

# IS FORWARD GRADIENT AN EFFECTIVE TOOL FOR EXPLAINING BLACK-BOX MODELS?

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Gradients are widely used to explain the decisions of deep neural networks. However, as models become deeper and more complex, computing gradients become challenging and sometimes infeasible, hindering traditional explanation methods. Recently, the forward gradient method has garnered attention for training structure-agnostic models with discontinuous objective functions. This method perturbs only the parameters of interest for gradient computation and optimization. Inspired by this, we investigate whether the forward gradient can be employed to explain black-box models. In this work, we use the likelihood ratio method to estimate output-to-input gradients and utilize them for the explanation of model decisions. Additionally, we propose block-wise computation techniques to enhance estimation accuracy. Extensive experiments to explain various models, including CNNs, LSTMs, Phi-3, and CLIP, in black-box settings validate the effectiveness of our method, demonstrating good gradient estimation and improved explainability under the black-box setting.

## 1 INTRODUCTION

Deep Neural Networks (DNNs) have achieved remarkable success across a range of applications, including autonomous driving (Grigorescu et al., 2020; Mozaffari et al., 2020; Huang & Chen, 2020), facial recognition (Mehdipour Ghazi & Kemal Ekenel, 2016; Fathallah et al., 2017; Mellouk & Handouzi, 2020), and clinical diagnostics (De Fauw et al., 2018; Van der Laak et al., 2021; Kermany et al., 2018). In these safety-critical domains, it is imperative that DNNs not only deliver high task performance but also provide transparent and understandable decision-making processes (Tambon et al., 2022; Borg et al., 2018). Extensive research has been devoted to demystifying DNNs, exploring various approaches such as counterfactual explanations (Mothilal et al., 2020; Slack et al., 2021), attribution and activation analysis (Achtibat et al., 2023; Lin et al., 2021), and saliency maps (Hu et al., 2023; Tjoa & Cuntai, 2022; Lorentz et al., 2021). Of these, gradients play a crucial role in decision explanation because of their robust theoretical foundations (Khorram et al., 2021; Kapishnikov et al., 2021) and superior performance. Moreover, gradient-based explanation methods are data-centric, offering greater adaptability to different model architectures. This flexibility is particularly promising in elucidating the decisions from black-box models, which are either proprietary and cloud-based (OpenAI, 2020) or computationally intensive (Radford et al., 2021).

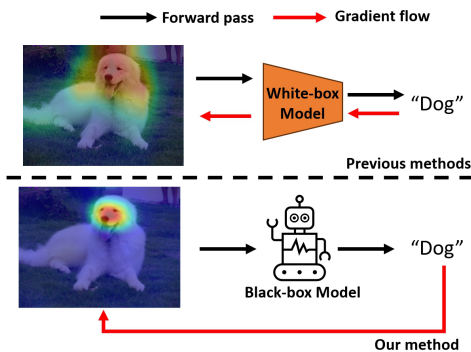


Figure 1: Previous methods require the full knowledge of studied models to compute the gradient for model explanation, which makes it impossible to be applied to explain the black-box model. Our proposed method only needs one forward pass of the studied model. It can directly get the gradient without relying on the backward pass and the knowledge of model architectures, thus enabling the black-box model decision.

The use of gradients in model explanation can be broadly divided into two categories: class activation mapping (Zhou et al., 2016; Chattopadhyay et al., 2018; Wang et al., 2020; Omeiza et al., 2019) and saliency map (Sundararajan et al., 2017; Erion et al., 2021; Lundstrom et al., 2022; Yang et al., 2023). Methods based on class activation map up-sample the feature map inner models for model explanation. Gradient-based methods derive the gradients of inputs relative to the output, thereby highlighting the model’s areas of interest. These two approaches are intuitive and efficient but become infeasible when the model architecture is inaccessible.

In black-box machine learning scenarios (Bodria et al., 2023), where access to full knowledge of computation graphs for optimization is unavailable, the problem becomes a Zeroth-order Optimization (ZO) problem (Chen et al., 2017). Among various ZO strategies (Chen et al., 2019; Rando et al., 2024; Kozak et al., 2023; Vemula et al., 2019), the likelihood ratio method (Peng et al., 2022) stands out. It pushes the parameters out of the loss function by introducing noise, allowing for unbiased gradient estimation. This method has shown competitive performance in training diverse neural networks, comparable to that achieved with backpropagation (Rumelhart et al., 1986). Such promising results inspire us to think about the plausibility of the likelihood ratio method for gradient estimation on inputs, facilitating black-box explanations through saliency maps. Motivated by this, we propose to use the likelihood ratio method to bridge the gap between gradient-based explanation methods and black-box scenarios, as illustrated in Fig. fig. 1, where we skip the access of backward pass of the model to get the gradient and explain the model decision.

Applying the likelihood ratio method to develop interpretation methods under black-box conditions is a non-trivial task, presenting two significant challenges. First, there is the question of how to estimate gradients on various characteristics of data using the likelihood ratio method. Many existing studies on zeroth-order optimization assume full knowledge of the local computational structure (Peng et al., 2022; Jiang et al., 2023), a luxury not available in our context where the whole model’s structure, even the first layer for processing the input, is unknown. Second, we must address how to minimize the estimation variance of gradients. The likelihood ratio method, by its nature of noise injection, tends to have a high variance in gradient estimation (Peng et al., 2022; Jiang et al., 2023), a problem exacerbated by the high dimensionality, such as those in the  $\mathbb{R}^{224 \times 224 \times 3}$  space of the ImageNet-1K dataset. Accurately estimating gradients in such high-dimensional spaces is challenging.

In this study, we introduce a unified framework for explaining black models using gradients. Our approach begins with deriving estimated gradients of inputs under a black-box model set using the likelihood ratio method. Then, we propose a blockwise computation pattern to address the challenge of high variance in gradient estimation. In evaluation, we integrate our framework to explain various models, including a series of computer vision models, text classification models, and the large language model Phi (Abdin et al., 2024). We also verify the feasibility and generalization ability in explaining multi-modal models against spurious correlations. There are three key findings revealed in our work which can be summarized as follows:

- **✔Interpretability of forward gradients:** Compared to vanilla gradients, forward gradients offer comparable interpretability with greater scalability in explaining model decisions.
- **✔Bridging the explanation gap:** Forward gradients provide a crucial link between white-box explanation methods and the challenges of explaining black-box models.
- **✘High computation and memory-consumption:** Despite their explanatory power, forward gradients suffer from high computational and memory costs, due to large estimation variance. This issue can only be mitigated by using a large number of copies during computation.

## 2 RELATED WORK

### 2.1 GRADIENTS FOR MODEL EXPLANATION

Gradients are widely used in explaining models, including the class activation mapping (CAM) (Zhou et al., 2016; Chattopadhyay et al., 2018; Wang et al., 2020; Omeiza et al., 2019) and the saliency map (Sundararajan et al. (2017); Erion et al. (2021); Lundstrom et al. (2022); Yang et al. (2023)).

CAM is first introduced in (Zhou et al., 2016), which employs the global average pooling neuron activations to weakly localize the objects in the inputs for saliency mapping. The GradCAM (Selvaraju et al., 2017) inherits the idea of CAM and combines the model gradients with the activations of

its internal neurons to compute the saliency maps. The GradCAM has been adopted and improved by a lot of later works. Such as GradCAM++(Chattopadhyay et al., 2018), SmoothGrad(Smilkov et al., 2017), Smooth-Grad++ (Omeiza et al., 2019), and Grad-CAM++(Chattopadhyay et al., 2018) to further boost the visualization ability. CAM requires the full knowledge of model architecture and feature maps during the inference process, thus failing to work under the black-box setting(Belharbi et al., 2022; Linardatos et al., 2020; Jiang et al., 2021).

The generation of saliency maps (Sundararajan et al., 2017; 2016) requires the gradients of the decision as inputs, which treat the model as a unified module regardless of the inner structure. Compared with directly leveraging the vanilla gradients, the use of integrate gradients is more popular. These methods gather the gradients of images along a path from the input to a reference image. The reference image represents what it looks like when the features present in the input are missing. Sundararajan et al. (2017) first proposed the use of the path attribution strategy, which can preserve certain desirable axiomatic properties for the explanations. This work inspires others to further investigations (Erion et al., 2021; Lundstrom et al., 2022; Yang et al., 2023).

## 2.2 ZEROth-ORDER OPTIMIZATION

Zeroth-order optimization (Wang et al., 2018; Golovin et al., 2019; Chen et al., 2019) is proposed to address challenges in optimizing the non-differential criteria. It typically uses random perturbation on input to estimate the gradient and follows the stochastic gradient descent to update the parameters of interest. Popular zeroth-order optimization methods include the simultaneous perturbation stochastic approximation (SPSA) (Spall, 2000; 1997; Maryak & Chin, 2001; Granichin & Amelina, 2014), evolutionary strategy (ES) (Salimans et al., 2017), and likelihood ratio (LR) (Peng et al., 2022; Jiang et al., 2023). SPSA individually perturbs each parameter and estimates the gradients. The low efficiency largely hinders its application in large-scale optimization problems. ES and LR are proposed to optimize deep learning models, which inject the noise into neuron weights and outputs respectively. Among these, LR stands out for its compelling performance in accurate gradient estimation (Jiang et al., 2023). In our work, we employ the LR to compute the gradient without reliance on the chain-rule to explain model decisions.

## 3 METHODOLOGY

We take the crafting of saliency map in computer vision as an example of using gradients to explain the model decision. In section 3.1, we present the preliminaries of the saliency map and why the application to black-box models is limited. In section 3.2, we introduce our framework based on likelihood ratio gradient estimation to explain black-box models. In section 3.3, we propose the blockwise computation manner that effectively enhances the quality of the saliency map generation, especially for high-dimensional inputs and large black-box models.

### 3.1 PRELIMINARIES

**Functions of gradients as saliency maps.** Given any image-label pair example  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ , with  $\mathcal{X} \subset \mathbb{R}^d$  representing the domain of input images and  $\mathcal{Y} = \{1, \dots, C\}$  denoting the set of possible label classes, let us consider a classification neural network,  $f: \mathcal{X} \rightarrow \mathbb{R}^C$ , which assigns a predicted class activation  $\hat{y}_c = f_c(x)$  for each class  $c \in \mathcal{Y}$ . Numerous studies Sundararajan et al. (2017); Erion et al. (2021); Lundstrom et al. (2022); Yang et al. (2023) have investigated the use of gradient-based saliency maps to visualize and elucidate the predictive mechanisms of the network  $f$ . Provided that the neural network  $f$  is differentiable almost everywhere, a gradient-based saliency map for a class  $c$  can be represented as a function dependent on the gradients of  $f$  w.r.t. inputs at dimension  $c$ , and a particular input instance  $x_0$  as follows,

$$M_c(f, x_0) = S(g_c, x_0), \quad (1)$$

where  $g_c(x_0) = \nabla_x f_c(x)|_{x=x_0}$  is the gradients of  $f$  w.r.t. inputs at  $x_0$ .

**Inaccessibility of black-box models' gradients.** The gradients of neural networks w.r.t. inputs are critical, as they quantify the impact of a tiny change in the input on the output and, therefore, can be linked intuitively to the importance of various input features. However, the reliance on these gradients poses a significant challenge in the context of black-box models, where the gradients are inaccessible

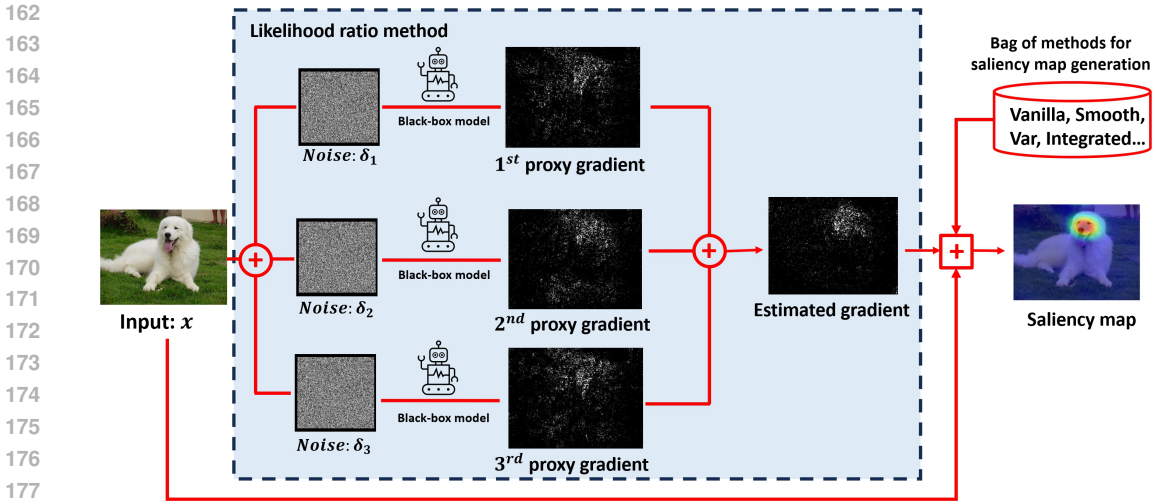


Figure 2: While the gradients of the black-box model can not be directly fetched, we use the likelihood ratio method to compute the proxy gradient of each noisy copy and obtain the estimated gradient. We incorporate the estimated gradient with gradient-based saliency map generation methods for the model explanation.

or even the internal workings are obscured. In fact, any use of black-box modules in the model would obstruct gradient computation due to the inapplicability of the chain rule, a fundamental principle in differential calculus. This limitation severely restricts the application of gradient-based saliency map techniques in the black-box setting. Therefore, the key problem in leveraging saliency map methods to interpret such models is to escape the inaccessibility of the black-box module to compute the gradients w.r.t. inputs.

### 3.2 LIKELIHOOD RATIO METHOD FOR SALIENCY MAP OF BLACK-BOX MODELS

In this section, we introduce a novel, unified framework for generating gradient-based saliency maps to interpret the decision of black-box models. As depicted in Fig. fig. 2, our framework encompasses three steps. Initially, we introduce perturbations to the images by injecting the noise before the model’s forward pass. Subsequently, we compute the proxy gradient of each perturbed copy respectively and then average them to approximate the true gradient. Finally, we employ the estimated gradient to craft the saliency map for the model explanation. The following paragraphs will detail the three steps in our proposed framework.

**Injecting noise to inputs.** Within the context of a black-box model  $f$ , which is differentiable almost everywhere but whose internal gradients are inaccessible, consider that we employ a technique of noise injection, in which we add small random noise  $z$  into the input  $x$ . Intuitively, if the noise has a neutral mean (zero) and a variance small enough, the expectation of gradient w.r.t. the noise-added input,  $\mathbb{E}_z(g_c(x + z))$ , would be close enough to the true gradient  $g_c(x)$ . Practically, one can easily verify that as the standard deviation  $\sigma$  of the noise approaches zero, the expectation of gradient w.r.t. the noise-added input converges to the true gradient, *i.e.*,  $\lim_{\sigma \rightarrow 0^+} \mathbb{E}_z(g_c(x + z)) = g_c(x)$ , implying the use of  $\mathbb{E}_z(g_c(x + z))$  as a viable approximation of  $g_c(x)$ . This technique allows us to control the precision loss of expected noise-adding gradients by selecting noise  $z$  with a distribution close to 0. In many situations, such as the Gaussian noise  $\mathcal{N}(0, \sigma^2 \mathbb{I})$ , it is equivalent to adjusting the  $\sigma$ .

**Likelihood ratio gradient estimator.** Directly computing  $\mathbb{E}_z(g_c(x + z))$  still necessitates access to the gradient, which exactly contradicts our objective of addressing the black-box setting where such information is unavailable. Therefore, we propose our likelihood ratio gradient estimator for black-box models. Initially, by forwarding the noise-added input into the model, we obtain the corresponding class activation  $f_c(x_0 + z)$ . Then, we compute the proxy gradient for each noise-added input by multiplying the class activation with the negative gradient of the noise’s log probability density function. Finally, we average  $n$  samples of proxy gradients to form our likelihood ratio

216 gradient estimator. Mathematically, this can be formalized as:

$$217 \hat{g}_c^{LR}(x_0) := \frac{1}{n} \sum_i^n (-f_c(x_0 + z_i) \nabla_z \ln \mu_z(z_i)) \quad (2)$$

220 where  $n$  is the number of perturbed sample, *i.e.*, copies,  $\{z_i\}_{i=1}^n$  represents the i.i.d. random noise,  
221 and  $\mu_z(\cdot)$  denotes the probability density function of the injected  $z$ . Notably, when the injected  
222 noise  $z_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2 \mathbb{I})$ , the likelihood ratio gradient estimator can be simplified into  $\hat{g}_c^{LR}(x_0) =$   
223  $\frac{1}{n\sigma^2} \sum_i^n f_c(x_0 + z_i) z_i$ .

224 Additionally, we introduce the following theorem 1, establishing the foundation for our likelihood  
225 ratio gradient estimator for black-box models.

227 **Theorem 1** Assume that  $\lim_{|\zeta| \rightarrow \infty} f_c(x_0 + \zeta) \mu_z(\zeta) = 0$  for any input  $x_0$ . Let  $z$  denote the random  
228 noise with the same distribution as  $\{z_i\}_1^n$ . Then, we have

$$229 \mathbb{E}(\hat{g}_c^{LR}(x_0)) = \mathbb{E}_z(g_c(x_0 + z)). \quad (3)$$

231 **Proof 1** Notice that a direct corollary from  $\lim_{|\zeta| \rightarrow \infty} f_c(x_0 + \zeta) \mu_z(\zeta) = 0$  is

$$232 \int_{\mathbb{R}^d} \nabla_\zeta f_c(x_0 + \zeta) \mu_z(\zeta) d\zeta + \int_{\mathbb{R}^d} f_c(x_0 + \zeta) \nabla_\zeta \mu_z(\zeta) d\zeta \quad (4)$$

$$233 = \int_{\mathbb{R}^d} \nabla_\zeta (f_c(x_0 + \zeta) \mu_z(\zeta)) d\zeta = 0.$$

234 Therefore, we can derive

$$235 \mathbb{E}_z(g_c(x_0 + z)) = \int_{\mathbb{R}^d} \nabla_x f_c(x)|_{x=x_0+\zeta} \mu_z(\zeta) d\zeta \quad \triangleright \text{by definition}$$

$$236 = \int_{\mathbb{R}^d} \nabla_\zeta f_c(x_0 + \zeta) \mu_z(\zeta) d\zeta \quad \triangleright \text{by change of variable}$$

$$237 = - \int_{\mathbb{R}^d} f_c(x_0 + \zeta) \nabla_\zeta \mu_z(\zeta) d\zeta \quad \triangleright \text{by eq. (4)}$$

$$238 = \mathbb{E}_z(-f_c(x_0 + z) \nabla_z \ln \mu_z(z)) \quad \triangleright \text{by derivative of ln}$$

$$239 = \mathbb{E}(\hat{g}_c^{LR}(x_0)). \quad \triangleright \text{by noise}$$

240 theorem 1 indicates that the expectation of our proposed gradient estimator is the same as the  
241 expectation of gradient w.r.t. the noise-added input under a certain asymptotic growth condition, which  
242 is commonly satisfied in practice. This, together with the previous discussion about  $\mathbb{E}_z(g_c(x_0 + z))$ ,  
243 substantiates the feasibility of utilizing  $\hat{g}_c^{LR}$  as an estimator of the true gradient. While this approach  
244 introduces a certain bias, as discussed, such impacts can be effectively mitigated by carefully  
245 controlling the standard deviation of the injected noise.

246 **Integrating with gradient-based saliency map methods.** By injecting random noise into the input  
247 alongside our likelihood ratio gradient estimator, we successfully circumvent the reliance on direct  
248 access to the gradients of black-box models. This approach allows for the seamless integration of any  
249 gradient-based saliency map technique within the context of black-box models.

250 To illustrate, consider an arbitrary gradient-based method that generates saliency map  $M_c$  for a  
251 model  $f$  and an image  $x$  from a function  $S$  of the model’s gradient w.r.t. inputs and that image,  
252  $M_c(f, x) = S(g_c, x)$ , as depicted in eq. (1). For instance, in vanilla gradient Simonyan et al.  
253 (2014) we have  $S(g_c, x) = g_c(x)$  while integrated gradient Sundararajan et al. (2017) defines  
254  $S(g_c, x) = (x - x') \odot \int_0^1 g_c(x' + \alpha(x - x')) d\alpha$ , where  $x'$  is a predetermined baseline. Then, we  
255 substitute  $g_c$  in the function  $S$  with our likelihood ratio gradient estimator  $\hat{g}_c$  defined in eq. (2),  
256 thereby enabling saliency map generation in the black-box setting, formalized as follows:

$$257 M_c^{LR}(f, x) = S(\hat{g}_c^{LR}, x). \quad (5)$$

### 258 3.3 BLOCKWISE COMPUTATION FOR SCALABILITY

259 **Challenges with estimation variance.** The aforementioned likelihood ratio method relies on the  
260 correlations between injected noise and final outputs to estimate the gradients of interested parameters  
261

or the inputs. Nevertheless, as the model becomes more complex and the input dimension  $d$  increases—a typical scenario in scaled-up models—the estimation variance,  $\text{Var}([\hat{g}_c^{\text{LR}}]_i)$ ,  $i \in \{1, \dots, d\}$ , becomes unbearable. This makes it difficult to implement in real-world applications, notably including the craft of saliency maps. Despite the requirement for only coarse gradients, high estimation variance in saliency maps associated with large images can yield vastly inconsistent attribution conclusions.

**Enhancing variance reduction through blockwise estimator** To make it scalable for interpreting the model decisions on high dimensional inputs, we further introduce a blockwise adaption of the likelihood ratio method for saliency map generation. Concretely, this approach involves initially selecting a small segment (referred to as a “block”) randomly on the original image. We adopt the sampling that ensures each pixel in the image has an equal probability  $q$  of being covered by the block. Then, we inject noise exclusively into the block area and leave the remainder of the image unchanged. Subsequently, we calculate the average of likelihood ratio proxy gradients across multiple random blocks and noise instances to form the resulting blockwise likelihood ratio gradient estimator, formalized as follows:

$$\hat{g}_c^{\text{BLR}}(x_0) := \frac{1}{nq} \sum_{i=1}^n (-f_c(x_0 + J_i \odot z_i) \nabla_z \ln \mu_z(J_i \odot z_i)), \quad (6)$$

where  $\{J_i\}_1^n$  represents the masks for random blocks—specifically,  $[J_i]_k = 1$  if the pixel  $k$  is covered by the block and 0 otherwise. Similarly, when  $\{z_i\}_1^n$  are Gaussian noise, we can simplify the blockwise likelihood ratio gradient estimator into  $\hat{g}_c^{\text{BLR}}(x_0) = \frac{1}{nq\sigma^2} \sum_{i=1}^n (f_c(x_0 + J_i \odot z_i) J_i \odot z_i)$ .

Under the same asymptotic growth condition, we derive the following theorem 2 for the blockwise likelihood ratio gradient estimator. This theorem leads us to a direct but critical corollary,  $\lim_{\sigma \rightarrow 0^+} \mathbb{E}(\hat{g}_c^{\text{BLR}}) = g_c(x)$ , ensuring that blockwise estimator’s efficacy in accurately estimating the gradient of black-box models with controlled standard deviation of integrated noise. The proof of the corollary is provided in detail in appendix B.1.

**Theorem 2** Assume that  $\lim_{|\zeta| \rightarrow \infty} f_c(x_0 + \zeta) \mu_z(\zeta) = 0$  for any input  $x_0$ . Then, for any noise  $z$  that is independent between dimensions, we have

$$\mathbb{E}(\hat{g}_c^{\text{BLR}}(x_0)) = \mathbb{E}_{z,J} \left( \frac{1}{q} J \odot g_c(x_0 + J \odot z) \right). \quad (7)$$

**Proof 2** Notice that when noise  $z$  is independent between dimensions, we have  $\nabla_\zeta \frac{\mu_z(\zeta)}{\mu_z(J \odot \zeta)} = 0$ . Then, we can obtain

$$\begin{aligned} \mathbb{E}(\hat{g}_c^{\text{BLR}}(x_0)) &= \mathbb{E}_{z,J} \left( -\frac{1}{q} f_c(x_0 + J \odot z) \nabla_z \ln \mu_z(J \odot z) \right) \\ &= \mathbb{E}_J \left( -\frac{1}{q} \int_{\mathbb{R}^d} f_c(x_0 + J \odot \zeta) (\nabla_\zeta \mu_z(J \odot \zeta)) \frac{\mu_z(\zeta)}{\mu_z(J \odot \zeta)} d\zeta \right) \\ &= \mathbb{E}_J \left( -\frac{1}{q} \int_{\mathbb{R}^d} f_c(x_0 + J \odot \zeta) \nabla_\zeta \mu_z(\zeta) d\zeta \right) \\ &= \mathbb{E}_J \left( \frac{1}{q} \int_{\mathbb{R}^d} (J \odot \nabla_x f_c(x)|_{x_0+J \odot \zeta}) \mu_z(\zeta) d\zeta \right) \\ &= \mathbb{E}_{z,J} \left( \frac{1}{q} J \odot g_c(x_0 + J \odot z) \right). \end{aligned} \quad (8)$$

While both the standard and blockwise likelihood ratio estimators are equivalent in terms of expected value, the implementation of blockwise computation introduces a significant advantage in variance reduction across each dimension, if the number of times that noise is added to that dimension is maintained. Formally, given the same number of copies for each input dimension in the gradient estimator computed in the standard manner in eq. (2) and the blockwise manner in eq. (6), we have

$$\text{Var}([\hat{g}_c^{\text{BLR}}(x_0)]_i) < \text{Var}([\hat{g}_c^{\text{LR}}(x_0)]_i), \quad \forall i \in \{1, \dots, d\}. \quad (9)$$

The proof and detailed analysis can be found in appendix B.2. This property supports us in reducing the gradient estimation variance for saliency map generation. We define our blockwise likelihood ratio method for saliency map generation in the same way as eq. (5) but leveraging  $\hat{g}_c^{\text{BLR}}$ .

### 3.4 THE SELECTION OF THE CLASS ACTIVATION $f_c$

In the previous discussion, a gradient-based saliency map harnesses the gradient of the class activation w.r.t. inputs,  $g_c(x_0) = \nabla_x f_c(x)|_{x=x_0}$ , to explain the decision of neural networks. While we have addressed the “undifferentiable” problem due to inaccessibility for black-box models by utilizing the likelihood principle mentioned above, there remains an open question: *what if we can’t access the class activation  $f_c$  of black-box models?* Here, we propose solutions for three possible scenes in real-world applications.

**Soft- & Hard-label decision.** When logits or confidence of models are given, we can directly use them to estimate the gradient of the loss function, *i.e.*, cross-entropy, to the inputs for saliency map generation in our framework. When only the predicted label is given, *i.e.*, the hard label, we can not directly compute the loss function for gradient estimation. To address this issue, we propose to compute the nearest path distance between the predicted class and ground truth in WorldNet as the surrogate criteria to evaluate the classification performance.

**Text-based decision.** (Multimodal-) language models output the texts for users’ queries, making it difficult to define the performance criteria for gradient estimation directly. To address this problem, we propose a prompt construction strategy. Specifically, when querying the decision of one image, we will additionally ask the model to respond to our answer by following the format of “The input image belongs to [MASK]”, where the mask is filled by taking one category of WordNet. Then, we will compute the nearest distance between two nodes, *i.e.*, one that comes from the model and the ground truth, for performance evaluation and gradient estimation.

## 4 EXPERIMENT

In our experiments, we set up three different tasks to verify the use of forward gradient in explaining model decisions: the saliency map generation for computer vision models, sentiment analysis for language models, and bias interpretation for vision-language models. For a vision-involved explanation, we leverage 5,000 fine-grained noisy blocks to estimate the forward gradient. For text explanation, we leverage 1,000–96,000 noisy copies of the original inputs to estimate the forward gradients. We use the auto-grad function provided by the PyTorch library to compute the vanilla gradient as the white-box explanation. All experiments are conducted in a single A6000 GPU with 48 GB VRAM.

### 4.1 SALIENCY MAP GENERATION ON VISUAL MODELS

We first evaluate the performance on the task of generating the saliency maps on the ImageNet-1K dataset (Deng et al., 2009). We compare the forward gradient with the vanilla gradient. The vanilla gradient is fetched by the auto-grad in the PyTorch library, while the forward gradient has no access to the model information except the input and output.

**Visualization.** We present the results in fig. 3. It can be observed that all saliency maps generated by the vanilla gradient are more susceptible to inherent noise in the image, indicating weaker interpretability in explaining model decisions. In contrast, the forward gradients naturally smooth out the noise by injecting perturbations, thereby providing better explainability. This can be evidenced by the main objects are significantly highlighted on the saliency maps generated by the forward gradient compared with the vanilla gradient. Also, the saliency maps generated for different models by the forward gradients also depict large difference. While only one fish is highlighted in explaining the AlexNet, nearly all fishes are highlighted in explaining the ResNet-50. This demonstrates the good scalability of the forward gradient in explaining various architectures.

**Deletion&Insertion.** Beyond the visualization results, we also use the deletion and insertion game to demonstrate the effectiveness of the forward gradient in explaining the models. In the deletion task, we remove the most important features and observe the drop in the model’s confidence, indicating the crucial role of these features—*where a lower deletion score is better*. In the insertion task, we reintroduce the important features into a blank image and track the increase in the model’s confidence, with a *higher insertion score indicating better performance*.

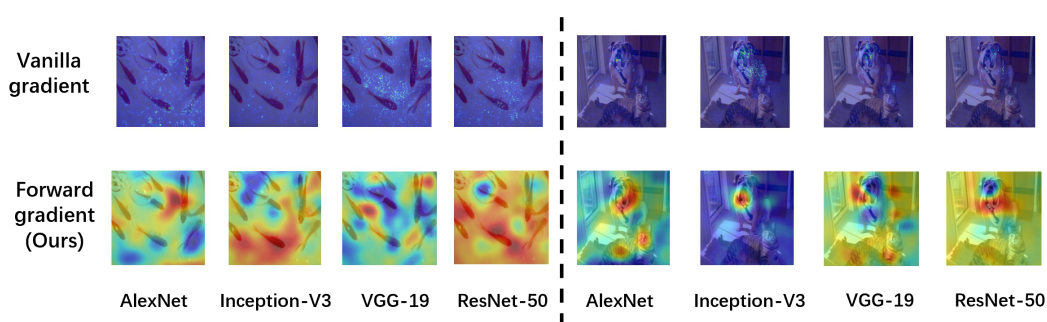


Figure 3: Visualization comparison between the saliency maps of four different models generated by vanilla gradient and forward gradient (ours). The lightness indicates the strength of model’s interest.

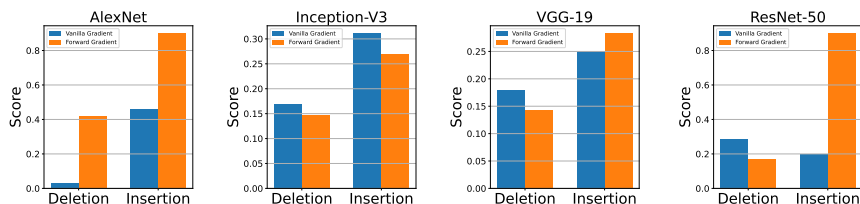


Figure 4: Quantitative comparison between the vanilla gradient and forward gradient in explaining model decisions. We play the deletion ( $\downarrow$ )&insertion ( $\uparrow$ ) game on four models to compare the performance.

The quantitative comparison results are shown in fig. 4. In general, though without any access to the model architecture, the forward gradient outperforms the vanilla gradient in explaining three of four models (AlexNet $\times$ , Inception-V3 $\checkmark$ , VGG-19 $\checkmark$ , ResNet-50 $\checkmark$ ), as indicated by the scores. Dive into the results of AlexNet, we can find the cause of failure from fig. 3, where the highlighted area has some position shift to some degree with the main object, which could be due to the estimation variance in gradient estimation. On the other hand, Although the vanilla gradient achieves a better score for AlexNet, the highlighted areas in fig. 3 are not sufficiently explainable. This could be attributed to adversarial perturbations rather than the semantic areas indicating the model’s interests.

## 4.2 EXPLAINING TEXT CLASSIFICATION MODELS ON SENTIMENT ANALYSIS

In this study, we delve deeper into the interpretability of the forward gradient method in natural language processing (NLP) tasks. Specifically, we train an attention-based bidirectional LSTM model for text classification using the Stanford Sentiment Treebank-2 (SST-2) dataset (Socher et al., 2013). The model incorporates pre-trained GloVe embeddings for word representation and achieves an overall accuracy of 82.4% on the test set. We focus on explainability by analyzing gradients with respect to the GloVe word embeddings rather than the word indices. Consistent with previous experimental setups, we compare the forward gradient method’s performance with the vanilla gradient, computed using PyTorch’s Auto-grad feature, to highlight differences in interpretability and insights.

As shown in fig. 5, we present the analysis of four sentences. The word importance identified by the vanilla gradient and forward gradient shows high consistency, differing primarily in magnitude. For example, in the first sentence, both gradients indicate that the word ‘must’ plays a crucial role in the final decision. However, a notable difference is observed in their treatment of the word ‘believed.’ While the vanilla gradient assigns it a less significant role, the forward gradient analysis suggests its considerable importance in classifying the sentence as positive. We argue that this observation aligns more closely with human understanding of natural language processing tasks.



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

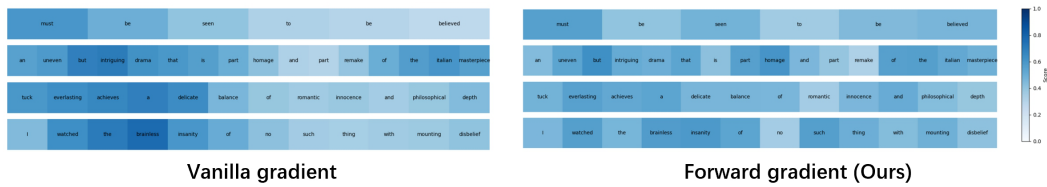


Figure 5: Comparison of vanilla gradient and forward gradient in explaining the text models from the word embedding level. A darker color indicates a higher importance approved by the gradient.

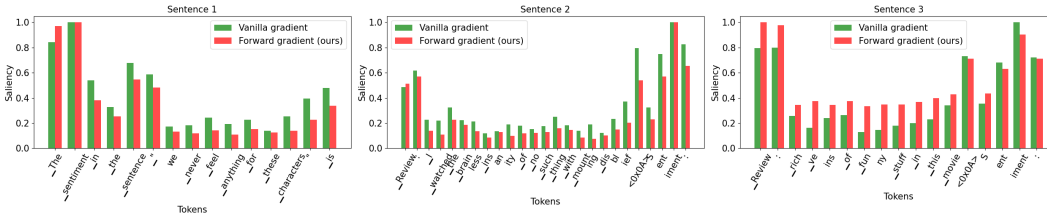


Figure 6: Comparison of calculated saliency scores and estimated saliency scores on LLM. Our method shows the same trend of variation as the calculated saliency scores.

### 4.3 EXPLAINING THE PHI-3

We evaluated our method on Phi-3 mini (Abdin et al., 2024). We prompt the LLM to determine the sentiment of three movie reviews obtained from the SST-2 dataset (Socher et al., 2013). The saliency map estimated by our method aligns with the overall trend observed in the saliency map generated using the vanilla gradient method. However, the noise introduced in our estimation process makes it less contrastive, meaning there is a less pronounced difference in saliency score between tokens with high and low saliency scores. Nevertheless, this effect can be mitigated by subtracting a bias from the saliency scores and rescaling the map. As shown in fig. 6, our estimated saliency map closely aligns with the saliency map generated using the vanilla gradient method.

### 4.4 BIAS INTERPRETATION FOR VISION-LANGUAGE MODELS

We move on to a more challenging task, using the forward gradient to explain vision-language models on a biased dataset. Specifically, we focus on the CLIP model, employing the ViT-B/32 architecture as the feature encoder. The dataset used for this study is Hard ImageNet, which contains many spurious correlations that can mislead the model’s decision-making. In the absence of explicit feedback signals from downstream tasks, we use the cosine similarity between the target class label and the image features to estimate the forward gradient.

We select five classes, including the howler monkey, seat belt, space bar, dog sled, and balance beam, and present the interpretation results in fig. 7. We provide the text embeddings with the ground-truth labels and compute the forward gradient as previously described to identify which parts contribute to the model’s decision. From the results, we observe that the CLIP encoder performs well without bias

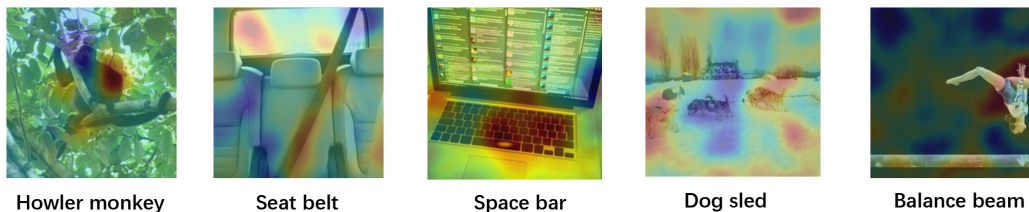


Figure 7: Leverage the forward gradient to explain the CLIP on Hard-ImageNet.

for the 'howler monkey' and 'space bar' classes. However, it is influenced by spurious correlations in the remaining three classes. For instance, while CLIP correctly identifies 'seat belt,' it also associates the concept with the outside sky. In the case of the 'dog sled' class, CLIP links the concept primarily with the dog, neglecting the sled. A similar phenomenon occurs with the 'balance beam' class, where the model focuses mainly on the athlete, causing the balance beam to be overlooked, resulting in misalignment between the visual and text embeddings. On the one hand, these results reveal that CLIP is susceptible to spurious correlations, even when trained on a large-scale, multi-modal dataset. On the other hand, we can see that the forward gradient serves as a powerful tool for analyzing and understanding failures in black-box models in practical applications.

#### 4.5 ABLATION STUDY

we provide an ablation study on the involved gradient estimation part. On text models, we use the cosine similarity between the estimated gradient and vanilla gradient to study the impact of the number of word embedding copies on estimation accuracy and use the insertion score to study the influence of block-wise computation in explaining vision models.

Table 1: Gradient estimation accuracy by varying the number of copies. Table 2: Insertion score using the blockwise technique with different numbers of blocks.

# copies	10	20	30	40	50	100	# Blocks ( $\times 10^3$ )	1	2	3	4	5
Similarity (%)	13.1	19.7	23.3	24.2	29.6	36.3	Insertion score (%)	51.33	59.54	60.08	60.26	57.20

**The gradient estimation accuracy using likelihood ratio method.** On a text model, we vary the number of copies for word embedding to study the performance. As shown in Tab. 1, we can see that the accuracy of the gradient estimation is improved when the number of word embedding copies is increased. While it only has a similarity of 13.1% for the estimated gradient with the ground-truth, it can be improved to 36.3% when taking 100 copies into the estimation process. Thus, there is a trade-off between computation requirement and robust interpretation of model decision.

**The blockwise computation for variance reduction.** We report the ablation study of the use of blockwise computation in Tab. 2. We vary the number of small blocks in gradient estimation from 1, 000 to 10, 000. The size of blocks is set to 10% of the original image size. With increasing the number of blocks, the insertion score can be improved, indicating a better capacity in explaining models. However, with large enough number of blocks, the gradient estimation can get stuck in a local optima due to the balance between estimation redundancy and variance. Specifically, it achieves the peak performance when selecting 4, 000 as the number of blocks in computation, while a larger number copies 5, 000 causes a performance degradation of 3.06%.

## 5 CONCLUSION

The likelihood ratio method makes it possible to train neural networks with only forward passes. Inspired by this, we design a pipeline to employ the likelihood ratio method in gradient estimation, particularly for interpreting model decisions in black-box scenarios. To address the issue of large gradient estimation variance, we propose a blockwise computation technique that balances computational cost with task performance. We explore the potential of the forward gradient in explaining models across three tasks: saliency map generation, sentiment analysis on text models, and interpreting spurious correlations learned by vision-language models. The experimental results demonstrate the effectiveness and scalability of the forward gradient in explaining black-box models. This study paves the way for future research in explaining black-box models.

## LIMITATIONS

While the forward gradient offers a method to estimate gradients without relying on the chain rule—facilitating the explanation of black-box models based on gradient information—it is computationally expensive. For larger models, this approach may require thousands of iterations to achieve an acceptable estimation. Developing effective variance reduction methods is crucial not only for training models without backpropagation (BP) but also for improving the explainability of black-box models.

## REFERENCES

- 540  
541  
542 Marah I Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany  
543 Hassan Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, Alon Benhaim,  
544 Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio César Teodoro Mendes, Weizhu  
545 Chen, Vishrav Chaudhary, Parul Chopra, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon,  
546 Ronen Eldan, Dan Iter, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Rus-  
547 sell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis,  
548 Dongwoo Kim, Mahmoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yanzhi Li, Chen  
549 Liang, Weishung Liu, Xihui (Eric) Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi,  
550 Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant,  
551 Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olli Saarikivi, Amin Saied, Adil  
552 Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Xia Song, Olatunji Ruwase,  
553 Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael Wyatt, Can Xu, Jiahang Xu, Son-  
554 ali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang,  
555 Li Lyna Zhang, Yi Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable  
556 language model locally on your phone. Technical Report MSR-TR-2024-12, Microsoft, Au-  
557 gust 2024. URL [https://www.microsoft.com/en-us/research/publication/  
phi-3-technical-report-a-highly-capable-language-model-locally-on-your-phone/](https://www.microsoft.com/en-us/research/publication/phi-3-technical-report-a-highly-capable-language-model-locally-on-your-phone/).
- 558 Reduan Achtibat, Maximilian Dreyer, Ilona Eisenbraun, Sebastian Bosse, Thomas Wiegand, Wojciech  
559 Samek, and Sebastian Lapuschkin. From attribution maps to human-understandable explanations  
560 through concept relevance propagation. *Nature Machine Intelligence*, 5(9):1006–1019, 2023.
- 561 Soufiane Belharbi, Aydin Sarraf, Marco Pedersoli, Ismail Ben Ayed, Luke McCaffrey, and Eric  
562 Granger. F-cam: Full resolution class activation maps via guided parametric upscaling. In  
563 *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3490–  
564 3499, 2022.
- 565 Francesco Bodria, Fosca Giannotti, Riccardo Guidotti, Francesca Naretto, Dino Pedreschi, and  
566 Salvatore Rinzivillo. Benchmarking and survey of explanation methods for black box models.  
567 *Data Mining and Knowledge Discovery*, pp. 1–60, 2023.
- 568 Markus Borg, Cristofer Englund, Krzysztof Wnuk, Boris Duran, Christoffer Levandowski, Shenjian  
569 Gao, Yanwen Tan, Henrik Kaijser, Henrik Lönn, and Jonas Törnqvist. Safely entering the deep:  
570 A review of verification and validation for machine learning and a challenge elicitation in the  
571 automotive industry. *arXiv preprint arXiv:1812.05389*, 2018.
- 572 Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-  
573 cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *2018*  
574 *IEEE winter conference on applications of computer vision (WACV)*, pp. 839–847. IEEE, 2018.
- 575 Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order  
576 optimization based black-box attacks to deep neural networks without training substitute models.  
577 In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pp. 15–26, 2017.
- 578 Xiangyi Chen, Sijia Liu, Kaidi Xu, Xingguo Li, Xue Lin, Mingyi Hong, and David Cox. Zo-  
579 adamm: Zeroth-order adaptive momentum method for black-box optimization. *Advances in neural*  
580 *information processing systems*, 32, 2019.
- 581 Jeffrey De Fauw, Joseph R Ledsam, Bernardino Romera-Paredes, Stanislav Nikolov, Nenad Tomasev,  
582 Sam Blackwell, Harry Askham, Xavier Glorot, Brendan O’Donoghue, Daniel Visentin, et al.  
583 Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature medicine*,  
24(9):1342–1350, 2018.
- 584 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
585 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
586 pp. 248–255. Ieee, 2009.
- 587 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
588 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
589 pp. 248–255. Ieee, 2009.
- 590 Gabriel Erion, Joseph D Janizek, Pascal Sturmfels, Scott M Lundberg, and Su-In Lee. Improving  
591 performance of deep learning models with axiomatic attribution priors and expected gradients.  
592 *Nature machine intelligence*, 3(7):620–631, 2021.
- 593

- 594 Abir Fathallah, Lotfi Abdi, and Ali Douik. Facial expression recognition via deep learning. In *2017*  
595 *IEEE/ACS 14th International Conference on Computer Systems and Applications (AICCSA)*, pp.  
596 745–750. IEEE, 2017.
- 597 Daniel Golovin, John Karro, Greg Kochanski, Chansoo Lee, Xingyou Song, and Qiuyi Zhang. Gradi-  
598 entless descent: High-dimensional zeroth-order optimization. *arXiv preprint arXiv:1911.06317*,  
599 2019.
- 600 Oleg Granichin and Natalia Amelina. Simultaneous perturbation stochastic approximation for  
601 tracking under unknown but bounded disturbances. *IEEE Transactions on Automatic Control*, 60  
602 (6):1653–1658, 2014.
- 603 Sorin Grigorescu, Bogdan Trasnea, Tiberiu Cocias, and Gigel Macesanu. A survey of deep learning  
604 techniques for autonomous driving. *Journal of Field Robotics*, 37(3):362–386, 2020.
- 605 Brian Hu, Paul Tunison, Brandon Richard Webster, and Anthony Hoogs. Xaitk-saliency: An open  
606 source explainable ai toolkit for saliency. *Proceedings of the AAAI Conference on Artificial*  
607 *Intelligence*, 37(13):15760–15766, Sep. 2023. doi: 10.1609/aaai.v37i13.26871. URL <https://ojs.aaai.org/index.php/AAAI/article/view/26871>.
- 608 Yu Huang and Yue Chen. Autonomous driving with deep learning: A survey of state-of-art technolo-  
609 gies. *arXiv preprint arXiv:2006.06091*, 2020.
- 610 Jinyang Jiang, Zeliang Zhang, Chenliang Xu, Zhaofei Yu, and Yijie Peng. One forward is enough for  
611 neural network training via likelihood ratio method. In *The Twelfth International Conference on*  
612 *Learning Representations*, 2023.
- 613 Peng-Tao Jiang, Chang-Bin Zhang, Qibin Hou, Ming-Ming Cheng, and Yunchao Wei. Layercam: Ex-  
614 ploring hierarchical class activation maps for localization. *IEEE Transactions on Image Processing*,  
615 30:5875–5888, 2021.
- 616 Andrei Kapishnikov, Subhashini Venugopalan, Besim Avci, Ben Wedin, Michael Terry, and Tolga  
617 Bolukbasi. Guided integrated gradients: An adaptive path method for removing noise. In  
618 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5050–  
619 5058, 2021.
- 620 Daniel S Kermany, Michael Goldbaum, Wenjia Cai, Carolina CS Valentim, Huiying Liang, Sally L  
621 Baxter, Alex McKeown, Ge Yang, Xiaokang Wu, Fangbing Yan, et al. Identifying medical  
622 diagnoses and treatable diseases by image-based deep learning. *cell*, 172(5):1122–1131, 2018.
- 623 Saeed Khorram, Tyler Lawson, and Li Fuxin. igos++ integrated gradient optimized saliency by  
624 bilateral perturbations. In *Proceedings of the Conference on Health, Inference, and Learning*, pp.  
625 174–182, 2021.
- 626 David Kozak, Cesare Molinari, Lorenzo Rosasco, Luis Tenorio, and Silvia Villa. Zeroth-order  
627 optimization with orthogonal random directions. *Mathematical Programming*, 199(1):1179–1219,  
628 2023.
- 629 Yi-Shan Lin, Wen-Chuan Lee, and Z Berkay Celik. What do you see? evaluation of explainable  
630 artificial intelligence (xai) interpretability through neural backdoors. In *Proceedings of the 27th*  
631 *ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 1027–1035, 2021.
- 632 Pantelis Linardatos, Vasilis Papastefanopoulos, and Sotiris Kotsiantis. Explainable ai: A review of  
633 machine learning interpretability methods. *Entropy*, 23(1):18, 2020.
- 634 Joe Lorentz, Thomas Hartmann, Assaad Moawad, Francois Fouquet, and Djamila Aouada. Explaining  
635 defect detection with saliency maps. In *International Conference on Industrial, Engineering and*  
636 *Other Applications of Applied Intelligent Systems*, pp. 506–518. Springer, 2021.
- 637 Daniel D Lundstrom, Tianjian Huang, and Meisam Razaviyayn. A rigorous study of integrated  
638 gradients method and extensions to internal neuron attributions. In *International Conference on*  
639 *Machine Learning*, pp. 14485–14508. PMLR, 2022.

- 648 John L Maryak and Daniel C Chin. Global random optimization by simultaneous perturbation  
649 stochastic approximation. In *Proceedings of the 2001 American control conference.(Cat. No.*  
650 *01CH37148)*, volume 2, pp. 756–762. IEEE, 2001.
- 651
- 652 Mostafa Mehdipour Ghazi and Hazim Kemal Ekenel. A comprehensive analysis of deep learning  
653 based representation for face recognition. In *Proceedings of the IEEE conference on computer*  
654 *vision and pattern recognition workshops*, pp. 34–41, 2016.
- 655
- 656 Wafa Mellouk and Wahida Handouzi. Facial emotion recognition using deep learning: review and  
657 insights. *Procedia Computer Science*, 175:689–694, 2020.
- 658
- 659 Ramaravind K Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning classifiers  
660 through diverse counterfactual explanations. In *Proceedings of the 2020 conference on fairness,*  
661 *accountability, and transparency*, pp. 607–617, 2020.
- 662
- 663 Sajjad Mozaffari, Omar Y Al-Jarrah, Mehrdad Dianati, Paul Jennings, and Alexandros Mouzakitis.  
664 Deep learning-based vehicle behavior prediction for autonomous driving applications: A review.  
665 *IEEE Transactions on Intelligent Transportation Systems*, 23(1):33–47, 2020.
- 666
- 667 Daniel Omeiza, Skyler Speakman, Celia Cintas, and Komminist Weldermariam. Smooth grad-cam++:  
668 An enhanced inference level visualization technique for deep convolutional neural network models.  
669 *arXiv preprint arXiv:1908.01224*, 2019.
- 670
- 671 OpenAI. Introducing chatgpt-3: The largest openai model yet. [https://openai.com/blog/  
chatgpt-3/](https://openai.com/blog/chatgpt-3/), 2020.
- 672
- 673 Yijie Peng, Li Xiao, Bernd Heidergott, L Jeff Hong, and Henry Lam. A new likelihood ratio method  
674 for training artificial neural networks. *INFORMS Journal on Computing*, 34(1):638–655, 2022.
- 675
- 676 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
677 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
678 models from natural language supervision. In *International conference on machine learning*, pp.  
8748–8763. PMLR, 2021.
- 679
- 680 Marco Rando, Cesare Molinari, Lorenzo Rosasco, and Silvia Villa. An optimal structured zeroth-  
681 order algorithm for non-smooth optimization. *Advances in Neural Information Processing Systems*,  
682 36, 2024.
- 683
- 684 David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by  
685 back-propagating errors. *nature*, 323(6088):533–536, 1986.
- 686
- 687 Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a  
688 scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- 689
- 690 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh,  
691 and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local-  
692 ization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626,  
2017.
- 693
- 694 Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:  
695 Visualising image classification models and saliency maps. In Yoshua Bengio and Yann LeCun  
696 (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada,*  
697 *April 14-16, 2014, Workshop Track Proceedings*, 2014.
- 698
- 699 Dylan Slack, Anna Hilgard, Himabindu Lakkaraju, and Sameer Singh. Counterfactual explanations  
700 can be manipulated. *Advances in neural information processing systems*, 34:62–75, 2021.
- 701
- 702 Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad:  
removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.

- 702 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and  
703 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank.  
704 In David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard (eds.),  
705 *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp.  
706 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.  
707 URL <https://aclanthology.org/D13-1170>.
- 708 James C Spall. A one-measurement form of simultaneous perturbation stochastic approximation.  
709 *Automatica*, 33(1):109–112, 1997.
- 710 James C Spall. Adaptive stochastic approximation by the simultaneous perturbation method. *IEEE*  
711 *transactions on automatic control*, 45(10):1839–1853, 2000.
- 712 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Gradients of counterfactuals. *arXiv preprint*  
713 *arXiv:1611.02639*, 2016.
- 714 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In  
715 *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.
- 716 Florian Tambon, Gabriel Laberge, Le An, Amin Nikanjam, Paulina Stevia Nouwou Mindom, Yann  
717 Pequignot, Foutse Khomh, Giulio Antoniol, Ettore Merlo, and François Laviolette. How to certify  
718 machine learning based safety-critical systems? a systematic literature review. *Automated Software*  
719 *Engineering*, 29(2):38, 2022.
- 720 Erico Tjoa and Guan Cuntai. Quantifying explainability of saliency methods in deep neural networks  
721 with a synthetic dataset. *IEEE Transactions on Artificial Intelligence*, 2022.
- 722 Jeroen Van der Laak, Geert Litjens, and Francesco Ciompi. Deep learning in histopathology: the  
723 path to the clinic. *Nature medicine*, 27(5):775–784, 2021.
- 724 Anirudh Vemula, Wen Sun, and J Bagnell. Contrasting exploration in parameter and action space:  
725 A zeroth-order optimization perspective. In *The 22nd International Conference on Artificial*  
726 *Intelligence and Statistics*, pp. 2926–2935. PMLR, 2019.
- 727 Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and  
728 Xia Hu. Score-cam: Score-weighted visual explanations for convolutional neural networks. In  
729 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*,  
730 pp. 24–25, 2020.
- 731 Yining Wang, Simon Du, Sivaraman Balakrishnan, and Aarti Singh. Stochastic zeroth-order opti-  
732 mization in high dimensions. In *International conference on artificial intelligence and statistics*,  
733 pp. 1356–1365. PMLR, 2018.
- 734 Peiyu Yang, Naveed Akhtar, Zeyi Wen, and Ajmal Mian. Local path integration for attribution.  
735 *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(3):3173–3180, Jun. 2023. doi: 10.  
736 1609/aaai.v37i3.25422. URL [https://ojs.aaai.org/index.php/AAAI/article/  
737 view/25422](https://ojs.aaai.org/index.php/AAAI/article/view/25422).
- 738 Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep  
739 features for discriminative localization. In *Proceedings of the IEEE conference on computer vision*  
740 *and pattern recognition*, pp. 2921–2929, 2016.
- 741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A ALGORITHM

We present the algorithm details in 1. To explain the decision of the black-box model  $f$  on the inputs  $x$ , we first generate  $n$  copies as  $\hat{x}$ . Then, we sample  $n$  independent random noise from the normal distribution to perturb the copies. Next, with the target label  $y$ , we compute the loss function value for each noisy copy and get the proxy gradient. Last, we compute the average of proxy gradients, get the estimated gradient, and employ the white-box saliency map generation method to craft  $s$ .

---

### Algorithm 1: Likelihood ratio method for saliency map generation

---

**input** : Image  $x$  with label  $y$ , black model  $f$ , loss function  $\mathcal{L}(\cdot, \cdot)$ , number of copies  $n$ , while-box saliency map generation method  $S(\cdot)$

**output** : The saliency map  $s$  to interpret the decision of  $f$  on  $x$

```

1  $\hat{x} \leftarrow$  generate  $n$  copies of  $x$ ;
2  $\delta \leftarrow$  sample  $n$  i.i.d. noise from the normal distribution;
3  $\hat{x} \leftarrow \hat{x} + \delta$ ;
  // perturb the inputs
4  $\hat{l} \leftarrow \mathcal{L}(f(\hat{x}), y)$ ;
5  $\hat{g} \leftarrow \hat{l} \cdot \delta$ ;
  // compute the proxy gradient
6  $g \leftarrow \frac{1}{n} \sum_{i=1}^n \hat{g}_i$ ;
  // compute the estimated gradient
7  $s \leftarrow S(g)$ 

```

---

## B THEORY ANALYSIS

### B.1 PROOF OF THE COROLLARY

We provide the proof of the corollary we present in 3.3. Mathematically, the corollary is formalized as

**Corollary 1** Assume the model’s gradient  $g_c$  is bounded and let the clockwise likelihood ratio gradient estimator  $\hat{g}_c^{BLR}$  as defined in 6. Then for any input  $x_0$ , we have

$$\lim_{\sigma \rightarrow 0^+} \mathbb{E}(\hat{g}_c^{BLR}(x_0)) = g_c(x_0) \quad (10)$$

**Proof 3** Consider arbitrary Gaussian noise r.v. sequence  $\{z_i\}_1^\infty$  that satisfies  $\mathbb{E}(z_i) = 0$  and  $\lim_{i \rightarrow \infty} \Sigma_i = 0$ .

810 *It follows that*

$$\begin{aligned}
811 & z_i \xrightarrow{L_2} 0 \\
812 & \Rightarrow z_i \xrightarrow{P} 0 \\
813 & \Rightarrow J \odot z_i \xrightarrow{P} 0 \\
814 & \Rightarrow x_0 + J \odot z_i \xrightarrow{P} x_0 \\
815 & \Rightarrow g_c(x_0 + J \odot z_i) \xrightarrow{P} g_c(x_0) \\
816 & \Rightarrow \frac{1}{q} J \odot g_c(x_0 + J \odot z_i) \xrightarrow{P} \frac{1}{q} J \odot g_c(x_0) \\
817 & \Rightarrow \frac{1}{q} J \odot g_c(x_0 + J \odot z_i) \xrightarrow{L_1} \frac{1}{q} J \odot g_c(x_0) \\
818 & \Rightarrow \lim_{i \rightarrow \infty} \mathbb{E}_{z_i, J} \left( \frac{1}{q} J \odot g_c(x_0 + J \odot z_i) \right) = \mathbb{E}_{z_i, J} \left( \frac{1}{q} J \odot g_c(x_0) \right) \\
819 & \Rightarrow \lim_{i \rightarrow \infty} \mathbb{E}_{z_i, J} \left( \frac{1}{q} J \odot g_c(x_0 + J \odot z_i) \right) = g_c(x_0)
\end{aligned} \tag{11}$$

820 *Recall that by 2, we have  $\mathbb{E}(\hat{g}_c^{BLR}) = \mathbb{E}_{z, J}(\frac{1}{q} J \odot g_c(x_0 + J \odot z))$ . This, together with the arbitrariness*

821 *of  $\{z_i\}_1^\infty$  as long as  $\lim_{i \rightarrow \infty} \Sigma_i = 0$ , indicates that*

$$822 \lim_{\sigma \rightarrow 0^+} \mathbb{E}(\hat{g}_c^{BLR}(x_0)) = g_c(x) \tag{12}$$

## 823 B.2 ANALYSIS OF VARIANCE REDUCTION FROM THE BLOCKWISE COMPUTATION

824 In this discussion, we consider the Gaussian noise  $z \sim N(0, \sigma^2 \mathbb{I})$ . Other situations are similar. Then,

825 we have simplified estimators as follows:

$$826 \hat{g}_c^{LR}(x_0) = \frac{1}{n^{LR} \sigma^2} \sum_i^{n^{LR}} f_c(x_0 + z_i) z_i \tag{13}$$

$$827 \hat{g}_c^{BLR}(x_0) = \frac{1}{n^{BLR} q \sigma^2} \sum_{i=1}^{n^{BLR}} (f_c(x_0 + J_i \odot z_i) J_i \odot z_i). \tag{14}$$

828 Given the same number of copies for each input dimension in the gradient estimator computed

829 in the standard manner, *i.e.*,  $n^{LR} = n^{BLR} q$ , we need to prove  $\text{Var}([\hat{g}_c^{LR}(x_0)]_i) < \text{Var}([\hat{g}_c^{BLR}(x_0)]_i)$ ,

830  $\forall i \in \{1, \dots, d\}$ . To prove this, we only need to prove  $\text{Var}([f_c(x_0 + z)]_i) < \text{Var}([f_c(x_0 + J \odot z)]_i \mid [J]_i = 1)$ . Using multivariate Taylor expansion, when  $\sigma$  is small, we can say

$$831 f_c(x_0 + z) = f_c(x_0) + \nabla_x f|_{x=x_0} \cdot z. \tag{15}$$

832 It follows that

$$833 \mathbb{E}(f_c(x_0 + z)z) = 0 + \nabla_x f|_{x=x_0} \sigma^2 \mathbb{I} = \sigma^2 \nabla_x f|_{x=x_0}. \tag{16}$$

834 Then, we have

$$835 \text{Var}([f_c(x_0 + z)]_i) \tag{17}$$

$$836 = \mathbb{E}([f_c(x_0 + z)]_i^2) - \sigma^4 [\nabla_x f|_{x=x_0}]_i^2 \tag{18}$$

$$837 = \mathbb{E}(f_c^2(x_0 + z) [z]_i^2) - \sigma^4 [\nabla_x f|_{x=x_0}]_i^2 \tag{19}$$

$$838 = \mathbb{E}((f_c(x_0) + \nabla_x f|_{x=x_0} \cdot z)^2 [z]_i^2) - \sigma^4 [\nabla_x f|_{x=x_0}]_i^2 \tag{20}$$

$$839 = f_c^2(x_0) \sigma^2 + 2\sigma^4 [\nabla_x f|_{x=x_0}]_i^2 + \sigma^4 \sum_{k=1, k \neq i}^d [\nabla_x f|_{x=x_0}]_k^2 \tag{21}$$

840 Similarly, we have

$$841 \mathbb{E}(f_c(x_0 + J \odot z) J \odot z \mid [J]_i = 1) = \sigma^2 [\nabla_x f|_{x=x_0}]_i \tag{22}$$



864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

and then

$$\text{Var}([f_c(x_0 + J \odot z)J \odot z]_i \mid [J]_i = 1) \tag{23}$$

$$= \mathbb{E}([f_c(x_0 + J \odot z)J \odot z]_i^2 \mid [J]_i = 1) - \sigma^4 [\nabla_x f|_{x=x_0}]_i^2 \tag{24}$$

$$\approx f_c^2(x_0)\sigma^2 + 2\sigma^4 [\nabla_x f|_{x=x_0}]_i^2 + q\sigma^4 \sum_{k=1, k \neq i}^d [\nabla_x f|_{x=x_0}]_k^2 \tag{25}$$

$$< \text{Var}([f_c(x_0 + z)z]_i), \tag{26}$$

if  $\exists k \in \{1, \dots, d\}, k \neq i$ , such that  $[\nabla_x f|_{x=x_0}]_k \neq 0$ .