIA FT (7.6M); GEMMA-2 (2B)	86.4; 95.6	HF	512	SENTIMENT
	SST-2 (2)	BERT-TINY FT (4.4M)	98.0	HF	59	SENTIMENT
VISION	IMAGENET-1K (1000)	RESNET18 (11.7M)	89.1	TORCHVISION	50176	OBJECT
	PATH (9)	MEDCNN (235.2K)	84.3	LOCAL	784	PATHOLOGY
	DERMA (7)	MEDCNN (234.9K)	73.2	LOCAL	784	DERMATOLOGY
	MNIST (10)	LENET (61.7K)	97.7	LOCAL	784	DIGIT
	FMNIST (10)	LENET (61.7K)	87.7	LOCAL	784	FASHION
TABULAR	ADULT (2)	3-LAYER MLP (11.7K); LR (28)	84.6; 83.3	OPENXAI	13	INCOME
	COMPAS (2)	3-LAYER MLP (11.1K); LR (16)	85.0; 85.3	OPENXAI	7	RECIDIVISM
	AVILA (12)	2-LAYER MLP (3.5K)	80.8	LOCAL	10	LETTER

6.1 Measuring Empirical Reliability

To investigate the empirical reliability of **GEF** and **Fast-GEF** evaluations compared to singular approaches, we perform meta-evaluation, which is the practice of evaluating the evaluation method itself. To this end, we adopt the meta-evaluation methodology from Hedström et al. (2023a), which bypasses the lack of ground truth labels by focusing on *metric consistency* (“does this evaluation method produce similar results under consistent conditions?”). For this, two practical meta-evaluative tests are performed: the Input Perturbation Test (IPT), and the Model Perturbation Test (MPT). Each test returns a meta-consistency (MC) score (see Equation 20), which ranges between $[0, 1]$. Higher values indicate greater reliability. Full meta-evaluation scoring methodology is provided in Appendix A.3. As a sanity check, we also show in Appendix A.8.3 that our proposed evaluators assign low scores to different random control variants, where other metrics fail to do so.

Setup. We benchmark three evaluation methods per criterion. In the robustness category, we include *Relative Input Stability* (RIS), *Relative Representation Stability* (RRS), *Relative Output Stability* (ROS) (Agarwal et al., 2022a). In the sensitivity category, we include *Model Parameter Randomisation Test* (MPRT) (Adebayo et al., 2018), *Smooth MPRT* (sMPRT), and *Efficient MPRT* (EMPRT) (Hedström et al., 2024). In the faithfulness category, we include *Faithfulness Correlation* (FC) (Bhatt et al., 2020), *Pixel-Flipping* (PF) (Bach et al., 2015), and *Region-Perturbation* (Samek et al., 2017). All metrics are mathematically described in Appendix A.4.5. To ensure comparability with the original publication (Hedström et al., 2023a), we run meta-evaluation on the same set of tasks, which includes ImageNet (Russakovsky et al., 2015), MNIST (LeCun et al., 2010) and fMNIST (Xiao et al., 2017) datasets with architectures such as ResNets (He et al., 2016) and LeNets (LeCun et al., 1998) architectures. Each metric evaluates GRAD, SAL, G-CAM, SHAP-G explanations. Further results, and details are provided in Tables A.1, and A.2, and Appendix A.8.5.

Results. Figure 6 (A) shows that our proposed unified methods (GEF and Fast-GEF) achieve the highest overall MC scores, averaged over both MPT, and IPT tests. Our unified methods significantly outperform the most comparable evaluation approach, the faithfulness metrics, which also use Z perturbation steps, with average MC scores of 0.733 compared to 0.601. Although no evaluation method achieves a perfect score (*i.e.*, MC=1), the unified methods still perform comparably to robustness metrics, and surpass sensitivity metrics, with average scores of 0.727, and 0.673, respectively. These results are encouraging as they show that unified methods can achieve high reliability, even when explanation behaviour is evaluated under multiple model conditions, unlike robustness, and sensitivity metrics that focus on a single perspective. While the ROS metric has the highest individual score, this is not statistically significant, and it only offers a limited view of explanation quality. Figure 6 (B) shows that unified metrics excel in MPT, while robustness metrics perform slightly better in IPT. These score differences correspond to robustness metrics using input perturbations, and unified metrics relying on model perturbations. Further details are provided in Appendix A.8.4.

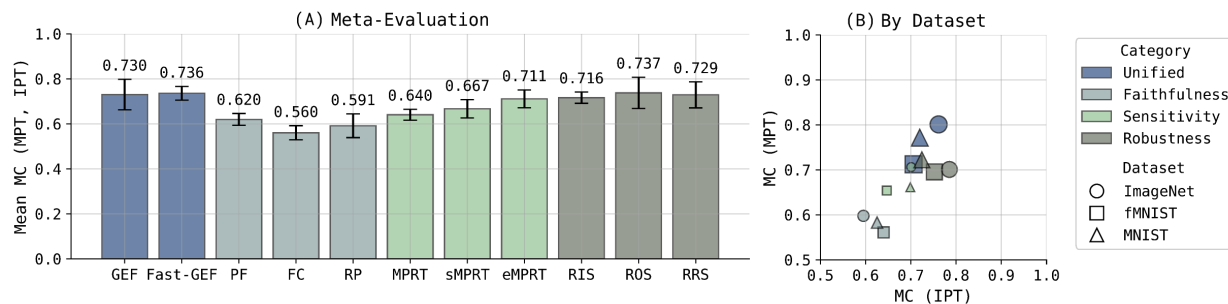


Figure 6: Meta-evaluation, and comparison to established explanation evaluation methods. (A) shows the mean MC scores across MPT, and IPT, aggregated over all datasets, with the error bars showing the standard deviation. (B) displays MC scores aggregated by the test type, and dataset, where the size of the scatter point denotes the standard deviation. GEF scores are computed for fMNIST, and MNIST datasets due to computational constraints.

6.2 Cross-Evaluating Local, and Global Methods

While local, and global explanations serve distinct purposes, and provide different insights *w.r.t.* their model, it is beneficial to compare them side-by-side in a unified view, as they often rely on similar methodological components, such as network gradients (LeCun et al., 1998; Olah et al., 2017). The absence of general-purpose evaluations has however so far prevented such comparison. GEF and Fast-GEF effectively fill this gap, facilitating a first, cross-domain comparative faithfulness benchmarking between global, and local methods. Extended results are provided in Appendix A.8.4.

Figures 7 and 8 provide an overview of cross-domain results for tabular, and vision tasks. For all tabular tasks, GEF estimates are computed. For vision tasks, due to the high computational cost of global methods, Fast-GEF is used to allow for a fair comparison to local methods. As shown in Figure 7, no explanation method is perfectly faithful to its model (*i.e.*, no score equals 1) nor consistently outperforms others across tested tasks. This variation aligns with most benchmarking studies of local linear approximation methods, which rarely identify a single winning method (Hedström et al., 2024; Hesse et al., 2024). Among tested

global feature visualisation methods, MACO generally outperforms FO variants, consistent with Fel et al. (2024). Comparing the faithfulness scores of DV, MACO, and FO reveals that more optimisation steps do not necessarily result in higher explanation faithfulness. All tested methods significantly outperform the random baseline (RAN), which serves as the theoretical lower bound. As expected, RAN produces faithfulness scores centered around zero. In Figure 8 (A), and (B), we observe that RAN explanation distortion quantities are flat, *i.e.*, independent of the model distortion. Tables A.4, and A.5 in Appendix A.8 present the result of Figure 7.

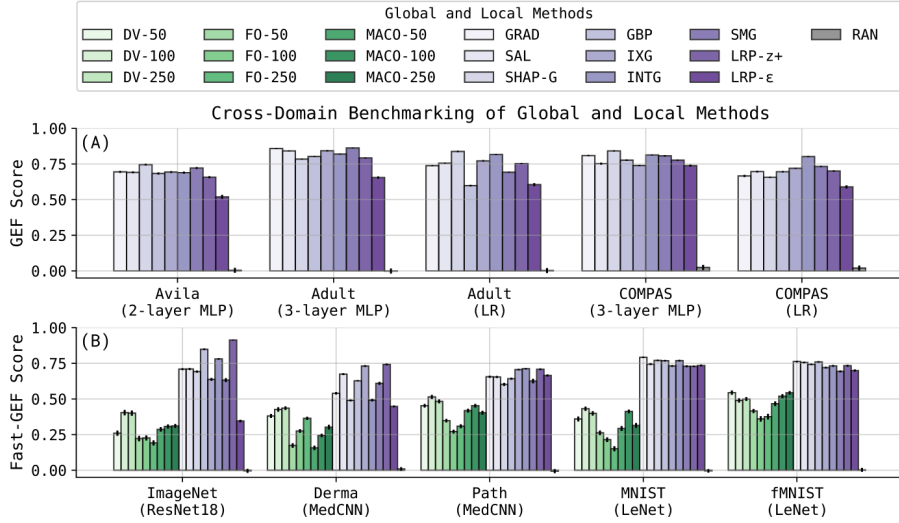


Figure 7: GEF and Fast-GEF results on (A) local across tabular, and (B) local versus global methods across vision tasks. The error bar shows the standard error, *i.e.*, $\frac{\sigma}{\sqrt{N}}$, where σ is the standard deviation, and N is the sample size.

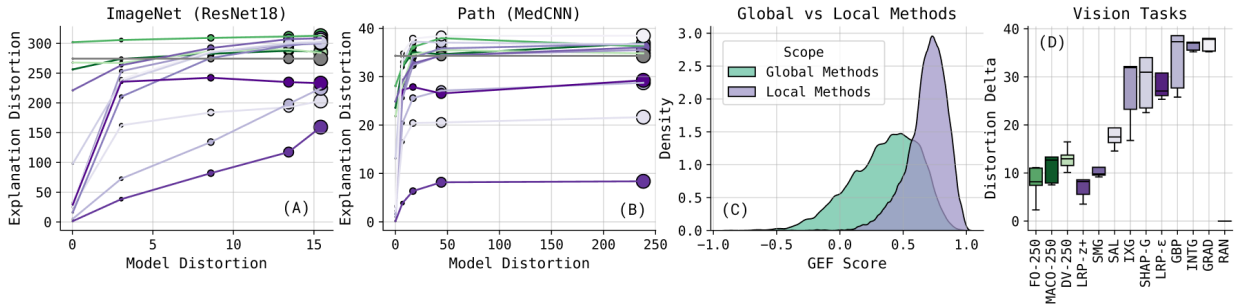


Figure 8: Fast-GEF results for vision tasks. (A), and (B) plot the model, and explanation distortion for ImageNet (ResNet18), and Path (MedCNN) along the perturbation path with $Z = 5$ perturbation steps. Here, global methods (DV, MACO, FO) are selected with 250 optimisation steps. (C) displays the distribution of Fast-GEF for local, and global methods, aggregated over all vision tasks. (D) reports the aggregated difference in explanation distortion between start $z = 1$, and end $z = 5$.

Local Methods are Moderately Aligned. Despite local methods showing imperfect, and highly varying scores across tested models, and datasets, most GEF estimates in tabular tasks, and Fast-GEF estimates in vision exceed 0.5, suggesting that the explanation retains some alignment with its model. This is not surprising given that parameter scaling directly effects the model’s curvature, to which local gradient-based methods are highly sensitive (Dombrowski et al., 2019), thereby instantaneously influencing their responsiveness to perturbation.

Figure 7 (A) shows that some local methods produce distortion outputs nearly monotonically related to its model, particularly at lower magnitudes (*i.e.*, a $z \leq 3$). This finding nuances studies by Adebayo et al. (2018), which provide single-point sensitivity estimates, conclusively reporting low reactivity to parameter randomisation in local methods. Corroborating recent rebuttal works (Yona & Greenfeld, 2021; Sundararajan

& Taly, 2018; Binder et al., 2022) that challenges stark claims of method failure (Adebayo et al., 2018), we find that gradient-based methods are moderately faithful.

Global Methods are Constrained by Regulariser. Figure 8 (C) shows aggregate **Fast-GEF** scores, indicating that global feature visualisation methods typically are less faithful compared to local linear approximation methods. These differences in faithfulness estimates may be attributed to the global methods’ inherent reliance on optimisation procedure (Olah et al., 2017), and NN’s ability to retain its learned features despite perturbation via parameter perturbation (Binder et al., 2022). For reference, DV applies multiple regularisation techniques directly to the image, such as Gaussian blur, and cropping regions based on norm, and pixel contribution, while MACO, and FO regularise the frequency domain representation, with MACO adding an extra layer of regularisation via a predefined magnitude template. As observed in Figure 8 (A), and (B), despite model perturbation, explanation distortions stay relatively flat, with lower distortion deltas compared to most local methods, as displayed in Figure 7 (D). A strongly regularised optimisation procedure may inherently limit the faithfulness of global methods, in favour of a maximally activated neuron response.

6.3 Evaluating LLMs as Post-hoc Explainers

While researchers have recently begun exploring the potential of using LLMs as post-hoc explainers, there is still limited theoretical understanding, and empirical evidence on the general faithfulness of such approach. Can an LLM which is inherently decoupled from the model it seeks to explain, provide faithful outcomes? In our evaluation, we prompt Gemma-2B-IT for a top- K token explanation for a given input, and prediction pair for datasets characterised by short tokenized lengths, *i.e.*, 59 for SST-2, and 128 for SMS Spam. The post-processed binary explanation vectors are then evaluated with **GEF** and **Fast-GEF**. See Appendix A.8.5 for further details, and extended results.

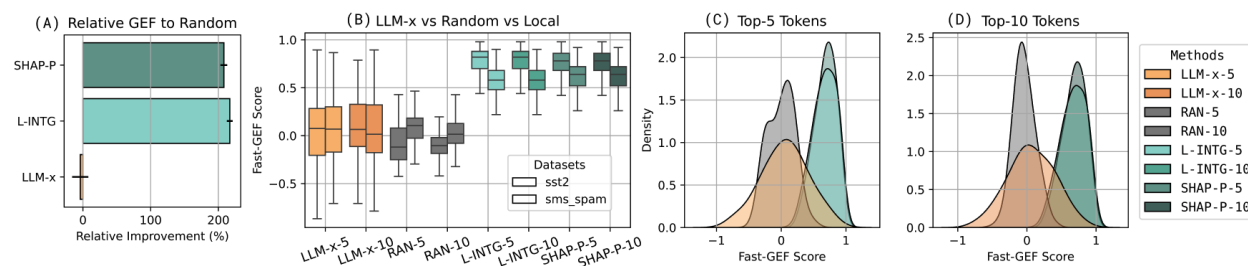


Figure 9: **GEF** (with $M=3$), and **Fast-GEF** results on different top- K explanation NLP tasks. (A) shows the percentage improvement in **GEF** scores relative to RAN, aggregated over all tasks, with error bars showing the standard error. (B) shows the results in the form of box plots for the two datasets with SST-2 (*left*), and SMS Spam (*right*). (C), and (D) show the distribution of **Fast-GEF** scores for top-5, and top-10 explanations respectively, both aggregated over all tasks.

LLM-x Explanations Comparable to Random. Our **GEF** and **Fast-GEF** results in Figure 9 (A), and (B) show that Gemma-2B-IT as an explainer is (i) significantly less faithful than local methods such as SHAP-P, and L-INTG, and (ii) similarly unfaithful as random explainers RAN-5 or RAN-10, on both SST-2, and SMS Spam classification tasks. Figure 9 (C), and (D) demonstrate that these findings generalise over both top-5, and top-10 tokens tasks, aggregated over both datasets. Our results, showing that LLM-x explanations are not more faithful than random, differ from the encouraging results reported by Kroeger et al. (2023), who found GPT-4 to be as faithful as local methods in identifying top- K tokens for tabular tasks. This divergence may naturally stem from variations in the experimental setup, including the specific explanation task, LLM used, methodology to evaluate faithfulness, and prompting strategies, however, it also underscores that the faithfulness of LLM-x is still an open research question. To fully understand the potential of LLMs as explainers, further research with additional LLMs would be beneficial.

6.4 Measuring Faithfulness of Sparse Autoencoders

SAEs are gaining attention for their claimed ability to construct interpretable “monosemantic” features of a given layer of an LLM. Yet, their general faithfulness remains underexplored (Makelov et al., 2024; Mallen

& Belrose, 2024). To this end, we evaluate SAE explanations for Gemma-2-2B model on $N = 250$ samples on the IMDb dataset (Maas et al., 2011). Here, the Gemma-2-2B model is repurposed as a binary classifier by extracting the logits of the “positive” or “negative” classes at the final token position of the prompt. Appendix A.4.4 provides more details.

High Faithfulness Independent of SAE Width. Figure 10 (A) demonstrate that the SAE explanations generally are faithful *w.r.t.* the model’s intermediate representations. **Fast-GEF** scores are consistently above 0.75 except for fluctuations in layers 1, and 16 – 19. No significant difference is observed between 16K, and 65K widths, suggesting that the width of the encoding, *i.e.*, the capacity of the SAE, does not correlate with explanation faithfulness. Moreover, Figure 10 (B) shows that although sparsity of the latent activations decreases in later layers, it does not influence its faithfulness (*i.e.*, $\rho = 0.023$, computed with Spearman Rank correlation). This suggests that a higher activation in SAE latents does not necessarily relate to its measured faithfulness. This raises the question of whether the faithfulness of SAE explanations is inherently scalable within different statistical or qualitative contexts where SAEs are studied. Figure 10 (C)-(G) further illustrates how explanation distortions vary across layers, with values (*y-axis*) increasing in the middle layers.

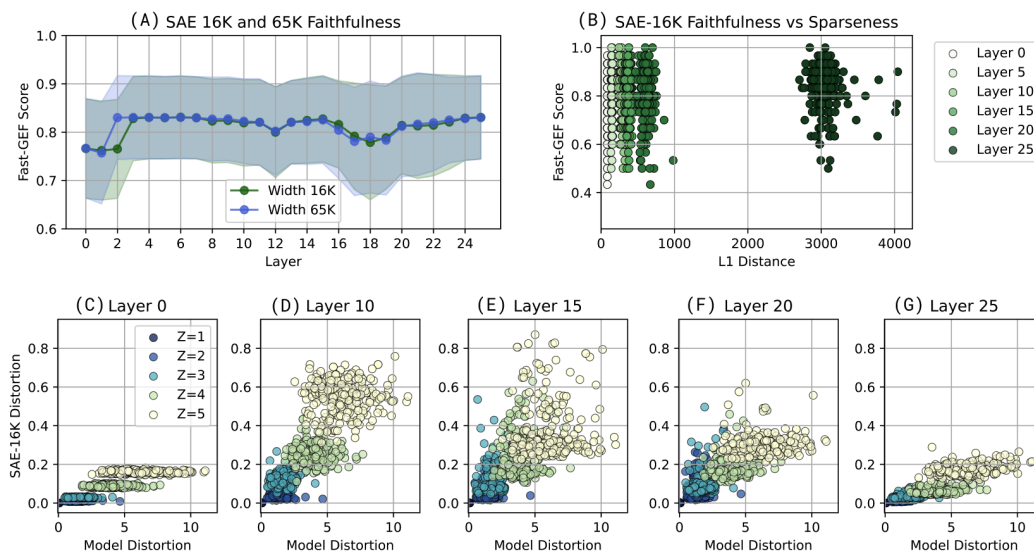


Figure 10: **Fast-GEF** results for SAE explanations on IMDb dataset, using $Z = 5$ perturbation levels, and $M = 3$ models. (A) shows **Fast-GEF** faithfulness scores across layers. (B) shows the sparsity as defined by L1 distance against the **Fast-GEF** indicating no relationship. (C)-(G) illustrates how SAE distortions develop across model layers, coloured by perturbation level.

7 Discussion: Where Are We Going?

With the evolving landscape of interpretability, redefining both the conceptual framework, and the geometric foundations of explanation faithfulness is important. Our work puts forward a long-overdue unification of robustness, sensitivity, and faithfulness evaluations, providing a novel, and urgently needed, revised approach (Definition 5) to evaluate the direct alignment between explanation, and model functions (Section 3). In this work, we address the fundamental flaws of many existing evaluations: a systematic overlook of the intrinsic geometry of non-linear spaces (Section 4). Our solution offers a threshold-free, fair comparison of functional distortions, making our approach not just another evaluation method but a necessary foundation for future interpretability research (Section 5).

Novel, Empirical Insights. In a first-ever cross-domain faithfulness benchmarking of global, and local explanations on vision, tabular, and NLP tasks (Section 6), we learn that tested local explanation methods generally are moderately faithful. We find that global feature visualisation methods are comparatively less faithful, which is an important understanding considering the recent evidence pointing to their general

susceptibility to adversarial manipulation (Geirhos et al., 2023; Bareeva et al., 2024a). While it would be valuable to compare our findings with existing studies, to our knowledge, there is no direct study on the faithfulness of feature visualisations. Existing evaluations focus on the alignment with human preferences or improvement on a downstream task (Borowski et al., 2021; Zimmermann et al., 2021; Krishna et al., 2023; Bareeva et al., 2024b) or similarity to natural samples of the explained class (Fel et al., 2024). Our findings on generalised faithfulness thus provide complementary insights into the quality of feature visualisation as *model* explainers.

Additionally, due to the recent interest LLMs as potential post-hoc explainers (Krishna et al., 2023; Kroeger et al., 2023), we study their faithfulness. We find no improved faithfulness compared to random explanations, and encourage more investigation on this question. Finally, we observe that residual stream SAEs on Gemma-2-2B exhibit generally high faithfulness, with the width having limited influence (*i.e.*, 16K or 65K). Further investigation is required to fully understand the potential of SAEs, and LLM-x as *generally faithful* explainers.

7.1 Limitations

While the results in our paper allow us to claim that our proposed method is more sound geometrically (Section 4), more reliable empirically (Section 6.1), and easier to use practically (Section 5), our evaluation alone does not imply that the explanation quality is sufficient. Without ground truth labels, we cannot assess the statistical validity of an explanation function. An explanation may be estimated to be *generally faithful* but still lack intrinsic value (Bhattacharjee & von Luxburg, 2024) or interpretable qualities (Bordt & von Luxburg, 2024). The need for a thorough, application-grounded assessment of explanation quality that asserts value on a downstream task (Krishna et al., 2023; Lanham et al., 2023) is not eliminated when using GEF. Evaluation using synthetic models with known ground truth (Carmichael & Scheirer, 2023) could complement our proposal.

7.2 Future Work

There are several exciting geometric, and empirical questions worth exploring. The geometric considerations in GEF suggest a deeper examination of the computational trade-offs of computing accurate pullbacks on individual explanation functions, specifically in comparing global versus local methods. In future work, there is opportunity to build on the growing body of research in ML that draws from geometry, and related topics in higher mathematics to deepen our understanding of NNs, and problems to which they are applied (Stephenson et al., 2021; Burns & Tang, 2023; Papamarkou et al., 2024). Recent theoretical studies on LLMs, and transformer models (Hoogland et al., 2024; Burns, 2024) have illustrated how neural activations may arrive at, and utilise “superpositional” encoding strategies (Elhage et al., 2022), which prominently feature considerations or findings of a geometric or topological nature. Continued development of general frameworks, and theories that conceptualise NNs in terms of geometry, and topology (Bianchini & Scarselli, 2014; Hauser & Ray, 2017; Naitzat et al., 2020; Benfenati & Marta, 2023a;b; Burns & Fukai, 2023) will likely facilitate a deeper understanding of both explanations, and evaluations, particularly in relation to the underlying mathematical characteristics of data, optimisation processes, and learned functions.

Recent advances in manifold geometry have introduced tools to analyse how input data modulates internal processing through perturbations (Kvinge et al., 2023). Exploring how explanation faithfulness varies with training data, and how it intersects with the geometric characteristics of the model presents an exciting direction. We also expect models optimised with non-Euclidean methods (Fei et al., 2023) to reveal stronger differences between GEF and Fast-GEF, providing new opportunities to study the interplay between geometry, and faithfulness in explainability.

Lastly, we plan to expand our benchmarking scope to include natural activation-maximisation explanations (Borowski et al., 2021), concept-based explanations like INVERT (Bykov et al., 2023), and non-classification tasks. Given that pullback calculations can be computationally prohibitive for high-dimensional explanations, and highly parameterised models, exploring ways to speed up the Jacobian calculation (Equation 19), and employ adaptive noise schedules would be valuable.

Broader Impact Statement

Interpretability, or XAI, is widely acknowledged as essential for responsible ML. This paper critically examines current evaluation methods from unifying, and geometric perspectives, and proposes improvements. While negative societal impacts are improbable, overreliance on any single evaluation method is not advised.

Acknowledgments

This work was partly funded by the German Ministry for Education and Research (BMBF) through the project Explaining 4.0 (ref. 01IS200551). Additionally, this work was supported by the European Union’s Horizon Europe research and innovation programme (EU Horizon Europe) as grant TEMA (101093003); the European Union’s Horizon 2020 research and innovation programme (EU Horizon 2020) as grant iToBoS (965221); the German Research Foundation (DFG) as research unit KI-FOR 5363 (project ID: 459422098); the state of Berlin within the innovation support programme ProFIT (IBB) as grant BerDiBa (10174498); and BIFOLD (refs. 01IS18025A, 01IS18037A); T.B. thanks Vasiliki Lontou for helpful discussions.

References

- Julius Adebayo, Justin Gilmer, Michael Muelly, Ian J. Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 9525–9536, 2018.
- Chirag Agarwal, Nari Johnson, Martin Pawelczyk, Satyapriya Krishna, Eshika Saxena, Marinka Zitnik, and Himabindu Lakkaraju. Rethinking stability for attribution-based explanations. *CoRR*, abs/2203.06877, 2022a.
- Chirag Agarwal, Satyapriya Krishna, Eshika Saxena, Martin Pawelczyk, Nari Johnson, Isha Puri, Marinka Zitnik, and Himabindu Lakkaraju. OpenXAI: Towards a transparent evaluation of model explanations. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022b.
- Chirag Agarwal, Sree Harsha Tanneru, and Himabindu Lakkaraju. Faithfulness vs. plausibility: On the (un)reliability of explanations from large language models. *CoRR*, abs/2402.04614, 2024.
- AlignmentResearch. robust-llm-pythia-imdb-14m-mz-ada-v3 (revision da044f3), 2024.
- Tiago A. Almeida, Jose Maria Gomez Hidalgo, and Akebo Yamakami. Contributions to the study of sms spam filtering: New collection and results. In *Proceedings of the 2011 ACM Symposium on Document Engineering (DOCENG’11)*, 2011.
- David Alvarez-Melis and Tommi S. Jaakkola. Towards robust interpretability with self-explaining neural networks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 7786–7795, 2018a.
- David Alvarez-Melis and Tommi S. Jaakkola. Towards robust interpretability with self-explaining neural networks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 7786–7795, 2018b.
- Kenza Amara, Rita Sevastjanova, and Mennatallah El-Assady. Challenges and opportunities in text generation explainability. In Luca Longo, Sebastian Lapuschkin, and Christin Seifert (eds.), *Explainable Artificial Intelligence*, pp. 244–264, Cham, 2024. Springer Nature Switzerland. ISBN 978-3-031-63787-2.

- Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Towards better understanding of gradient-based attribution methods for deep neural networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.
- Christopher J. Anders, David Neumann, Wojciech Samek, Klaus-Robert Müller, and Sebastian Lapuschkin. Software for dataset-wide xai: From local explanations to global insights with Zennit, CoRelAy, and ViRelAy. *CoRR*, abs/2106.13200, 2021.
- Vijay Arya, Rachel K. E. Bellamy, Pin-Yu Chen, Amit Dhurandhar, Michael Hind, Samuel C. Hoffman, Stephanie Houde, Q. Vera Liao, Ronny Luss, Aleksandra Mojsilovic, Sami Mourad, Pablo Pedemonte, Ramya Raghavendra, John T. Richards, Prasanna Sattigeri, Karthikeyan Shanmugam, Moninder Singh, Kush R. Varshney, Dennis Wei, and Yunfeng Zhang. One explanation does not fit all: A toolkit and taxonomy of AI explainability techniques, 2019.
- Pepa Atanasova, Oana-Maria Camburu, Christina Lioma, Thomas Lukasiewicz, Jakob Grue Simonsen, and Isabelle Augenstein. Faithfulness tests for natural language explanations. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pp. 283–294. Association for Computational Linguistics, 2023.
- Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7), 2015.
- David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *J. Mach. Learn. Res.*, 11:1803–1831, 2010.
- Dilyara Bareeva, Marina MC Höhne, Alexander Warnecke, Lukas Pirch, Klaus Robert Muller, Konrad Rieck, and Kirill Bykov. Manipulating feature visualizations with gradient slingshots. In *ICML 2024 Next Generation of AI Safety Workshop*, 2024a.
- Dilyara Bareeva, Galip Ümit Yolcu, Anna Hedström, Niklas Schmolenski, Thomas Wiegand, Wojciech Samek, and Sebastian Lapuschkin. Quanda: An interpretability toolkit for training data attribution evaluation and beyond, 2024b. URL <https://arxiv.org/abs/2410.07158>.
- Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI: <https://doi.org/10.24432/C5XW20>.
- I. Bellido and E. Fiesler. Do backpropagation trained neural networks have normal weight distributions? In Stan Gielen and Bert Kappen (eds.), *ICANN '93*, pp. 772–775, London, 1993. Springer London. ISBN 978-1-4471-2063-6.
- Alessandro Benfenati and Alessio Marta. A singular riemannian geometry approach to deep neural networks i. theoretical foundations. *Neural Networks*, 158:331–343, 2023a. ISSN 0893-6080. doi: 10.1016/j.neunet.2022.11.022.
- Alessandro Benfenati and Alessio Marta. A singular riemannian geometry approach to deep neural networks ii. reconstruction of 1-d equivalence classes. *Neural Networks*, 158:344–358, 2023b. ISSN 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2022.11.026>.
- Jose Manuel Benitez, Juan Luis Castro, and Ignacio Requena. Are artificial neural networks black boxes? *IEEE Trans. Neural Networks*, 8(5):1156–1164, 1997.
- Pietro Berkes and Laurenz Wiskott. On the analysis and interpretation of inhomogeneous quadratic forms as receptive fields. *Neural Comput.*, 18(8):1868–1895, 2006.
- Umang Bhatt, Adrian Weller, and José M. F. Moura. Evaluating and aggregating feature-based model explanations. In Christian Bessiere (ed.), *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pp. 3016–3022. ijcai.org, 2020.

- Robi Bhattacharjee and Ulrike von Luxburg. Auditing local explanations is hard, 2024.
- Monica Bianchini and Franco Scarselli. On the complexity of neural network classifiers: A comparison between shallow and deep architectures. *IEEE Transactions on Neural Networks and Learning Systems*, 25(8):1553–1565, 2014. doi: 10.1109/TNNLS.2013.2293637.
- Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. Language models can explain neurons in language models. *URL <https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html>*. (Date accessed: 14.05.2023), 2, 2023.
- Alexander Binder, Leander Weber, Sebastian Lapuschkin, Grégoire Montavon, Klaus-Robert Müller, and Wojciech Samek. Shortcomings of top-down randomization-based sanity checks for evaluations of deep neural network explanations. *CoRR*, abs/2211.12486, 2022.
- Stefan Blücher, Johanna Vielhaben, and Nils Strodthoff. Decoupling pixel flipping and occlusion strategy for consistent xai benchmarks. *arXiv preprint arXiv:2401.06654*, 2024.
- Douglas G Bonett and Thomas A Wright. Sample size requirements for estimating pearson, kendall and spearman correlations. *Psychometrika*, 65:23–28, 2000.
- Sebastian Bordt and Ulrike von Luxburg. Statistics without interpretation: A sober look at explainable machine learning. *CoRR*, abs/2402.02870, 2024.
- Judy Borowski, Roland Simon Zimmermann, Judith Schepers, Robert Geirhos, Thomas S. A. Wallis, Matthias Bethge, and Wieland Brendel. Exemplary natural images explain CNN activations better than state-of-the-art feature visualization. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.
- Lennart Brocki and Neo Christopher Chung. Evaluation of interpretability methods and perturbation artifacts in deep neural networks. *CoRR*, abs/2203.02928, 2022.
- Lukas Brunke, Prateek Agrawal, and Nikhil George. Evaluating input perturbation methods for interpreting CNNs and saliency map comparison. In *Computer Vision – ECCV 2020 Workshops*, pp. 120–134. Springer International Publishing, 2020.
- Thomas F Burns. Semantically-correlated memories in a dense associative model. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 4936–4970. PMLR, 21–27 Jul 2024.
- Thomas F Burns and Tomoki Fukai. Simplicial hopfield networks. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=_QLsH8gatwx.
- Thomas F Burns and Robert Tang. Detecting danger in gridworlds using Gromov’s Link Condition. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
- Kirill Bykov, Anna Hedström, Shinichi Nakajima, and Marina M.-C. Höhne. Noisegrad - enhancing explanations by introducing stochasticity to model weights. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, February 22 - March 1, 2022*, pp. 6132–6140. AAAI Press, 2022.
- Kirill Bykov, Laura Kopf, Shinichi Nakajima, Marius Kloft, and Marina MC Höhne. Labeling neural representations with inverse recognition. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

- Zachariah Carmichael and Walter Scheirer. How well do feature-additive explainers explain feature-additive predictors? In *XAI in Action: Past, Present, and Future Applications*, 2023.
- Prasad Chalasani, Jiefeng Chen, Amrita Roy Chowdhury, Xi Wu, and Somesh Jha. Concise explanations of neural networks using adversarial training. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1383–1391. PMLR, 2020.
- Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 782–791. Computer Vision Foundation / IEEE, 2021.
- Yu-Neng Chuang, Guanchu Wang, Chia-Yuan Chang, Ruixiang Tang, Fan Yang, Mengnan Du, Xuanting Cai, and Xia Hu. Large language models as faithful explainers. *CoRR*, abs/2402.04678, 2024.
- Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- Ian Covert, Chanwoo Kim, and Su-In Lee. Learning to estimate shapley values with vision transformers. *CoRR*, abs/2206.05282, 2022.
- Sanjoy Dasgupta, Nave Frost, and Michal Moshkovitz. Framework for evaluating faithfulness of local explanations. In *International Conference on Machine Learning*, pp. 4794–4815. PMLR, 2022.
- Ann-Kathrin Dombrowski, Maximilian Alber, Christopher J. Anders, Marcel Ackermann, Klaus-Robert Müller, and Pan Kessel. Explanations can be manipulated and geometry is to blame. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 13567–13578, 2019.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022.
- Dumitru Erhan, Y. Bengio, Aaron Courville, and Pascal Vincent. Visualizing higher-layer features of a deep network. *Technical Report, Université de Montréal*, 01 2009.
- Yanhong Fei, Xian Wei, Yingjie Liu, Zhengyu Li, and Mingsong Chen. A survey of geometric optimization for deep learning: From euclidean space to riemannian manifold. *arXiv preprint arXiv:2302.08210*, 2023.
- Thomas Fel, Thibaut Boissin, Victor Boutin, Agustin Picard, Paul Novello, Julien Colin, Drew Linsley, Tom Rousseau, Rémi Cadène, Laurent Gardes, and Thomas Serre. Unlocking feature visualization for deeper networks with magnitude constrained optimization. *Advances in Neural Information Processing Systems*, 36, 2024.
- Robert Geirhos, Roland S. Zimmermann, Blair L. Bilodeau, Wieland Brendel, and Been Kim. Don’t trust your eyes: on the (un)reliability of feature visualizations. *CoRR*, abs/2306.04719, 2023.
- Marzyeh Ghassemi, Luke Oakden-Rayner, and Andrew L Beam. The false hope of current approaches to explainable artificial intelligence in health care. *The Lancet Digital Health*, 3(11):e745–e750, 2021. ISSN 2589-7500. doi: 10.1016/S2589-7500(21)00208-9.
- Riccardo Guidotti. Evaluating local explanation methods on ground truth. *Artif. Intell.*, 291:103428, 2021. doi: 10.1016/J.ARTINT.2020.103428.

- Peter Hase, Harry Xie, and Mohit Bansal. The out-of-distribution problem in explainability and search methods for feature importance explanations. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 3650–3666, 2021.
- Johannes Haug, Stefan Zürn, Peter El-Jiz, and Gjergji Kasneci. On baselines for local feature attributions. *CoRR*, abs/2101.00905, 2021.
- Michael Hauser and Asok Ray. Principles of riemannian geometry in neural networks. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 770–778. IEEE Computer Society, 2016.
- Anna Hedström, Philine Lou Bommer, Kristoffer Knutsen Wickstrøm, Wojciech Samek, Sebastian Lapuschkin, and Marina M.-C. Höhne. The meta-evaluation problem in explainable AI: identifying reliable estimators with metaquantus. *Trans. Mach. Learn. Res.*, 2023, 2023a.
- Anna Hedström, Leander Weber, Daniel Krakowczyk, Dilyara Bareeva, Franz Motzkus, Wojciech Samek, Sebastian Lapuschkin, and Marina M.-C. Höhne. Quantus: An explainable ai toolkit for responsible evaluation of neural network explanations and beyond. *Journal of Machine Learning Research*, 24(34): 1–11, 2023b.
- Anna Hedström, Leander Weber, Sebastian Lapuschkin, and Marina Höhne. A fresh look at sanity checks for saliency maps. In *Explainable Artificial Intelligence*, pp. 403–420, Cham, 2024. Springer Nature Switzerland.
- Robin Hesse, Simone Schaub-Meyer, and Stefan Roth. Benchmarking the attribution quality of vision models. *CoRR*, abs/2407.11910, 2024. doi: 10.48550/arXiv.2407.11910.
- Jesse Hoogland, George Wang, Matthew Farrugia-Roberts, Liam Carroll, Susan Wei, and Daniel Murfet. The developmental landscape of in-context learning. *arXiv preprint arXiv:2402.02364*, 2024.
- Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability methods in deep neural networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 9734–9745, 2019.
- Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=F76bwRSLeK>.
- Alon Jacovi and Yoav Goldberg. Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness? In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pp. 4198–4205. Association for Computational Linguistics, 2020.
- Neil Jethani, Adriel Saporta, and Rajesh Ranganath. Don’t be fooled: label leakage in explanation methods and the importance of their quantitative evaluation. In Francisco J. R. Ruiz, Jennifer G. Dy, and Jan-Willem van de Meent (eds.), *International Conference on Artificial Intelligence and Statistics, 25-27 April 2023, Palau de Congressos, Valencia, Spain*, volume 206 of *Proceedings of Machine Learning Research*, pp. 8925–8953. PMLR, 2023.
- Niklas Koenen and Marvin N. Wright. Toward understanding the disagreement problem in neural network feature attribution. *CoRR*, abs/2404.11330, 2024.

- Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. Captum: A unified and generic model interpretability library for pytorch. *CoRR*, abs/2009.07896, 2020.
- Narine Kokhlikyan, Vivek Miglani, Bilal Alsallakh, Miguel Martin, and Orion Reblitz-Richardson. Investigating sanity checks for saliency maps with image and text classification. *CoRR*, abs/2106.07475, 2021.
- Laura Kopf, Philine Lou Bommer, Anna Hedström, Sebastian Lapuschkin, Marina MC Höhne, and Kirill Bykov. Cosy: Evaluating textual explanations of neurons. In *ICML 2024 Workshop on Mechanistic Interpretability*, 2024.
- Satyapriya Krishna, Tessa Han, Alex Gu, Javin Pombra, Shahin Jabbari, Steven Wu, and Himabindu Lakkaraju. The disagreement problem in explainable machine learning: A practitioner’s perspective. *CoRR*, abs/2202.01602, 2022.
- Satyapriya Krishna, Jiaqi Ma, Dylan Z Slack, Asma Ghandeharioun, Sameer Singh, and Himabindu Lakkaraju. Post hoc explanations of language models can improve language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Nicholas Kroeger, Dan Ley, Satyapriya Krishna, Chirag Agarwal, and Himabindu Lakkaraju. Are large language models post hoc explainers? *CoRR*, abs/2310.05797, 2023.
- Henry Kvinge, Grayson Jorgenson, Davis Brown, Charles Godfrey, and Tegan Emerson. Internal representations of vision models through the lens of frames on data manifolds. In *NeurIPS 2023 Workshop on Symmetry and Geometry in Neural Representations*, 2023.
- Himabindu Lakkaraju, Dylan Slack, Yuxin Chen, Chenhao Tan, and Sameer Singh. Rethinking explainability as a dialogue: A practitioner’s perspective. *CoRR*, abs/2202.01875, 2022.
- Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamile Lukosiute, Karina Nguyen, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Larson, Sam McCandlish, Sandipan Kundu, Saurav Kadavath, Shannon Yang, Thomas Henighan, Timothy Maxwell, Timothy Telleen-Lawton, Tristan Hume, Zac Hatfield-Dodds, Jared Kaplan, Jan Brauner, Samuel R. Bowman, and Ethan Perez. Measuring faithfulness in chain-of-thought reasoning. *CoRR*, abs/2307.13702, 2023.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, 1998.
- Yann LeCun, Corinna Cortes, and CJ Burges. MNIST handwritten digit database, 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- John M Lee. *Smooth manifolds*. Springer, 2012.
- Tom Lieberum, Senthoran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2, 2024. URL <https://arxiv.org/abs/2408.05147>.
- Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 4765–4774, 2017.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics.

- Aleksandar Makelov, Georg Lange, and Neel Nanda. Towards principled evaluations of sparse autoencoders for interpretability and control. In *ICLR 2024 Workshop on Secure and Trustworthy Large Language Models*, 2024. URL <https://openreview.net/forum?id=MHIX9H8aYF>.
- Alex Mallen and Nora Belrose. Balancing label quantity and quality for scalable elicitation, 2024.
- Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Cristian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, and et al. Gemma: Open models based on gemini research and technology. *CoRR*, abs/2403.08295, 2024.
- Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognit.*, 65:211–222, 2017. doi: 10.1016/j.patcog.2016.11.008.
- Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Methods for interpreting and understanding deep neural networks. *Digit. Signal Process.*, 73:1–15, 2018.
- Niels J. S. Morch, Ulrik Kjems, Lars Kai Hansen, Claus Svarer, Ian Law, Benny Lautrup, Stephen C. Strother, and Kelly Rehm. Visualization of neural networks using saliency maps. In *Proceedings of International Conference on Neural Networks (ICNN’95), Perth, WA, Australia, November 27 - December 1, 1995*, pp. 2085–2090. IEEE, 1995.
- Gregory Naitzat, Andrey Zhitnikov, and Lek-Heng Lim. Topology of deep neural networks. *Journal of Machine Learning Research*, 21(184):1–40, 2020.
- Michael Neely, Stefan F. Schouten, Maurits J. R. Bleeker, and Ana Lucic. Order in the court: Explainable AI methods prone to disagreement. *CoRR*, abs/2105.03287, 2021.
- Anphi Nguyen and Maria Rodriguez Martinez. On quantitative aspects of model interpretability. *CoRR*, abs/2007.07584, 2020.
- Hoa Nguyen. Activation maximization. <https://github.com/Nguyen-Hoa/Activation-Maximization>, 2020.
- Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2(11):e7, 2017.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023.
- Theodore Papamarkou, Tolga Birdal, Michael M. Bronstein, Gunnar E. Carlsson, Justin Curry, Yue Gao, Mustafa Hajj, Roland Kwitt, Pietro Lio, Paolo Di Lorenzo, Vasileios Maroulas, Nina Miolane, Farzana Nasrin, Karthikeyan Natesan Ramamurthy, Bastian Rieck, Simone Scardapane, Michael T Schaub, Petar Veličković, Bei Wang, Yusu Wang, Guowei Wei, and Ghada Zamzmi. Position: Topological deep learning is the new frontier for relational learning. In *Forty-first International Conference on Machine Learning*, 2024.
- ProPublica. Compas recidivism risk score data and analysis, 2016. URL <https://github.com/propublica/compas-analysis>.
- Luyu Qiu, Yi Yang, Caleb Chen Cao, Jing Liu, Yueyuan Zheng, Hilary Hei Ting Ngai, Janet H. Hsiao, and Lei Chen. Resisting out-of-distribution data problem in perturbation of XAI. *CoRR*, abs/2107.14000, 2021.

- Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should I trust you?": Explaining the predictions of any classifier. In Balaji Krishnapuram, Mohak Shah, Alexander J. Smola, Charu C. Aggarwal, Dou Shen, and Rajeev Rastogi (eds.), *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pp. 1135–1144. ACM, 2016a.
- Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should I trust you?": Explaining the predictions of any classifier. In Balaji Krishnapuram, Mohak Shah, Alexander J. Smola, Charu C. Aggarwal, Dou Shen, and Rajeev Rastogi (eds.), *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pp. 1135–1144. ACM, 2016b.
- Laura Rieger and Lars Kai Hansen. IROF: a low resource evaluation metric for explanation methods. *CoRR*, abs/2003.08747, 2020.
- Manuel Romero. bert-tiny-finetuned-sms-spam-detection (revision 012e268), 2024.
- Yao Rong, Tobias Leemann, Vadim Borisov, Gjergji Kasneci, and Enkelejda Kasneci. A consistent and efficient evaluation strategy for attribution methods. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 18770–18795. PMLR, 2022.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vis.*, 115(3):211–252, 2015.
- Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE Trans. Neural Networks Learn. Syst.*, 28(11):2660–2673, 2017.
- Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *Int. J. Comput. Vis.*, 128(2):336–359, 2020.
- Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *CoRR*, abs/1605.01713, 2016.
- Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 3145–3153. PMLR, 2017.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In Yoshua Bengio and Yann LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Workshop Track Proceedings*, 2014.
- Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- Irwin Sobel, Gary Feldman, et al. A 3x3 isotropic gradient operator for image processing. *a talk at the Stanford Artificial Project in*, 1968:271–272, 1968.

- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp. 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.
- Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin A. Riedmiller. Striving for simplicity: The all convolutional net. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Workshop Track Proceedings*, 2015.
- Claudio Stefano, Francesco Fontanella, Marilena Maniaci, and Alessandra Freca. Avila. UCI Machine Learning Repository, 2018. DOI: 10.24432/C5K02X.
- Cory Stephenson, suchismita padhy, Abhinav Ganesh, Yue Hui, Hanlin Tang, and SueYeon Chung. On the geometry of generalization and memorization in deep neural networks. In *International Conference on Learning Representations*, 2021.
- Pascal Sturmfels, Scott Lundberg, and Su-In Lee. Visualizing the impact of feature attribution baselines. *Distill*, 2020. doi: 10.23915/distill.00022. <https://distill.pub/2020/attribution-baselines>.
- Mukund Sundararajan and Ankur Taly. A note about: Local explanation methods for deep neural networks lack sensitivity to parameter values. *CoRR*, abs/1806.04205, 2018.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 3319–3328. PMLR, 2017.
- Zeren Tan and Yang Tian. Robust explanation for free or at the cost of faithfulness. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 33534–33562. PMLR, 2023.
- Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- VityaVitalich. bert-tiny-sst2 (revision 2e14b76), 2023.
- Matthew Watson, Bashar Awwad Shiekh Hasan, and Noura Al Moubayed. Agree to disagree: When deep learning models with identical architectures produce distinct explanations. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 875–884, January 2022.
- Kristoffer K Wickstrøm, Marina M. C. Höhne, and Anna Hedström. From flexibility to manipulation: The slippery slope of xai evaluation. *CoRR*, abs/INSERT, 2024.
- Frank Wilcoxon. Individual comparisons by ranking methods. *Biometrics Bulletin*, 1(6):80–83, 1945. ISSN 00994987.

- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. Association for Computational Linguistics.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing Ni. Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image classification. *Scientific Data*, 10(1):41, 2023.
- Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Sai Suggala, David I. Inouye, and Pradeep Ravikumar. On the (in)fidelity and sensitivity of explanations. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 10965–10976, 2019.
- Gal Yona and Daniel Greenfeld. Revisiting sanity checks for saliency maps. *CoRR*, abs/2110.14297, 2021.
- Jason Yosinski, Jeff Clune, Anh Mai Nguyen, Thomas J. Fuchs, and Hod Lipson. Understanding neural networks through deep visualization. *CoRR*, abs/1506.06579, 2015.
- Jerrold H Zar. Spearman rank correlation. *Encyclopedia of biostatistics*, 7, 2005.
- Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In David J. Fleet, Tomás Pajdla, Bernt Schiele, and Tinne Tuytelaars (eds.), *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I*, volume 8689 of *Lecture Notes in Computer Science*, pp. 818–833. Springer, 2014a.
- Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In David J. Fleet, Tomás Pajdla, Bernt Schiele, and Tinne Tuytelaars (eds.), *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I*, volume 8689 of *Lecture Notes in Computer Science*, pp. 818–833. Springer, 2014b.
- Roland S. Zimmermann, Judy Borowski, Robert Geirhos, Matthias Bethge, Thomas S. A. Wallis, and Wieland Brendel. How well do feature visualizations support causal understanding of CNN activations? In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 11730–11744, 2021.

Appendix

The Appendix is organised as follows theoretical considerations (Section A.1), implementation notes for **GEF** and **Fast-GEF** (Section A.2), details on the **MetaQuantus** framework (Section A.3), the general experimental setup (Section A.4), analysis of model assumptions (Section A.5), alignment patterns extended results (Section A.6), ablation experimental results (Section A.7), extended results from individual experiments (Section A.8), and notation tables (Section A.9).

A.1 Theoretical Considerations

The following subsections provide detailed proofs, extensions, and discussions surrounding the GEF criterion.

A.1.1 GEF: Penalising Random Explanations

Following the discussion in Section 3, a quality estimator should be able to identify unfaithful explanations. In the following, we show that our proposed GEF criterion (Definition 5) recognise two specific types of unfaithful explanations: constant, and random explanations, which is independent of its model, by design.

Corollary 1 (Penalising Unfaithful Explanations) *Let Ψ^{GEF} be a quality estimator that yields estimates $q^{GEF} \in \mathbb{R}$, with $q^{GEF} = 0$ indicating a lack of generalised faithfulness. To be a valid measure of explanation quality, Ψ^{GEF} should assign low scores to both (I) constant, and (II) random explanations*

$$\text{Constant (I): } \forall \hat{e} : \mathbf{e} = \hat{e} \Rightarrow q^{GEF} = 0$$

$$\text{Random (II): } \forall \hat{e} : \hat{e} \sim \mathcal{U}(0, 1) \Rightarrow q^{GEF} = 0$$

where \hat{e} , and \mathbf{e} are perturbed, and unperturbed explanations, and $\mathcal{U}(0, 1)$ denotes a uniform distribution. The GEF estimate $q^{GEF} = \rho(\mathbf{d}_f, \mathbf{d}_\phi)$ (Definition 5) assigns low scores in the first, and the second case.

Proof. In case (I), the explanation does not change across perturbations, leading to an explanation distortion vector \mathbf{d}_ϕ that contains only zeros

$$\forall \hat{e}, z \in [1, Z] : \mathbf{e} = \hat{e} \Rightarrow \mathbf{D}_\phi^z = 0,$$

whereas the model’s distortion vector \mathbf{d}_f will contain non-zero values due to perturbations

$$\forall z \in [1, Z] : \mathbf{D}_f^z \neq 0.$$

Consequently, the correlation coefficient $\rho(\mathbf{d}_f, \mathbf{d}_\phi)$ will be zero with $q^{GEF} = 0$.

In case (II), the explanation distortion \mathbf{D}_ϕ^z will be approximately uniform across all perturbation steps since each perturbation is independently drawn from the same distribution:

$$\forall \hat{e}, z, j \in [1, Z], z \neq j : \hat{e} \sim \mathcal{U}(0, 1) \Rightarrow \mathbf{D}_\phi^z \approx \mathbf{D}_\phi^j,$$

whereas the model distortion \mathbf{D}_f^z will vary according to the degree of the perturbation

$$\forall \hat{e}, z, j \in [1, Z], z \geq j : \hat{e} \sim \mathcal{U}(0, 1) \Rightarrow \mathbf{D}_f^j \geq \mathbf{D}_f^z.$$

The lack of correlation between \mathbf{d}_f , and \mathbf{d}_ϕ results in a quality measure q^{GEF} that is equal to zero. This completes the proof.

A.1.2 GEF: Extension

To extend the applicability of GEF (Definition 5) to global methods that explain *any* neuron within a model, we adopt Kopf et al. (2024), and view the model f as a composition of two functions, $F : \mathcal{X} \rightarrow \mathcal{G}$, and $L : \mathcal{G} \rightarrow \mathcal{Y}$, such that $f = L \circ F$. Here $\mathcal{G} \subset \mathbb{R}^{c \times w^* \times h^*}$, where $c \in \mathbb{N}$ is the number of neurons in the layer, and $w^*, h^* \in \mathbb{N}$ represent the width, and height of the feature map, respectively. The function F , is referred to as the *feature extractor*. We redefine the model function as a chosen feature extractor, and replace y in Definition 5 with the activation of the c^{th} neuron such that *i.e.*, $y = F_c(\mathbf{x}, \theta) : \mathcal{X} \rightarrow \mathbb{R}^{w^* \times h^*}$. While the model’s output space \mathcal{Y} is replaced by \mathcal{G} , we similarly define the perturbed instance \hat{y} .

A.1.3 GEF: Derivation of Linear Case

Our definition of GEF is based on the observation that any distortion present in the model output space \mathcal{Y} , should be mirrored in the explanation space \mathcal{E} . Since neural networks are non-linear functions, a fair distortion in \mathcal{Y} , and \mathcal{E} , requires the introduction of the pullback (Section 4).

In the case of a linear model, however, the relationship between the distortion quantities \mathbf{D}_f , and \mathbf{D}_ϕ can be derived analytically. Here, the explanation is based on the first-order Taylor term, which is a linear approximation of the model’s behaviour, forming the foundation of many established explanation methods (*e.g.*, Montavon et al. (2017)). We proceed to derive this relationship explicitly below.

Proof. Consider f to be a linear model of the form $f(\mathbf{x}; \theta) = \theta\mathbf{x} + c$. The explanation is the parameter vector θ . We can derive the expected distortion $\mathbf{D}_f := \mathbb{E}_{\hat{\theta}_m} [(f(\mathbf{x}; \theta) - f(\mathbf{x}; \hat{\theta}_m))^2]$ (see Equation 6) where $m \in [1, M]$ denotes the number of perturbed models for a fixed perturbation magnitude ξ , *i.e.*, a step z .

$$\begin{aligned} \mathbf{D}_f^z &= (\theta\mathbf{x} + c)^2 - 2(\theta\mathbf{x} + c)\mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m\mathbf{x} + c] + \mathbb{E}_{\hat{\theta}_m} [(\hat{\theta}_m\mathbf{x} + c)^2] \\ \mathbf{D}_f^z &= \theta^2\mathbf{x}^2 + 2c\theta\mathbf{x} + 2c^2 - 2\mathbf{x}(\theta\mathbf{x} + c)\mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m] - 2c(\theta\mathbf{x} + c) + \mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m^2\mathbf{x}^2 + 2c\hat{\theta}_m\mathbf{x}] \\ \mathbf{D}_f^z &= \theta^2\mathbf{x}^2 - 2\mathbf{x}(\theta\mathbf{x} + c)\mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m] + 2c\mathbf{x}\mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m] + \mathbf{x}^2\mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m^2], \\ \mathbf{D}_f^z &= \theta^2\mathbf{x}^2 - 2\theta\mathbf{x}^2\mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m] + \mathbf{x}^2\mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m^2]. \end{aligned} \quad (14)$$

For the explanation distortion, a similar decomposition can be performed

$$\mathbf{D}_\phi^z = \theta^2 - 2\theta\mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m] + \mathbb{E}_{\hat{\theta}_m} [\hat{\theta}_m^2]. \quad (15)$$

By combining Equation 14 and 15, we arrive at

$$\mathbf{D}_\phi^z = \frac{1}{\mathbf{x}^2} \mathbf{D}_f^z, \quad (16)$$

We can construct the distortion vectors \mathbf{d}_ϕ , and \mathbf{d}_f , and for each entry Equation 16 holds. When ρ is defined as the Pearson correlation coefficient, we find the distortion of the model \mathbf{d}_f , and the distortion of the explanation function \mathbf{d}_ϕ to be perfectly correlated

$$\rho(\mathbf{d}_\phi, \mathbf{d}_f) = \frac{\text{cov}_\xi(\mathbf{d}_f, \mathbf{d}_\phi)}{\sqrt{\text{Var}_\xi(\mathbf{d}_f)}\sqrt{\text{Var}_\xi(\mathbf{d}_\phi)}},$$

which is equal to

$$\rho(\mathbf{d}_\phi, \mathbf{d}_f) = \frac{1/\mathbf{x}^2\text{Var}_\xi(\mathbf{d}_f)}{1/\mathbf{x}^2\text{Var}_\xi(\mathbf{d}_f)} = 1. \quad (17)$$

This proves that in a simplified scenario, the key assumption of correlated distortion quantities holds, *i.e.*, the model parameters θ provide a *perfectly faithful* explanation. Since monotonicity is a weaker condition than linearity, Equation 17 also holds when ρ is defined as the *Spearman Rank correlation* coefficient.

A.1.4 GEF: Influence of Z

The parameter Z represents the number of steps in the perturbation path, and consequently dictates how finely the model’s response will be captured by the GEF criterion (Definition 5). As such, selecting an appropriate value for Z is critical because it affects the interpretation of the results. A higher Z allows for a finer evaluation of how well an explanation aligns with the model’s behaviour under varying conditions. When using Spearman’s rank correlation coefficient as our measure of ρ , a larger Z generally stabilises the faithfulness score due to the reduction in confidence intervals with more samples (*i.e.*, $CI \sim \frac{1}{Z}$) (Bonett &

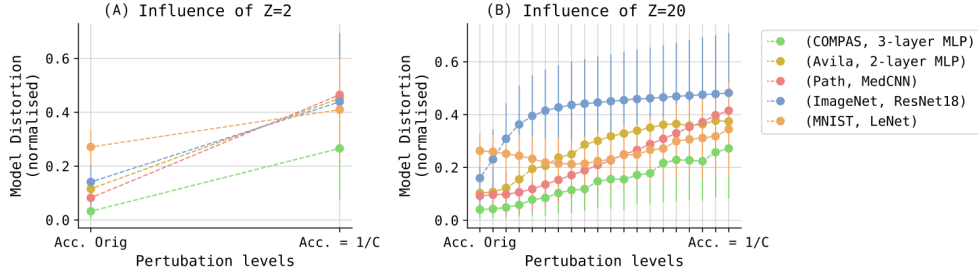


Figure A.1: Model distortion (normalised by its maximum value) with (A) showing $Z = 2$ perturbation steps, and (B) showing $Z = 20$ perturbation steps.

Wright, 2000). Nonetheless, this assumes a monotonic response from both the model, and the explanation, which may not be realistic (Section 4). If the model itself is not monotonic across perturbations, expecting the explanation to behave monotonically is also unrealistic.

Figure A.1 (A) ($Z = 2$), and (B) ($Z = 20$) demonstrate the violation of the monotonicity assumption, as we observe large error bars, and divergent behaviour for $Z = 20$, indicating non-monotonic responses. Accordingly, a moderate value of Z (Zar, 2005) is advised for meaningful measurement.

A.2 Notes on GEF and Fast-GEF Implementation

In the following, we provide details on the GEF algorithm.

A.2.1 Generate Perturbation Path

To generate the perturbation path of length Z , satisfying $y \neq \hat{y}$, $\forall y, \hat{y} \in \mathcal{Y}$, we computationally find the minimum noise level σ_z^2 at $z = Z$, such that the perturbed model’s accuracy (ACC) approximates $\frac{1}{C}$, where C is the number of classes, within a threshold, *i.e.*, $\epsilon \ll 1$. Here, $\text{ACC} = \frac{1}{N} \sum_{i=1}^N (f(\mathbf{X}_i; \theta) = \mathbf{Y}_i)$ where N is the number of samples in the test set, denoted \mathbf{X} . This is achieved by progressively increasing σ^2 , and applying it to the model according to Section 5.1, which process concludes when the model’s accuracy satisfies the condition $|\text{ACC} - \frac{1}{C}| < \epsilon$, thereby determining perturbation level for subsequent evaluation.

Compute Path Length. For a more faithful estimate of explanation distortion, for each step $z \in [0, Z]$, we compute \mathbf{D}_ϕ , defined as the path length $L(\gamma)$. We replace the integral in Equation 13 with a sum over T steps:

$$L(\gamma) = \sum_{t=1}^T de_t^T (J_f(\hat{\mathbf{e}}_t)^T J_f(\hat{\mathbf{e}}_t)) de_t, \quad (18)$$

where $de_t \in \mathbb{R}^V$ denotes the feature-wise difference in explanations *i.e.*, $(\mathbf{e} - \hat{\mathbf{e}}_t)$ with $\phi(f_{\hat{\theta}_t}, \dots) = \hat{\mathbf{e}}_t$, and $J_f(\hat{\mathbf{e}}_t) \in \mathbb{R}^{V \times C}$ is the Jacobian for fixed \mathbf{x} , and $f_{\hat{\theta}_t}$. To numerically approximate this Jacobian, for each step $t \in [0, T]$, we perturb the neural activations (*i.e.*, logits $\hat{\mathbf{y}}$) by adding infinitesimal noise. In practice, we sample from a Gaussian distribution $v_k \sim \mathcal{N}(0, 0.001)$ such that $\hat{\mathbf{y}}_k = \hat{\mathbf{y}} + v_k$, $k \in [1, K]$ times. After each perturbation, we recalculate the corresponding explanation $\phi(\hat{\mathbf{y}}_k, \dots) = \hat{\mathbf{e}}_k$. Elements of the Jacobian $J_f(\hat{\mathbf{e}}_t)$ are then computed as feature-wise difference between \mathbf{e} , and $\hat{\mathbf{e}}_k$:

$$\frac{\partial e_i}{\partial f_j} \approx \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K (e_j - \hat{e}_{j,k}) v_k^{-1}. \quad (19)$$

where i, j refers to the indices of the Jacobian $J_f(\hat{\mathbf{e}}_t)$.

Unless specified otherwise, we set M , Z , T , and K to 5 in all experiments. Please find Appendix A.7 for an ablation study motivating these hyperparameters.

A.3 Notes on MetaQuantus framework

For meta-evaluation, a two-step process is employed. First, two types of controlled perturbations are introduced: minor, and disruptive. These are designed to evaluate the metric’s resilience to noise (NR), and its sensitivity to adversarial conditions (AR), respectively. Specifically, these perturbations are applied in both the input, and model spaces, resulting in two distinct tests: the Input Perturbation Test (IPT), and the Model Perturbation Test (MPT)⁴. Second, the effects of the perturbations are measured in two meta-evaluative criteria: intra-consistency (**IAC**), and inter-consistency (**IEC**). Here, **IAC** refers to measuring the similarity in score distributions post-perturbation, and **IEC** refers to the occurrence of categorical ranking changes within a set of distinct explanation methods⁵. Each metric is then assigned a summarised meta-consistency score, denoted as $MC \in [0, 1]$:

$$MC = \left(\frac{1}{|\mathbf{m}^*|} \right) \mathbf{m}^{*T} \mathbf{m} \quad \text{where} \quad \mathbf{m} = \begin{bmatrix} \mathbf{IAC}_{NR} \\ \mathbf{IAC}_{AR} \\ \mathbf{IEC}_{NR} \\ \mathbf{IEC}_{AR} \end{bmatrix}, \quad (20)$$

with $\mathbf{m}^* \in \mathbb{R}^4$ representing an ideal quality estimator, essentially a vector of ones. A higher MC score, approaching 1, indicates superior reliability according to the defined evaluation criteria. Metrics that demonstrate both resilience to minor perturbations, and reactivity to disruptive changes achieve higher MC scores. We refer to the original publication (Hedström et al., 2023a) for further details on the elements in the meta-evaluation vector \mathbf{m} (Equation 20), and the framework in general.

A.4 General Experimental Setup

Here, we describe the models, datasets, tooling, hardware, explanation, and evaluation methods in this work.

A.4.1 Models, and Datasets

We employ various models for vision, text, and tabular tasks in our experiments. See Table 2.

- For vision classification, we use ImageNet-1K for object recognition (Russakovsky et al., 2015) with ResNet18 (He et al., 2016); Pathology, and Derma for medical image analysis with proposed Med-CNN architecture (Yang et al., 2023); and MNIST (LeCun et al., 2010), and fMNIST, (Xiao et al., 2017) for digit, and fashion recognition with LeNet (LeCun et al., 1998).
- For text classification, we use SMS Spam (Almeida et al., 2011) with a tiny, fine-tuned BERT model (Romero, 2024); IMDb (Maas et al., 2011) with Pythia (AlignmentResearch, 2024); and SST-2 (Socher et al., 2013) with a tiny, fine-tuned BERT model (VityaVitalich, 2023).
- For tabular classification, we use Adult (Becker & Kohavi, 1996) and, COMPAS (ProPublica, 2016), with 3-layer MLP; and Avila (Stefano et al., 2018) with 2-layer MLP.

All models that are not publicly accessible are released at GitHub repository at <https://github.com/annahedstroem/GEF>.

A.4.2 Tooling

Several libraries, and open-source implementations enabled this work, including **transformers** (Wolf et al., 2020), **OpenXAI** (Agarwal et al., 2022b), **Captum** (Kokhlikyan et al., 2020), **Zennit** (Anders et al., 2021),

⁴For the IPT, independent, and identically distributed (i.i.d.) additive uniform noise is applied, defined as $\hat{x}_i = x + \nu_i$, where $\nu_i \sim \mathcal{U}(\alpha, \beta)$. For the MPT, multiplicative Gaussian noise is applied to all network weights, represented as $\hat{\theta}_i = \theta \cdot \nu_i$ with $\nu_i \sim \mathcal{N}(\mu, \sigma^2)$. The hyperparameters $\alpha, \beta, \mu, \sigma^2$ follow the specifications of the original study (Hedström et al., 2023a).

⁵**IAC** provides a normalised p-value derived from the non-parametric *Wilcoxon signed-rank test* (Wilcoxon, 1945), comparing the original, and perturbed score distributions. For NR , similar distributions are expected, whereas for AR , the distributions are anticipated to differ. **IEC** counts ranking changes within explanation methods post-perturbation, with an ideal metric showing consistent rankings under minor noise (NR), and altered rankings under disruptive noise (AR).

Shap (Lundberg & Lee, 2017), Activation-Maximization (Nguyen, 2020), and Horama (Fel et al., 2024). For metric implementation, and meta-evaluation, we use the Quantus (Hedström et al., 2023b), and MetaQuantus (Hedström et al., 2023a) libraries, respectively.

A.4.3 Hardware

The experiments were conducted using two hardware configurations: a cluster with four Tesla V100S-PCIE-32GB GPUs, each offering 32 GB of memory, and a DGX-2 system featuring eight NVIDIA A100-SXM4-40GB GPUs, each with 40 GB of memory. Both setups support the NVIDIA driver version 535.161.07, and CUDA 12.2.

A.4.4 Explanation Methods

All the hyperparameters of the individual explanations methods, are listed in the main manuscript. Concerning the preprocessing, the signs of the attributions are maintained, unless the method algorithmically relies on it such as SAL. Note, that not every explanation method is suitable or intended to be used for all data modalities, and/ or model architectures. For example, GradCAM explanations are primarily designed for convolutional neural networks (CNN) models, and global feature visualisation methods are generally applied to vision tasks. We only report GEF and Fast-GEF results where appropriate.

Normalisation. We perform normalisation using the square root of the mean of the squared values (as detailed in the Appendix of (Binder et al., 2022)). This approach introduces less variance compared to normalisation techniques like scaling by the maximum value. It is defined as follows

$$\text{norm}(\mathbf{e}) = \frac{\mathbf{e}_{h,w}}{\left(\frac{1}{HW} \sum_{h',w'} \mathbf{e}_{h',w'}^2\right)^{1/2}}, \quad (21)$$

where H , and W represent the height, and width, respectively, and $\hat{\mathbf{e}}_{h,w}$ denotes the explanation value at the pixel location (h, w) ⁶.

LLM-x Methodology. In the following, we describe the methodology used to produce LLM-x explanations. An illustration is provided in Figure A.2.

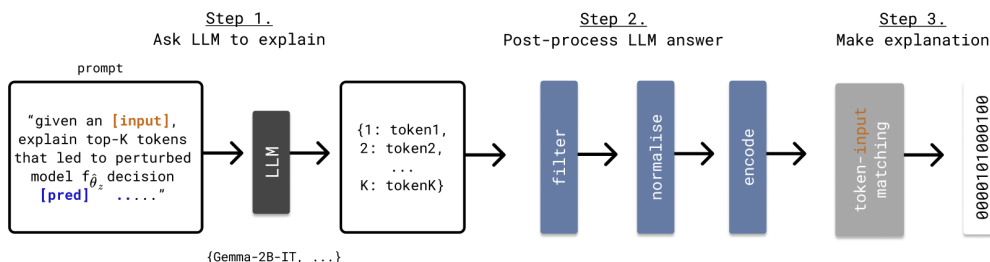


Figure A.2: A high-level overview of the three-step LLM-x methodology.

To generate LLM-x explanations, we use Gemma-2B-IT (Mesnard et al., 2024) as the explainer. For each instance, we create a prompt describing the task, softmax confidence before, and after perturbation, and the class labels. The prompt template introduces the task (*e.g.*, classifying **sms messages** or **sentiment analysis**), and uses synonyms for model descriptions (*e.g.*, "AI", "machine learning"), and perturbation types (*e.g.*, "adversarially manipulated", "perturbed with noise") to vary language. The softmax change is calculated, and added to the template and is described in the context of the model getting "more"

⁶This normalisation method ensures that the mean squared distance from zero of each explanation score equals one. Unlike other normalisation techniques that constrain attribution values to a predefined range—making them suitable for visualisation—this method retains a metric useful for comparing the distances across different explanation methods.

or "less" certain of a class label. The LLM is asked to return the top- K important tokens in a structured JSON format, ranking tokens from 1 to K . The temperature is set to 0 for deterministic outputs.

After prompting, invalid or non-JSON outputs are removed. The LLM-ranked tokens are normalised by lowercasing, removing punctuation. Then, these tokens (or words) are encoded with the original model’s tokenizer. Binary explanation vectors are created by matching the LLM-ranked tokens to the original input tokens, with a value of 1 for matching tokens, and 0 otherwise.

For full details, including the code, and prompt template, we refer to our GitHub repository at <https://github.com/annahedstroem/GEF>.

SAE Methodology. SAEs are designed to create sparse, interpretable representations of the internal activations of a LLM, while preserving their reconstruction. The activations of a given layer l , *i.e.*, $f_\theta(\mathbf{x})$ are encoded into a sparse latent vector $z(\mathbf{x})$, with latent dimensions larger than the internal representation,. Then, these decoded representations are reconstructed such that $\mathbf{g}(z) \approx f_\theta(\mathbf{x})$. This process is defined by the encoder and decoder functions

$$\mathbf{z}(f_\theta(\mathbf{x})) := \sigma(\mathbf{W}_{enc}f_\theta(\mathbf{x}) + \mathbf{b}_{enc}), \quad (22)$$

$$\mathbf{g}(z) := \mathbf{W}_{dec}z + \mathbf{b}_{dec}, \quad (23)$$

where σ enforces sparsity through activation functions like ReLU or JumpReLU (Lieberum et al., 2024), using $L1$ or $L0$ regularisation during training. The SAE explanations are generated by performing a forward pass through the SAE encodings, and storing the activated values of z .

A.4.5 Evaluation Methods

Next, we mathematically define the evaluation methods (or “metrics”) used in this work (Section 6.1).

Faithfulness. Within the faithfulness category, we evaluate three metrics, including, *Faithfulness Correlation* (FC) (Bhatt et al., 2020), *Pixel-Flipping* (PF) (Bach et al., 2015), and *Region-Perturbation* (RP) (Samek et al., 2017). FC is defined as follows

$$\Psi_{\text{FC}} = \underset{S \in |S| \subseteq d}{\text{corr}} \left(\sum_{i \in S} \phi(\mathbf{x}, f, \hat{y}; \lambda)_i, f(\mathbf{x}) - f(\mathbf{x}_{[x_s = \bar{x}_s]}) \right), \quad (24)$$

where $|S| \subseteq D$ is a subset of indices of a sample \mathbf{x} , \bar{x} is the chosen baseline value, and $\mathbf{x}_{[x_s = \bar{x}_s]}$ are the masked input, with randomly chosen indices.

PF returns a vector of prediction scores p_i corresponding to pixel replacements $i \in n$, which are sorted in descending order by the highest relevant pixel in the explanation $\phi(\mathbf{x}, f, \hat{y}; \lambda)$. To return one evaluation score per input sample, we calculate the area under the curve (AUC) as follows

$$\Psi_{\text{PF}} = \sum_{i=1}^n (\hat{y}_i + \hat{y}_{i+1}) \cdot \frac{p_{i+1} - p_i}{2} \quad (25)$$

where p_i , and p_{i+1} are the prediction values of the i^{th} , and $(i+1)^{\text{th}}$ perturbation step, and \hat{y}_i , and \hat{y}_{i+1} the corresponding network prediction.

RP follows the most-relevant-first perturbation strategy, creating consecutive perturbed samples \hat{y}_i, \hat{y}_{i+1} such that for \hat{y}_i perturbed pixels correspond to larger respective explanation values than the pixel perturbed in \hat{y}_{i+1} . Across each perturbation curve, the area over the curve (AOC) is calculated, and averaged across multiple masked inputs $\hat{\mathbf{x}}$ as follows

$$\Psi_{\text{RP}} = \frac{1}{L+1} \mathbb{E}_{(\hat{\mathbf{x}})} \left(\sum_{k=1}^L (\hat{y}_0 + \hat{y}_k) \right), \quad (26)$$

where L is the number of perturbed features in the input.

Robustness. Within the robustness category, we evaluate three metrics, including, *Relative Input Stability* (RIS), *Relative Representation Stability* (RRS), *Relative Output Stability* (ROS) (Agarwal et al., 2022a). RIS extends (Alvarez-Melis & Jaakkola, 2018b), which is a measure of how much the explanation changes *w.r.t.* the input under slight perturbation $\hat{\mathbf{x}} = \mathbf{x} + \mathbf{u}_i$. The change is measured as the l_p norm, and the RIS metric only considers perturbations that result in the same model prediction, *i.e.*, $f(\mathbf{x}) = f(\hat{\mathbf{x}})$. It is defined as follows

$$\Psi_{\text{RIS}} = \max_{\hat{\mathbf{x}}} \frac{\left\| \frac{\phi(\mathbf{x}, f, \hat{y}; \lambda) - \phi(\hat{\mathbf{x}}, f, \hat{y}; \lambda)}{\phi(\mathbf{x}, f, \hat{y}; \lambda)} \right\|_p}{\max\left(\left\| \frac{\mathbf{x} - \hat{\mathbf{x}}}{\mathbf{x}} \right\|_p, \epsilon_{\min}\right)}, \quad \forall \hat{\mathbf{x}} \in \mathcal{N}_\epsilon; f(\mathbf{x}) = f(\hat{\mathbf{x}}) \quad (27)$$

where $\epsilon_{\min} > 0$ ensures a non-zero denominator.

In contrast to the RIS metric, RRS considers the internal representation of the model $\mathcal{L}(\cdot)$ (*e.g.*, an output embedding), while maintaining similar perturbation conditions

$$\Psi_{\text{RRS}} = \max_{\hat{\mathbf{x}}} \frac{\left\| \frac{\phi(\mathbf{x}, f, \hat{y}; \lambda) - \phi(\hat{\mathbf{x}}, f, \hat{y}; \lambda)}{\phi(\mathbf{x}, f, \hat{y}; \lambda)} \right\|_p}{\max\left(\left\| \frac{\mathcal{L}_{\mathbf{x}} - \mathcal{L}_{\hat{\mathbf{x}}}}{\mathcal{L}_{\mathbf{x}}} \right\|_p, \epsilon_{\min}\right)}, \quad \forall \hat{\mathbf{x}} \in \mathcal{N}_\epsilon; f(\mathbf{x}) = f(\hat{\mathbf{x}}) \quad (28)$$

where $\epsilon_{\min} > 0$ ensures a non-zero denominator.

ROS makes similar adaptations as the RRS metric, assumes however that the model’s internal representations are not accessible. Instead the output logits $h(\mathbf{x})$, and $h(\hat{\mathbf{x}})$ are assessed

$$\Psi_{\text{ROS}} = \max_{\hat{\mathbf{x}}} \frac{\left\| \frac{\phi(\mathbf{x}, f, \hat{y}; \lambda) - \phi(\hat{\mathbf{x}}, f, \hat{y}; \lambda)}{\phi(\mathbf{x}, f, \hat{y}; \lambda)} \right\|_p}{\max\left(\|h(\mathbf{x}) - h(\hat{\mathbf{x}})\|_p, \epsilon_{\min}\right)}, \quad \forall \hat{\mathbf{x}} \in \mathcal{N}_\epsilon; f(\mathbf{x}) = f(\hat{\mathbf{x}}) \quad (29)$$

where $\epsilon_{\min} > 0$ ensures a non-zero denominator.

Sensitivity. Within the sensitivity category, we evaluate three metrics, including, *Model Parameter Randomisation Test* (MPRT) (Adebayo et al., 2018), *Smooth Model Parameter Randomisation Test* (sMPRT), *Efficient Model Parameter Randomisation* (eMPRT) (Hedström et al., 2024). MPRT measures the similarity between the original explanation \mathbf{e}_l , and the explanation $\hat{\mathbf{e}} := \phi(\mathbf{x}, \hat{f}_l^t, y)$ of the perturbed model \hat{f}_l^t randomised in a top-down fashion up to layer $l \in [L, L-1, \dots, 1]$

$$\hat{q}^{\text{MPRT}} = \rho(\mathbf{e}, \hat{\mathbf{e}}_l), \quad (30)$$

with similarity function $\rho : \mathbb{R}^D \times \mathbb{R}^D \mapsto \mathbb{R}$.

sMPRT computes a quality estimate $\hat{q} \in \mathbb{R}$ between explanations $\mathbf{e}_i := \phi(\hat{\mathbf{x}}_i, f, y; \lambda)$, and $\hat{\mathbf{e}}_{l,i} := \phi(\hat{\mathbf{x}}_i, \hat{f}_l^b, y; \lambda)$ averaged over $i \in [1, N]$ where $\hat{\mathbf{e}}_{l,i}$ corresponds to the perturbed model \hat{f}_l^b randomised in a bottom-down fashion up to layer $l \in [1, 2, \dots, L]$

$$\hat{q}^{\text{sMPRT}} = \rho\left(\frac{1}{N} \sum_{i=1}^N \mathbf{e}_i, \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{e}}_{l,i}\right), \quad (31)$$

with $\hat{\mathbf{x}}_i = \mathbf{x} + \eta_i$, and $\eta_i \sim \mathcal{N}(0, \sigma)$ with $\|\eta_i\|_p \leq \epsilon$ holding with high probability, for $\sigma, \epsilon \in \mathbb{R}$.

eMPRT measures the relative rise in the complexity of the explanation from a fully randomised model \hat{f} such that $\hat{\mathbf{e}} := \phi(\mathbf{x}, \hat{f}, y; \lambda)$:

$$\hat{q}^{\text{eMPRT}} = \frac{c(\hat{\mathbf{e}}) - c(\mathbf{e})}{c(\mathbf{e})} \quad (32)$$

where $c : \mathbb{R}^D \mapsto \mathbb{R}$ is a complexity function, *e.g.*, discrete entropy.

A.5 Analysing Violations of Model Assumptions

To understand whether perturbation techniques commonly employed for robustness, sensitivity, and faithfulness evaluations generally fulfill the critical assumptions of model distortion (Assumptions 1-3), we performed several experiments. To investigate how often model robustness, sensitivity, and faithfulness (Assumptions 1-3) hold versus fail in practice, we set up a simple experiment that tracks model, and explanation distortions, *i.e.*, \mathbf{D}_f , and \mathbf{D}_ϕ , while applying perturbation commonly used in evaluation such as additive Gaussian noise for robustness evaluation, top-down, and bottom-up layer-by-layer parameter randomisation for sensitivity evaluation, and cumulative masking for faithfulness evaluation.

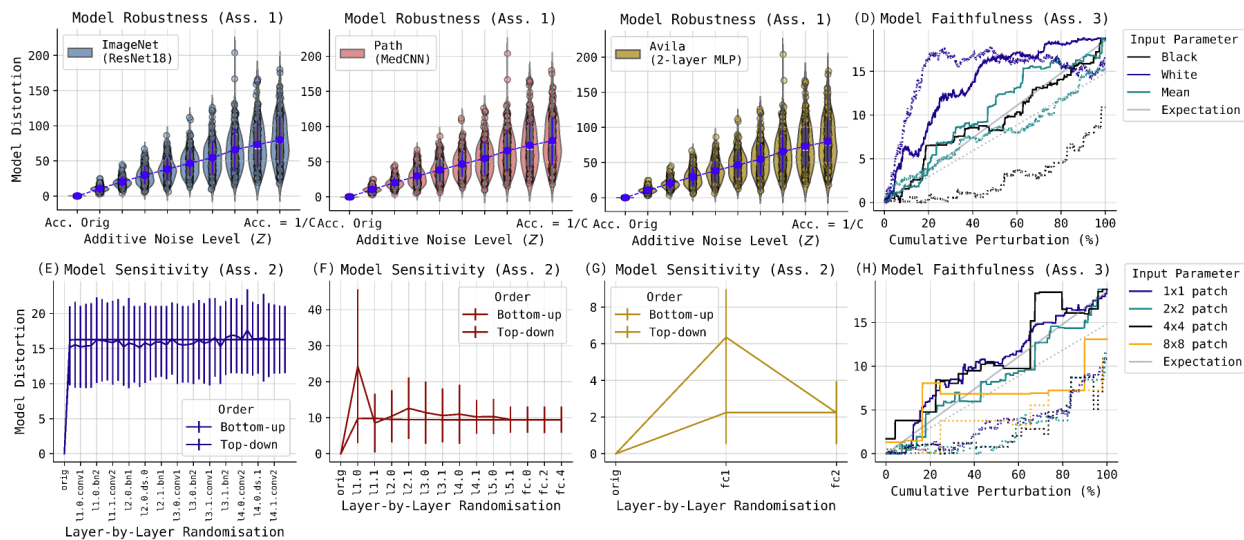


Figure A.3: Impact of model distortion (y-axis) over common perturbation types in robustness, sensitivity, and faithfulness evaluations, across different datasets, and NN architectures. (A), (B), and (C) depict the distribution of model distortions across different perturbation magnitudes of additive Gaussian noise for ImageNet (ResNet18), Path (MedCNN), and Avila (2-layer MLP), respectively. (D), (E), and (F) show the average, and standard deviation of model distortions over layer-wise top-down, and bottom-up randomisation for the same datasets (as indicated by colour). (G), and (H) display model distortions for randomly chosen MNIST (*solid* line), and fMNIST (*dashed* line) samples (LeNet) under cumulative perturbations using different patch sizes (1×1 , 2×2 , 4×4 , 8×8), and baseline replacement strategies (*black*, *white*, *mean*).

Model Robustness Under Additive Noise. To understand the extent to which model robustness (Assumption 1) is generally satisfied for robustness evaluation (Definition 2), we examine Figure A.3 (A), (B), and (C), and Figure A.4 (A) and (B). Here, the distribution of \mathbf{D}_f is visualised over $Z = 10$ input perturbation steps, showing how model distortion varies with increasing input perturbation magnitude, using additive Gaussian noise, *i.e.*, $\nu_i \sim \mathcal{N}(0, \sigma)$ to generate perturbed inputs $\hat{\mathbf{x}}_i = \mathbf{x} + \nu_i$, with σ increasing until the model behaves randomly (*i.e.*, accuracy = $1/C$). While the average trend (*blue* line) indicates that larger perturbation causes higher model distortion, sample-wise exceptions frequently appear. In Figure A.5, random sample trajectories reveal both correlated and uncorrelated patterns between perturbation levels and model distortions. This is a key observation, as it implies that model robustness cannot be assured by a general threshold without inspecting each evaluation sample individually.

Model Sensitivity Under Layer-by-Layer Randomisation. For sensitivity evaluations (Definition 3) to be meaningful, the model distortion caused by perturbation must be significant (Assumption 2). To test this practice, we perform consecutive layer-wise model parameter randomisation; in both a top-down (Adebayo et al., 2018), and bottom-up (Hedström et al., 2024) manner. From Figure A.3 (E), (F), and (G), and Figure A.4 (D) and (E), we observe that, although model distortion generally increases with layer-wise randomisation, there are exceptions of non-monotonicity (see, *e.g.*, Path, and Avila results in Figure A.3 (E), and (F), respectively). The high standard deviation (see the error bars) suggests that layer-wise randomisa-

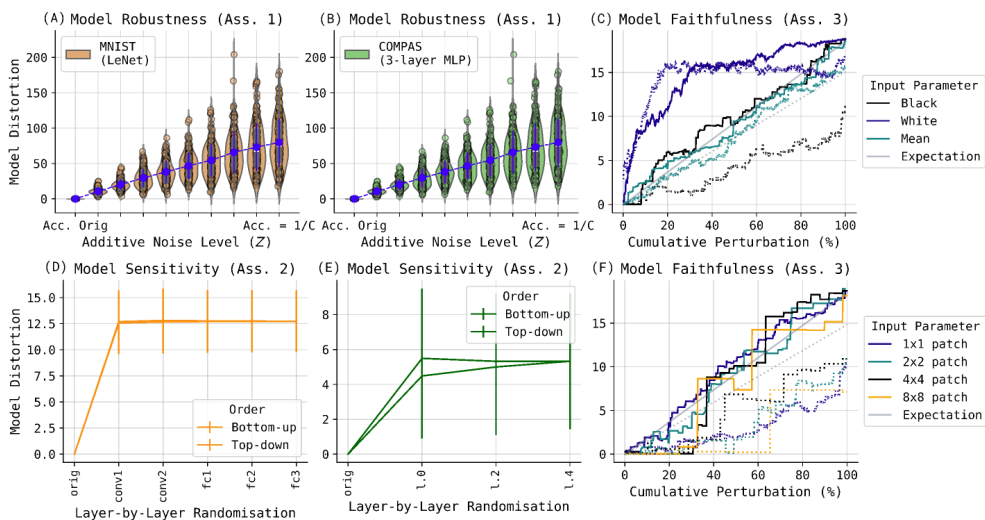


Figure A.4: Impact of model distortion (y -axis) over common perturbation types in robustness, sensitivity, and faithfulness evaluations, across different datasets, and NN architectures. (A), and (B) depict the distribution of model distortions across different perturbation magnitudes of additive Gaussian noise for MNIST (LeNet), and COMPAS (2-layer MLP), respectively. (D), and (E) show the average, and standard deviation of model distortions over layer-wise top-down, and bottom-up randomisation for the same datasets (as indicated by colour). (C), and (F) display model distortions for randomly chosen MNIST (*solid* line), and fMNIST (*dashed* line) samples (LeNet) under cumulative perturbations using different patch sizes (1×1 , 2×2 , 4×4 , 8×8), and baseline replacement strategies (*black*, *white*, *mean*).

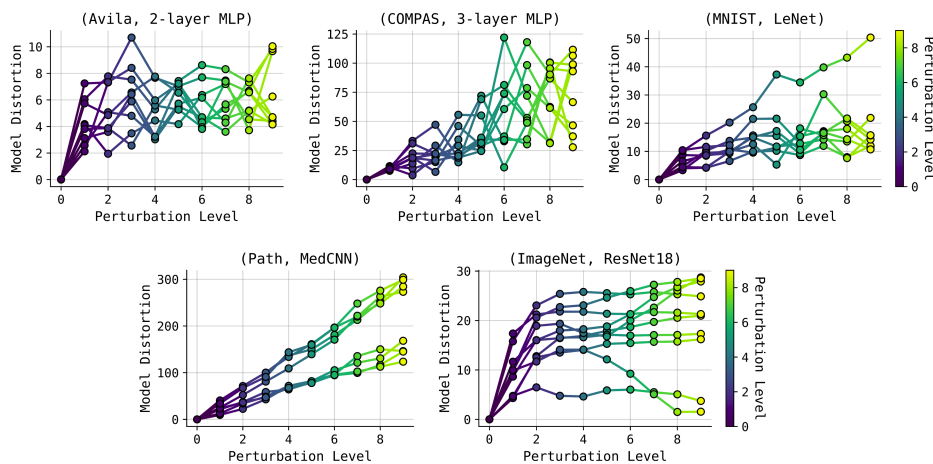


Figure A.5: Sample-wise trajectories of model distortion (y -axis) across different perturbation magnitudes of additive Gaussian noise. Each panel includes line plots across $N = 10$ samples, showing both correlated and uncorrelated outcomes.

tion fails to predictably dictate the degree of model distortion, undermining the assumption that significant model distortions will always occur in sensitivity evaluations.

Model Faithfulness Under Cumulative Input Perturbation. To investigate whether model distortion increases monotonically under cumulative input perturbation (Assumption 3), we measure D_f using a standard “pixel-flipping” faithfulness procedure (Bach et al., 2015). By randomising the perturbation order, the resulting faithfulness curve should reflect *only* the model’s response; any deviation from a linear trend suggests that Assumption 3 is failed. Observing Figure A.3 (D), and (H), and Figure A.4 (C) and (F), we see that neither patch size (top) nor replacement strategy (bottom) induces monotonic non-decreasing model behaviour. While these results are expected due to the model’s inherent nonlinearity, and OOD effects (Hase et al., 2021; Hesse et al., 2024), it is not accounted for in the faithfulness evaluation itself (Definition 4).

When genuine signals (*i.e.*, explanation quality) are not decoupled from noise (*i.e.*, non-monotonic model behaviour), interpretations may become biased (Hooker et al., 2019; Brocki & Chung, 2022; Brunke et al., 2020).

Together, these results reveal how easily, and systematically Assumptions 1-3 are violated by perturbation strategies commonly applied in practice (Section 2.1.1). Our findings are consequential as they demonstrate that the validity of existing robustness, sensitivity, and faithfulness evaluations (Definitions 2-4) are frequently undermined. As displayed in Figure A.3, there are many sample-wise exceptions where the perturbation magnitude, and the model distortion quantity are not strictly monotonically related, challenging the assumption that increased perturbations lead to proportionally greater distortions, and vice-versa.

A.5.1 Issues with Cumulative Input Perturbation

If small changes in input parameters, cause large variations in evaluation outcomes, evaluation reliability is compromised. Corroborating previous studies (Brunke et al., 2020; Brocki & Chung, 2022; Rong et al., 2022), the varied faithfulness curves in Figure A.3 (G), and (H), and Figure A.4 (C) and (F), demonstrate how input parameter choices, such as patch size or pixel value, can drastically influence the evaluation outcomes across tasks, *i.e.*, act as evaluation confounds (*cf.* the same parameter for MNIST *solid* line vs. fMNIST *dotted* line). These variations between tasks expose a simple, yet systematically overlooked issue in faithfulness evaluations: that parameter choices to perturb the input inherently introduce task-specific biases to the evaluation. Attempts to mitigate these biases—using inverse curves (Blücher et al., 2024) or assessing the OOD impact of perturbations (Qiu et al., 2021; Haug et al., 2021)—fail to address the core problem: that evaluation methods (Section 2.1.2) that require input parameters to be tuned according to its task, are inherently biased, impeding impartial comparisons across tasks, and explanation approaches.

A.6 Alignment Patterns, and Extended Results

Figure A.6 provides complementary results to Figure 3 in the main manuscript. The scatter points are coloured by perturbation magnitude up to $Z = 10$, using additive Gaussian noise. The varying but consistent overlaps of points of high and low perturbation magnitudes alongside the almost uniform distribution along the y-axis, illustrate that perturbation effects are not guaranteed to have a proportional effect on the model and explanation functions. Thus, using the perturbation magnitude as an indicator of the magnitude of which the explanation should change (as done in existing evaluations, see Definitions 2-3) is not reliable.

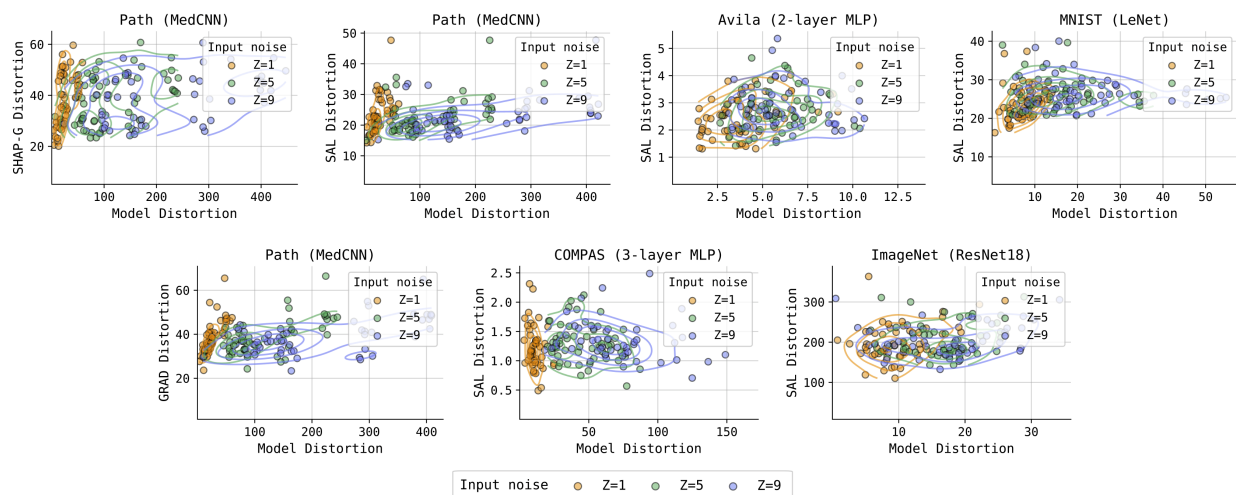


Figure A.6: Model (x-axis), and explanation distortions (y-axis) under varying levels of additive Gaussian input noise for vision and tabular tasks, as indicated in the titles. Scatter points represent individual samples, coloured by perturbation magnitude ($Z = 10$), with overlapping contours highlighting the relative alignment patterns. The individual plots contain SHAP-G, GRAD, and SAL explanations, as indicated by the y-axis labels.

Next, we show the change in the relationship between model distortion and explanation distortion when comparing the commonly used input perturbations with model perturbation (**Fast-GEF**) and when measuring explanation distortion and model distortion on the same geometry (**GEF**). Figure A.7 illustrates the relationship of model and explanation distortion across these three cases from left to right. Each contour plot includes $N = 10$ samples per perturbation magnitudes from $Z = 1$ to $Z = 5$. The scatter points are colored to one of three increasing perturbation magnitudes ($Z = 1$ to $Z = 3$). We can observe the distortion quantities across tasks with increasing model complexity, from a simpler tabular task (first row) to a highly parameterised model for a vision task (last row).

As expected, the first column (input perturbation) coincides with Figure A.6 and yields the same findings. However, in the following two columns, we can observe how **Fast-GEF** and **GEF** behave in practice. Most notably we observe that the alignment between model and explanation distortion changes from less complex to more complex tasks. Furthermore, we find that **Fast-GEF** tends to generate higher coherence compared to input perturbation, except for ImageNet, and that **GEF** yields the most coherent distortions. While these findings appear to support both our approaches, it is important to note that without access to ground truth, it is unclear whether the contour plots should necessarily show stronger coherence (*i.e.*, post-perturbed correlation), as it depends on the relationship between the explanation and model functions.

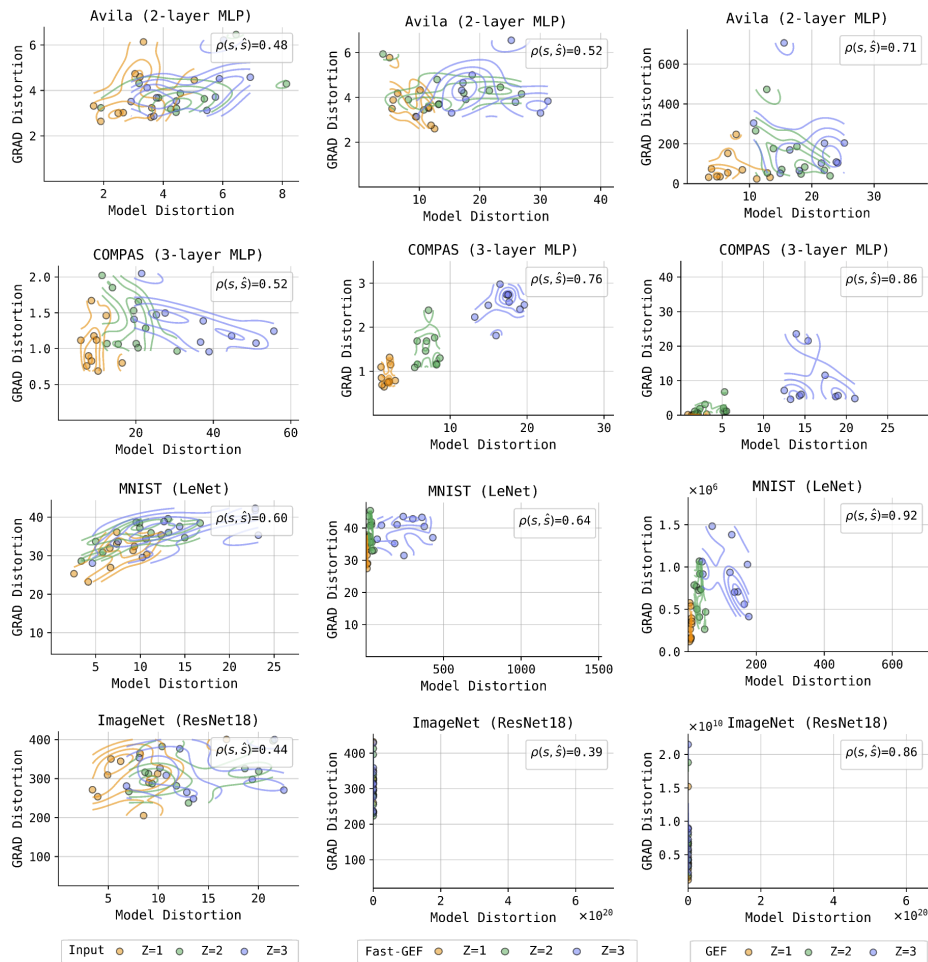


Figure A.7: Each plot shows the model (x-axis), and explanation distortions (y-axis) under different types of noise for *Gradient* (GRAD) explanation. The first column shows distortion outcomes after applying additive Gaussian input noise. The second and third columns show distortion outcomes after applying model parameter scaling (Section 5.1). The second column computes explanation distortion using **Fast-GEF** and the third column computes distortion using **GEF** (*i.e.*, with pullback mechanism). Scatter points represent individual samples, coloured by perturbation magnitude ($z=1, z=2, z=3$), with $Z = 5$ number of steps.

A.7 Ablation Study

To better understand the influence of the hyperparameters, on the proposed GEF evaluation method, we conducted an ablation study. We employed two tasks, *i.e.*, a tabular dataset (Avila) using a 2-layer MLP model with SAL explanations and a vision dataset (MNIST) using a LeNet model with 250 random explanations, sampled from a uniform distribution, *i.e.*, $\hat{e}_i = \mathcal{U}(0, 1)$. For each hyperparameter, *i.e.*, the number of perturbed models M , the length of the perturbation path Z , the number of summation steps T , and the number of samples K , we enumerated over values from 0 to 20, while fixing the others at a default value of 10. For each configuration, we recorded mean (solid line) and standard deviation (shaded area) of the model distortion, explanation distortion, Jacobian quantity, and the mean computation time.

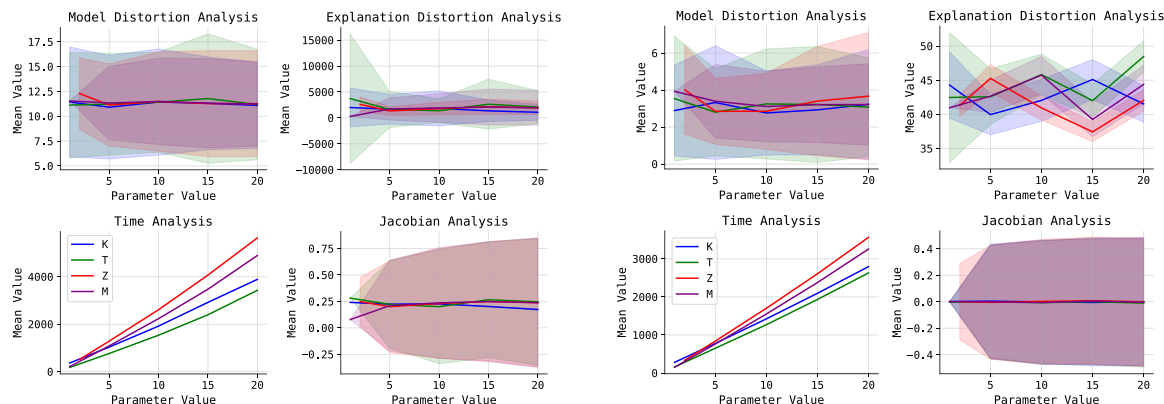


Figure A.8: Ablation study results across hyperparameters M , Z , T , and K for two tasks: (*left*) Saliency explanation on Avila (2-layer MLP), and (*right*) Random explanation on MNIST (LeNet). The mean value (*solid line*) and variance (*shaded area*) are reported. The time analysis is measured in seconds.

Figure A.8 demonstrates the results for both the tabular (*left*) and the vision task (*right*). As can be observed by the converging values of the standard deviation and means, the hyperparameters are resilient to key parameter changes once parameter values reach 5 or higher. The Jacobian variance reflects the curvature captured in its estimate. While the variance increases for parameter values above 5, the mean stabilises, indicating diminishing returns in capturing additional curvature. At parameter values of 5, the majority of the curvature is already captured, providing a practical trade-off between computational efficiency and quality of approximation. Considering computational time, all parameters lead to a linear, non-negligible increase. Among them, Z and M are identified as the primary drivers of time. Based on these experimental findings, setting $K = T = Z = M = 5$ balances computational efficiency and stability of the quality estimate.

A.8 Experiments, and Extended Results

This section provides descriptions of experimental setups, and extended results, including meta-evaluations, and agreement between different scoring methods. Additionally, we present further results for random control variant sanity checks, cross-domain benchmarking, and LLM-x methodology, and extended results.

A.8.1 Meta-Evaluation

To employ the scoring methodology (Section A.3), we used the pre-existing test suite available in the MetaQuantus library⁷ with their pre-defined hyperparameters.

MetaQuantus Hyperparameters. We applied these metrics over $K = 5$ perturbations, conducting 3 iterations with the test configurations specified in the library for two different sets of explanation methods, namely {GRAD, G-CAM}, and {SAL, SHAP-G}, which were evaluated by each metric. The explanation

⁷Find the library at <https://github.com/annahedstroem/MetaQuantus/>.

method groups were created by randomly selecting methods from the complete set of available methods, ensuring consistency across various experimental setups, such as dataset, and model combinations. In terms of choosing K , and the number of iterations, we followed the recommendations from the original study to keep the standard deviation between different sets relatively low. To ensure a fair comparison across metrics, all shared hyperparameters were assigned the same values.

Metrics Hyperparameters. All metrics have been implemented in **Quantus** (Hedström et al., 2023b). Different hyperparameters were chosen for the individual metrics based on the dataset. For the robustness metrics, we use 5 noisy samples, and employ additive Gaussian noise such that $\nu \sim \mathcal{N}(0, 0.001)$. For the faithfulness metrics, we use 28 features per perturbation step, and a patch size of 7 for the MNIST, and fMNIST datasets. For ImageNet, we set the number of features to 896, and the patch size to 28. For FC, similar to the robustness metrics, we let it run 5 times. For the sensitivity metrics, namely MPRT, and sMPRT, we use a noise magnitude of 0.01 for each sample, and sMPRT uses 5 samples in its calculation. For all sensitivity metrics, we use the Spearman rank correlation coefficient.

Table A.1: MC scores and standard deviation for unified, and faithfulness methods listed in A.4.5 for ImageNet, MNIST, and fMNIST datasets. The final row shows the mean score for each metric across the datasets. Values range between $[0, 1]$, with higher values indicating better outcomes. Due to computational constraints, **GEF** scores are only computed for fMNIST, and MNIST datasets.

	UNIFIED		FAITHFULNESS		
	GEF	Fast-GEF	PF	FC	RP
ImageNet	NAN \pm NAN	0.78 \pm 0.02	0.63 \pm 0.01	0.51 \pm 0.02	0.63 \pm 0.06
MNIST	0.75 \pm 0.07	0.74 \pm 0.03	0.61 \pm 0.04	0.63 \pm 0.03	0.59 \pm 0.03
fMNIST	0.71 \pm 0.07	0.71 \pm 0.03	0.63 \pm 0.01	0.50 \pm 0.04	0.58 \pm 0.09
Mean	0.73 \pm 0.07	0.74 \pm 0.03	0.62 \pm 0.02	0.56 \pm 0.03	0.59 \pm 0.06

Table A.2: MC scores and standard deviation for sensitivity, and robustness methods listed in A.4.5 for ImageNet, MNIST, and fMNIST datasets. The final row shows the mean score for each metric across the datasets. Values range between $[0, 1]$, with higher values indicating better outcomes.

	SENSITIVITY			ROBUSTNESS		
	MPRT	sMPRT	eMPRT	RIS	ROS	RRS
ImageNet	0.71 \pm 0.02	0.69 \pm 0.04	0.71 \pm 0.02	0.72 \pm 0.06	0.76 \pm 0.07	0.75 \pm 0.04
MNIST	0.63 \pm 0.02	0.66 \pm 0.04	0.76 \pm 0.03	0.73 \pm 0.02	0.70 \pm 0.09	0.74 \pm 0.09
fMNIST	0.63 \pm 0.01	0.67 \pm 0.05	0.67 \pm 0.05	0.70 \pm 0.02	0.77 \pm 0.06	0.70 \pm 0.03
Mean	0.64 \pm 0.02	0.67 \pm 0.04	0.71 \pm 0.03	0.72 \pm 0.03	0.74 \pm 0.07	0.73 \pm 0.05

Extended Results. In Tables A.1, and A.2, we provide the corresponding results for Figure 6.

A.8.2 Agreement between **GEF**, and **Fast-GEF**

To determine whether the simpler, computationally efficient **Fast-GEF** method can serve as an alternative to the more exact but computationally intensive **GEF** method, we compare the agreement between their respective faithfulness estimates. For a subset of explanation methods, and tasks (see Table 2), we thus compute scores, and rank explanation methods from R1 to RN. While it is expected that estimates from the two methods differ, a high agreement in a categorical ranking would make **Fast-GEF** a practical alternative in resource-constrained environments.

Results. Figure A.9 (A) visually compares how **GEF** and **Fast-GEF** ranks (x -axis) each explanation method in terms of increases (y -axis), highlighting the relative agreement between them. The explanations in the tabular, and text tasks show perfect ranking agreement. In the MNIST vision task, with minimal nominal differences, GRAD, and SHAP-G methods disagree in their ranking (R1, and R2), but such disagreement can be expected acknowledging the algorithmic similarity between these explanation methods. In the Derma vision task, the same pattern is observed, yet with a slightly larger difference for the global method FO-50. Interestingly, we observe that nominal differences are pronounced for global methods (DV-50, and FO-50), and that **Fast-GEF** tends to generate slightly lower faithfulness estimates *cf.* **GEF**.

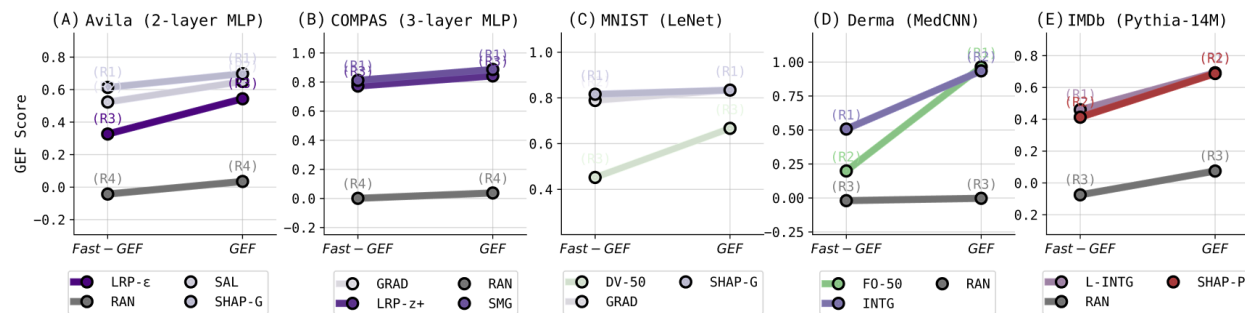


Figure A.9: (A) to (E) illustrates the GEF scores of GEF and Fast-GEF (with $M = 1$) for various explanation methods, and tasks. Explanation methods are ranked between R1 to RN, in descending order.

A.8.3 Scoring Control Variants

Next, we validate that both GEF and Fast-GEF assign low faithfulness scores to different control variant explanations. In our sanity checks, we evaluate explanations generated by uniform sampling, *i.e.*, $\hat{e}_i \sim \mathcal{U}(0, 1)$, a constant value, *i.e.*, $\hat{e}_i = \mathbf{0}$, and with a model-independent Sobel filter. For non-random reference, we evaluate GRAD explanations for the predicted class of the Derma task (see Table 2) (Sobel et al., 1968). For comparability, we extend this sanity check exercise to one metric per evaluative criteria, *i.e.*, FC (faithfulness), MPRT (sensitivity), and RIS (robustness). Hyperparameters are provided in Appendix A.8.1.

Table A.3: Evaluation scores of Derma (MedCNN) explanations for three random, and one regular (GRAD) explanation. The arrow (\uparrow, \downarrow) indicates whether higher or lower values are better. A nan value indicates that no score is produced.

EXPLANATION	GEF (\uparrow)	FAST-GEF (\uparrow)	FC (\uparrow)	MPRT (\downarrow)	RIS (\downarrow)
CONTROL VAR. CONSTANT	NAN \pm NAN	NAN \pm NAN	NAN \pm NAN	NAN \pm NAN	0.11 \pm 0.29
CONTROL VAR. RANDOM UNIFORM	-0.01 \pm 0.30	-0.01 \pm 0.22	-0.00 \pm 0.51	-0.00 \pm 0.04	3.21 \pm 2.89
CONTROL VAR. SOBEL FILTER	NAN \pm NAN	NAN \pm NAN	-0.01 \pm 0.50	1.00 \pm 0.00	82197.21 \pm 132718.26
GRAD	0.47 \pm 0.23	0.48 \pm 0.15	-0.05 \pm 0.49	0.01 \pm 0.04	1764.60 \pm 10007.26

Results. Table A.3 presents the results. Some metrics produce no values (nan), *e.g.*, when correlating identical vectors, and by that identify the unfaithful explanation. Fast-GEF, and GEF consistently assign low scores to random explanations, and high scores to non-random GRAD explanations, indicating their ability to identify the control explanations. Conversely, other metrics fail at least in one random test, either showing little discrepancy between regular, and control variants or even giving higher scores to the control. For instance, MPRT, and RIS score random uniform explanations as good or better than regular ones.

A.8.4 Cross-Domain Benchmarking

Extended Results. We benchmark various local, and global explanation methods with GEF and Fast-GEF. In Figure A.10, we extend the results in Figure 8.

The results presented in Figure 7 are provided in Tables A.4, and A.5.

A.8.5 LLM-x

In the following, we provided extended results of the LLM-x experiments.

Extended Results. The results presented in Figure 9 are provided in Tables A.6, and A.7.

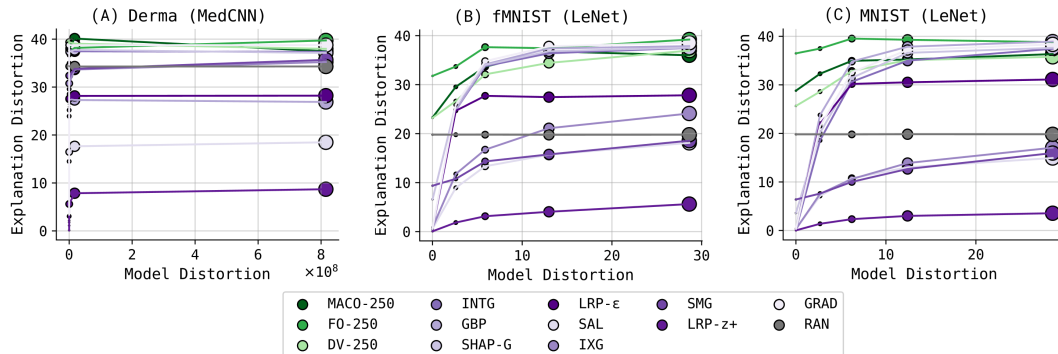


Figure A.10: **Fast-GEF** results for vision tasks. (A), (B), and (C) plot the model and explanation distortion for Derma (MedCNN), and fMNIST (LeNet), and MNIST (LeNet) along the perturbation path with $Z = 5$ perturbation steps. The size of the scatter point represents each perturbation steps, from 1 to 5.

Table A.4: **GEF** results on local methods for tabular tasks. Mean faithfulness scores, and standard errors are reported, with higher values indicating better quality.

	TASK	ADULT (3-LAYER MLP)	ADULT LR	AVILA (2-LAYER MLP)	COMPAS (3-LAYER MLP)	COMPAS LR
LOCAL METHODS	SMG	0.86 ± 0.00	0.69 ± 0.00	0.72 ± 0.01	0.81 ± 0.01	0.73 ± 0.01
	SHAP-G	0.78 ± 0.00	0.84 ± 0.01	0.75 ± 0.01	0.84 ± 0.00	0.66 ± 0.01
	SAL	0.84 ± 0.00	0.76 ± 0.00	0.69 ± 0.01	0.75 ± 0.01	0.70 ± 0.01
	RAN	-0.00 ± 0.02	0.00 ± 0.02	0.00 ± 0.02	0.02 ± 0.02	0.02 ± 0.02
	LRP- ϵ	0.66 ± 0.01	0.61 ± 0.01	0.52 ± 0.01	0.74 ± 0.01	0.59 ± 0.01
	LRP- z^+	0.79 ± 0.00	0.75 ± 0.00	0.66 ± 0.01	0.78 ± 0.00	0.70 ± 0.01
	IXG	0.84 ± 0.00	0.77 ± 0.01	0.69 ± 0.01	0.74 ± 0.00	0.72 ± 0.00
	INTG	0.82 ± 0.00	0.82 ± 0.00	0.69 ± 0.01	0.81 ± 0.00	0.80 ± 0.00
	GRAD	0.86 ± 0.00	0.74 ± 0.00	0.69 ± 0.01	0.81 ± 0.01	0.67 ± 0.01
	GBP	0.80 ± 0.00	0.60 ± 0.01	0.68 ± 0.01	0.78 ± 0.01	0.70 ± 0.01

Table A.5: **Fast-GEF** result on local methods for vision tasks. Mean faithfulness scores, and standard errors are reported, with higher values indicating better quality.

	TASK	DERMA MEDCNN	FMNIST LENET	IMAGENET-1K RESNET18	MNIST LENET	PATH MEDCNN
LOCAL METHODS	SMG	0.61 ± 0.01	0.69 ± 0.01	0.63 ± 0.01	0.73 ± 0.01	0.63 ± 0.02
	SHAP-G	0.49 ± 0.01	0.74 ± 0.01	0.69 ± 0.01	0.77 ± 0.01	0.60 ± 0.01
	SAL	0.67 ± 0.01	0.76 ± 0.01	0.71 ± 0.01	0.74 ± 0.01	0.65 ± 0.01
	RAN	0.01 ± 0.01	0.00 ± 0.01	-0.00 ± 0.01	-0.00 ± 0.01	-0.01 ± 0.01
	LRP- ϵ	0.45 ± 0.01	0.70 ± 0.01	0.35 ± 0.01	0.73 ± 0.01	0.66 ± 0.01
	LRP- z^+	0.74 ± 0.01	0.73 ± 0.01	0.91 ± 0.00	0.73 ± 0.01	0.71 ± 0.01
	IXG	0.73 ± 0.01	0.72 ± 0.01	0.64 ± 0.01	0.73 ± 0.01	0.71 ± 0.01
	INTG	0.49 ± 0.01	0.73 ± 0.01	0.78 ± 0.01	0.77 ± 0.01	0.71 ± 0.01
	GRAD	0.54 ± 0.01	0.76 ± 0.01	0.71 ± 0.01	0.79 ± 0.01	0.66 ± 0.01
	GBP	0.63 ± 0.01	0.76 ± 0.01	0.85 ± 0.01	0.77 ± 0.01	0.64 ± 0.01
GLOBAL METHODS	MACO-50	0.16 ± 0.02	0.47 ± 0.02	0.29 ± 0.02	0.29 ± 0.02	0.42 ± 0.01
	MACO-250	0.30 ± 0.02	0.54 ± 0.01	0.31 ± 0.01	0.31 ± 0.02	0.40 ± 0.01
	MACO-100	0.24 ± 0.01	0.52 ± 0.01	0.31 ± 0.01	0.41 ± 0.01	0.45 ± 0.01
	FO-50	0.17 ± 0.02	0.42 ± 0.01	0.22 ± 0.02	0.26 ± 0.02	0.35 ± 0.01
	FO-250	0.36 ± 0.01	0.38 ± 0.02	0.19 ± 0.02	0.15 ± 0.02	0.31 ± 0.01
	FO-100	0.28 ± 0.01	0.36 ± 0.02	0.23 ± 0.02	0.21 ± 0.02	0.27 ± 0.01
	DV-50	0.38 ± 0.01	0.54 ± 0.02	0.26 ± 0.02	0.36 ± 0.02	0.45 ± 0.01
	DV-250	0.44 ± 0.01	0.50 ± 0.01	0.40 ± 0.02	0.40 ± 0.02	0.48 ± 0.01
	DV-100	0.43 ± 0.02	0.49 ± 0.02	0.40 ± 0.02	0.43 ± 0.02	0.51 ± 0.01

Table A.6: **Fast-GEF** results on LLM-x, and local methods for top- K tasks. Mean faithfulness scores, and standard errors are reported, with higher values indicating better quality.

	TASK	SMS SPAM	SST2
		BERT-TINY FT	BERT-TINY FT
LOCAL METHODS	SHAP-P-5	0.62 ± 0.01	0.75 ± 0.01
	SHAP-P-10	0.62 ± 0.01	0.75 ± 0.01
	RAN-5	0.08 ± 0.01	-0.08 ± 0.01
	RAN-10	0.03 ± 0.01	-0.10 ± 0.01
	LLM-X-5	0.06 ± 0.02	0.05 ± 0.02
	LLM-X-10	0.05 ± 0.02	0.08 ± 0.02
	L-INTG-5	0.58 ± 0.01	0.77 ± 0.01
	L-INTG-10	0.58 ± 0.01	0.77 ± 0.01

Table A.7: **Fast-GEF** results on LLM-x, and local methods for top- K tasks. Mean faithfulness scores, and standard errors are reported, with higher values indicating better quality.

TASK	SMS SPAM	SST2
	BERT-TINY FT	BERT-TINY FT
LLM-X	-4.25 ± 8.42	-3.73 ± 7.72
L-INTG	195.49 ± 2.91	238.15 ± 2.98
SHAP-P	185.95 ± 3.06	230.27 ± 3.75

A.9 Notation Tables

All notations used in this paper is provided in the following.

Spaces, and Elements

\mathcal{X}, \mathbf{x}	The input space $\mathcal{X} \subseteq \mathbb{R}^D$ with a sample $\mathbf{x} \in \mathcal{X}$
\mathcal{F}, θ	The model space $\mathcal{F} \subseteq \mathbb{R}^U$ with parameters $\theta \in \mathcal{F}$
\mathcal{Y}, \mathbf{y}	The function output space $\mathcal{Y} \subseteq \mathbb{R}^C$ with logits $\mathbf{y} \in \mathcal{Y}; \mathbf{y} = [y_1, \dots, y_C]^T$ for C classes $y_c \in \mathcal{Y} \forall c \in [1, C]$
\mathcal{E}, \mathbf{e}	The explanation space $\mathcal{E} \subseteq \mathbb{R}^V$ with an explanation $\mathbf{e} \in \mathcal{F}$
\mathcal{Q}, q	The evaluation space $\mathcal{Q} \subseteq \mathbb{R}^M$ with a quality estimate $q \in \mathcal{Q}$
\mathcal{S}, \mathbf{s}	A set of spaces $\mathcal{S} \subset \{\mathcal{X}, \mathcal{F}, \mathcal{Y}, \mathcal{E}, \mathcal{Q}\}$ where $\mathcal{S} \subseteq \mathbb{R}^S, \mathcal{S} \in \mathbb{N}$ with $\mathbf{s} \in \mathcal{S}$
\mathcal{H}, \mathbf{h}	A subset of spaces $\mathcal{H} \subseteq \{\mathcal{F}, \mathcal{E}\}$ with $\mathbf{h} \in \mathcal{H}$
$\hat{\mathbf{s}}, \hat{\mathbf{x}}, \hat{\theta}, \hat{\mathbf{y}}, \hat{\mathbf{e}}$	A sample, input, parameters, logit, explanation, post-perturbation.

Functions

f	A classifier function $f : \mathcal{X} \rightarrow \mathcal{Y}$ with $f(\mathbf{x}; \theta) = \mathbf{y}$ (we refer f_θ as f), parameterised by θ
ϕ_L	A local explanation function $\phi_L : \mathcal{F} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^V$ with $\phi_L(f, \mathbf{x}, \mathbf{y}; \lambda) = \mathbf{e}$, parameterised by λ
ϕ_G	A global explanation function $\phi_G : \mathcal{F} \times \mathcal{Y} \rightarrow \mathbb{R}^V$ with $\phi_G(f, \mathbf{y}; \kappa) = \mathbf{e}$, parameterised by κ
ϕ	Collectively, denoting ϕ_L , and ϕ_G although they formally reside in different spaces
Ψ	An evaluation function $\Psi : \mathcal{E} \times \mathcal{X} \times \mathcal{F} \times \mathcal{Y} \rightarrow \mathbb{R}$ with $\Psi(\mathbf{e}, \mathbf{x}, f, \mathbf{y}; \tau) = \mathbf{q}$, parameterised by τ
\mathcal{P}_S	A perturbation function $\mathcal{P} : \mathcal{S} \rightarrow \mathcal{S}$ where $\mathcal{P}(\mathbf{s}; \omega)$ on space \mathcal{S}
δ	A general discrepancy function $\delta : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ with $\delta(\mathbf{s}, \hat{\mathbf{s}}) = \xi$, parameterised by $\omega \in \mathbb{R}$
k	A separate mapping function $k : \mathcal{S} \rightarrow \mathcal{H}$ mapping $\mathbf{s}, \hat{\mathbf{s}}$ to a distinct space \mathcal{H}
D_k	A functional distortion $D_k : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ with $D_k(\mathbf{s}, \hat{\mathbf{s}}) = \delta(k(\mathbf{s}), k(\hat{\mathbf{s}}))$
ρ	A correlation function with $\rho : \mathbb{R}^Z \times \mathbb{R}^Z \rightarrow \mathbb{R}$

Constants

C	The number of classes
D	The dimension of the input
W	The dimension of the parameter vector
V	The dimension of the explanation outputs
Z	The number of perturbation steps
K	The number of samples to approximate the Jacobian
T	The number of integral steps between two points, e , and \hat{e}
M	The number of models to average over in GEF , and Fast-GEF

Variables

ξ	The perturbation magnitude defined as the discrepancy $\delta(\mathbf{s}, \hat{\mathbf{s}}) = \xi$ between $\hat{\mathbf{s}}$, and \mathbf{s}
\mathbf{D}_f	The model distortion \mathbf{D}_f across parameter- $\mathbf{D}_f(\theta, \hat{\theta})$, and input perturbation $\mathbf{D}_f(\mathbf{x}, \hat{\mathbf{x}})$
\mathbf{D}_ϕ	The explanation distortion \mathbf{D}_ϕ across parameter- $\mathbf{D}_\phi(\theta, \hat{\theta})$, and input perturbation $\mathbf{D}_\phi(\mathbf{x}, \hat{\mathbf{x}})$
$\varepsilon_{\mathbf{D}_k}^{RO}$	The implicit upper boundary value with $\varepsilon^{RO} \in \mathbb{R}^+$, and $k \in \{\phi, f\}$ used in robustness
$\varepsilon_{\mathbf{D}_k}^{SE}$	The implicit lower boundary value with $\varepsilon^{SE} \in \mathbb{R}^+$, and $k \in \{\phi, f\}$ used in sensitivity
α	A boundary value for the perturbation magnitude, with $\alpha \in \mathbb{R}^+$
η_i	The Gaussian noise matrix with $\eta_i \sim \mathcal{N}(\mathbf{0}, \sigma_i^2 \mathbf{1})$
σ_z^2	The covariance scale of a Gaussian distribution with $\sigma_z^2 \in \mathbb{R}^+$ at z^{th} perturbation
J_f	The network Jacobian for fixed input \mathbf{x} , and model f , with $J_f \in \mathbb{R}^{V \times C}$, and elements $J_{i,j} = \frac{\partial e_i}{\partial f_j}$
\mathbf{g}	Pullback metric tensor based on the elementwise Jacobian with $\mathbf{g} \in \mathbb{R}^{V \times V}$
z	Index of perturbation steps with $z \in [1, Z]$
\mathbf{D}_f^z	The model distortion at perturbation step z with $\mathbf{D}_f^z := \mathbf{D}_f^z(\theta, \hat{\theta}_z)$
\mathbf{D}_ϕ^z	The explanation distortion at perturbation step z with $\mathbf{D}_\phi^z := \mathbf{D}_\phi^z(\theta, \hat{\theta}_z)$
\mathbf{d}_f	The vector of model distortion with Z steps, $\mathbf{d}_f = [\mathbf{D}_f^1, \mathbf{D}_f^2, \dots, \mathbf{D}_f^Z]$
\mathbf{d}_ϕ	The vector of explanation distortion with Z steps, $\mathbf{d}_\phi = [\mathbf{D}_\phi^1, \mathbf{D}_\phi^2, \dots, \mathbf{D}_\phi^Z]$