
Gradient strikes back: How filtering out high frequencies improves explanations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent years have witnessed an explosion in the development of novel prediction-
2 based attribution methods, which have slowly been supplanting older gradient-
3 based methods to explain the decisions of deep neural networks. However, it is still
4 not clear why prediction-based methods outperform gradient-based ones. Here, we
5 start with an empirical observation: these two approaches yield attribution maps
6 with very different power spectra, with gradient-based methods revealing more
7 high-frequency content than prediction-based methods. This observation raises
8 multiple questions: What is the source of this high-frequency information, and does
9 it truly reflect decisions made by the system? Lastly, why would the absence of
10 high-frequency information in prediction-based methods yield better explainability
11 scores along multiple metrics? We analyze the gradient of three representative
12 visual classification models and observe that it contains noisy information emanat-
13 ing from high-frequencies. Furthermore, our analysis reveals that the operations
14 used in Convolutional Neural Networks (CNNs) for downsampling appear to be a
15 significant source of this high-frequency content – suggesting aliasing as a possible
16 underlying basis. We then apply an optimal low-pass filter for attribution maps and
17 demonstrate that it improves gradient-based attribution methods. We show that (i)
18 removing high-frequency noise yields significant improvements in the explainabil-
19 ity scores obtained with gradient-based methods across multiple models – leading
20 to (ii) a novel ranking of state-of-the-art methods with gradient-based methods
21 at the top. We believe that our results will spur renewed interest in simpler and
22 computationally more efficient gradient-based methods for explainability.

23 1 Introduction

24 Explaining and interpreting the decision of AI architectures is an important area of research towards
25 enabling the development of more interpretable models. Explainability methods (XAI) aim to provide
26 insights into the strategies used by models to arrive at their decision. This is expected to lead to the
27 development of better models that are more accurate, robust, and better aligned with humans.

28 One of the first attribution methods proposed, “Saliency” [1], consists of back-propagating a model’s
29 decision back to an input image to highlight areas that most affected the final decision. The method
30 remains relatively simple and computationally efficient, but it is also known to be noisy and to lead to
31 attribution maps that are often hard to interpret. Multiple methods have been proposed since to try to
32 improve on these limitations. These methods fall broadly into two main classes. (i) Gradient-based
33 methods extend Saliency [1] by smoothing the resulting attribution maps [2–8]. However, these
34 so-called white-box methods require access to all the model’s components, which is not always

[†] The authors contributed equally.

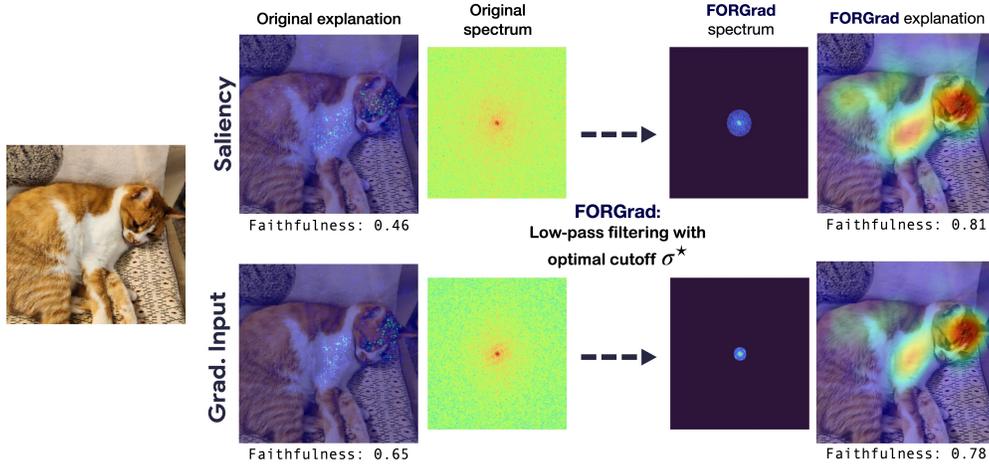


Figure 1: **Effect of FORGrad on gradient-based attribution methods.** We show that for an input image (left), the initial explanations from two gradient-based methods are plagued by noise as indicated by the high power in the high-frequency range of their respective spectra. Filtering the explanations with **FORGrad** yields improved explanations (right).

35 possible. Conversely, prediction-based methods, also called black-box methods [3, 9–11], alter the
 36 input of the model to produce an explanation based on the resulting change in the output. Those
 37 methods are computationally inefficient and are known to sometimes fail to capture all the diagnostic
 38 information, but they currently lead to the best fidelity scores across all explainability methods.
 39 Overall, there are dozens of attribution methods available but relatively little is understood about
 40 what makes certain methods more accurate than others.

41 Here, we start with the observation made across multiple studies [12, 4, 13] that the attribution maps
 42 derived with Saliency are very noisy. Generally, these maps highlight sparse pixel activations around
 43 a region of interest, and they are often hard to interpret. Because Saliency is simply the gradient of
 44 the score function with respect to the input, we suggest that the noise originates from the gradient
 45 itself: in other words, because the gradient is noisy, the explanation provided by Saliency is also
 46 noisy. To try to better understand the origin of this noise, we compare the Fourier power spectra of
 47 gradient-based methods (including Saliency) against prediction-based methods and observe that they
 48 differ quite markedly. We discern significant differences between the two classes of approaches, with
 49 gradient-based methods returning higher frequency content and prediction-based methods returning
 50 lower frequency content. In the remainder of this paper, we will show that:

- 51 • The gradient is indeed noisy, and this noise is especially present in the high-frequencies.
- 52 • We then look for the origin of these high frequencies in vision models. Our findings show
 53 that downsampling operations (via MaxPooling or strides) are the main sources of high
 54 frequencies, and training the model does not alleviate the issue.
- 55 • We then propose to repair Saliency – as well as other gradient-based methods – by introduc-
 56 ing **FORGrad** (FOurier Reparation of the Gradient). This method consists in estimating the
 57 optimal amount of high frequencies to remove per model to make gradient-based methods
 58 surpass the prediction-based family of attribution methods.

59 2 Related Work

60 **Attribution methods for black-box models** Various methods have been developed to compute
 61 importance scores for individual pixels or groups of pixels. For black-box (prediction-based) attribu-
 62 tion methods, the analytical form and potential internal states of the model are unknown. The first
 63 method, Occlusion [3], masks individual image regions, one at a time, using an occluding mask set to
 64 a baseline value. The corresponding prediction scores are assigned to all pixels within the occluded
 65 region, providing an easily interpretable explanation. However, occlusion fails to account for the
 66 joint (higher-order) interactions between multiple image regions. For instance, occluding two image
 67 regions individually may only have a minimal impact on the model’s prediction, such as removing a
 68 single eye or mouth component from a face. However, occluding these two regions together may lead

69 to a substantial change in the model’s prediction if these regions interact non-linearly, as expected
 70 in a deep neural network. Sobol [10], along with related methods such as LIME [11] and RISE [9],
 71 address this problem by randomly perturbing multiple regions of the input image simultaneously.
 72 Interestingly, recent studies, including RISE [9] and Sobol [10], have demonstrated that black-box
 73 attribution methods can rival and even surpass the commonly used white-box methods without relying
 74 on internal states.

75 **Attribution methods for white-box models** The gradient-based methods, that we propose to
 76 improve here, were first introduced in [14] and improved in [2–4]. They consist in explaining the
 77 decisions of a model by back-propagating the gradient from the output to the input, indicating which
 78 pixels affect the decision score the most. However, this family of methods is limited because they
 79 focus on the influence of individual pixels in an infinitesimal neighborhood in the input image. For
 80 instance, it has been shown that gradients often vanish when the prediction score to be explained is
 81 near the maximum value [6]. Integrated Gradient [6] and SmoothGrad [5] partially address this issue
 82 by accumulating gradients. Another family of attribution methods relies on the neural network’s
 83 activation, like CAM [7], which computes an attribution score based on a weighted sum of feature
 84 channel activities – right before the classification layer. GradCAM [8] extends CAM via the use of
 85 gradients, re-weighting each feature channel to take into account their importance for the predicted
 86 class. Nevertheless, the choice of the layer has a huge impact on the quality of the explanation. Our
 87 contribution proposes to overcome some of the mentioned issues by removing the noise present in
 88 the gradients in the form of high frequencies.

89 **Fourier analysis of vision models** Very little work has been proposed to analyze vision models
 90 and methods from a Fourier perspective. The closest, [15], used Fourier analysis to investigate the
 91 impact of DNNs optimization parameters and methods without a specific focus on vision.
 92 Additional work has focused on the analysis and development of adversarial attacks in the Fourier
 93 domain, [16, 17], while others [18–20] have proposed to defend against adversarial attacks by
 94 transforming the input image in the Fourier domain. Jo and Bengio [21] examined whether CNNs
 95 rely on high-level features by using Fourier-filtered images. None of the mentioned studies make a
 96 link between explainability and attribution methods with Fourier analysis.

97 3 Decomposing the gradient: An analysis of frequency content in attribution 98 methods

99 **Notations** We consider a general supervised learning setting, where a classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ maps
 100 images from an input space $\mathcal{X} \subseteq \mathbb{R}^{W \times H}$ to an output space $\mathcal{Y} \subseteq \mathbb{R}$. Let (x_1, \dots, x_N) be a set of
 101 images which contains N samples drawn from a probability distribution $\forall i \in \{1 \dots n\}, x_i \sim \mathcal{D}$.
 102 Moreover, we respectively denote \mathcal{F} and \mathcal{F}^{-1} the Discrete Fourier Transform (DFT) on $\mathbb{R}^{W \times H}$
 103 and its inverse. Therefore: $\forall x \in \mathcal{X}, \mathcal{F}(x) \in \mathbb{C}^{W \times H}$ and $(\mathcal{F}^{-1} \circ \mathcal{F})(x) = x$. Additionally, when we
 104 visualize the Fourier spectrum, we always shift the low-frequency components to the center of the
 105 spectrum. We recall that an attribution method is a function $\varphi : \mathcal{X} \rightarrow \mathbb{R}^{W \times H}$ that maps an input of
 106 interest to its corresponding importance scores $\varphi(x)$. Finally, we denote by $\varphi_\sigma(x)$ the attribution
 method where high frequencies have been filtered using a cutoff value of σ .

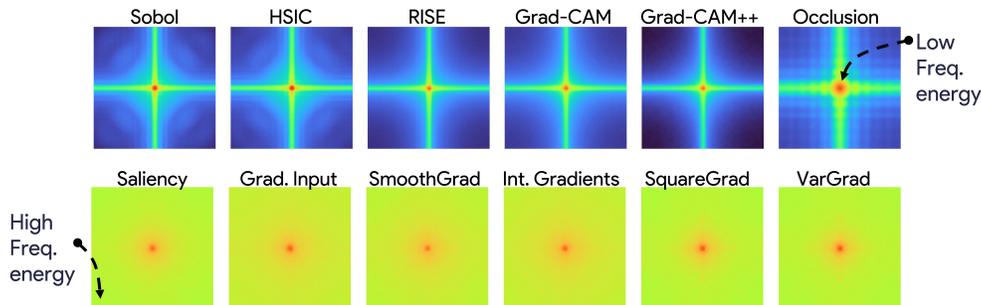


Figure 2: **Fourier footprint of attribution methods.** We show on the top row the Fourier spectrum of prediction-based attribution methods and of the gradient-based methods on the bottom row, computed with a ResNet50. The two families can be distinguished by methods but also by their signature in the Fourier domain. The former method has magnitudes largely concentrated in the low frequencies, while the latter is more spread out: it features non-trivial magnitudes almost everywhere, including in high frequencies.

108 **3.1 Different signatures for different categories of methods**

109 In this work, we analyze the Fourier signature of several attribution methods. To do so, we compute
 110 the feature map $\varphi(x)$ for most existing attribution methods on representative models of the literature
 111 (ResNet50 in Figure 2). From these importance maps, we extract the corresponding amplitude of the
 112 Fourier spectrum, $|(\mathcal{F} \circ \varphi)(x)|$. In Figure 2, we show the average power spectra, over 1,000 images,
 113 for an array of methods. Upon visual inspection, it is obvious that certain methods tend to emphasize
 114 higher frequencies in their explanations, while others concentrate on lower frequencies. Interestingly,
 115 these differences can be traced to the class of methods: Black-box methods, which do not rely
 116 on gradients, exhibit frequency footprints dominated by very low frequencies, whereas white-box
 117 methods exhibit footprints that extend into higher frequencies. To quantify our observations, we
 118 employ two metrics to measure the complexity of the attribution maps. The first metric employs
 119 a Laplacian-based operator [22, 23] that evaluates the presence of high frequencies in images by
 120 analyzing their second derivative. The second metric involves measuring the file size of the image
 121 after undergoing lossless compression [24, 25], which we refer to as ‘‘High-frequency content’’
 122 throughout this study (as it can be seen as a loose approximation of Kolmogorov complexity).
 123 Both metrics validate our visual observations,
 124 as depicted in Figure 3 (see Laplace quantity
 125 in appendix). It is evident that black-box meth-
 126 ods (shown in dark in the figure) exhibit fewer
 127 high frequencies compared to white-box meth-
 128 ods. This observation provides valuable insight
 129 into where these methods extract information
 130 from the model to compute their explanations.

131 **3.2 High-frequencies
 132 are just noise in the gradient**

133 Naturally, gradient-based methods will be sub-
 134 ject to the characteristics of the gradient itself.
 135 Consequently, when the gradient is subject to
 136 noise, the resulting explanation provided by
 137 such methods becomes similarly noisy. In light
 138 of this observation, we propose to demonstrate
 139 that the gradients obtained from three state-of-
 140 the-art models (ResNet50 [26], ViT [27], and
 141 ConvNeXT [28]) do indeed contain noise, pre-
 142 dominantly present in high-frequency compo-
 nents. To achieve this, we suggest an approach

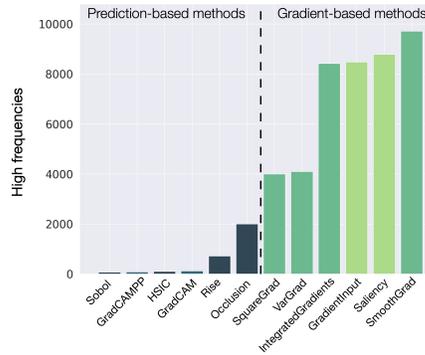


Figure 3: **High-frequency power in attribution methods.** High-frequency power present in the importance maps derived from different attribution methods. Prediction-based methods produce less high-frequency content than gradient-based methods.

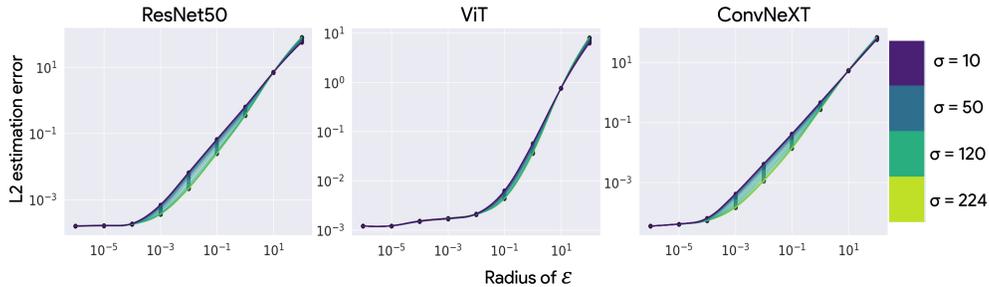


Figure 4: **Evidence for noise in the gradient.** We plot the residual of the first-order approximation of the model, that is $f(x + \epsilon) \approx f(x) + \epsilon \nabla_x f(x)$, with the gradient ∇f filtered at different bandwidths σ . We sample 100 values of ϵ uniformly on \mathbb{S}_{d-1} scaled by the radius, for 1,000 images of the validation set of ImageNet. If high-frequencies contained information necessary for a good linearization of the model then we would observe a gap between the curves of $\sigma = 224$ - where no filter is applied, vs. the curves where we apply a filter - $\sigma < 224$.

143 that involves selectively removing high-frequency gradient information by employing various fre-
 144 quency cutoffs σ . By computing residuals between $f(x + \epsilon)$ and its first order filtered by σ
 145 decomposition around x that we denote by $f(x) + \epsilon \nabla_x^\sigma f(x)$, we generate corresponding curves for
 146

147 different scales of ε , which are presented in Figure 4. As anticipated, our observations reveal that the
 148 curves exhibiting reduced high-frequency content (from $\sigma < 224$ to $\sigma = 10$) closely align with the
 149 one of the non-filtered gradient ($\sigma = 224$). In other words, the gradient remains approximately as
 150 informative, even after removing high-frequency information. This implies that the high-frequency
 151 content primarily contains noisy information within the gradient.

152 3.3 Investigating the mechanisms introducing noise

153 Next, we investigated the underlying operations responsible for the introduction of such content
 154 by computing the power of high-frequency content in the gradients at the level of all the layers.
 155 Notably, we observed a consistent trend in CNNs where high-frequency content tends to increase
 156 and jump at each block, indicative of downsampling operation through strided convolutions or
 157 pooling. This observation aligns with the findings of [29, 30], suggesting that downsampling
 158 operations via MaxPooling or strided convolution can introduce noise. To verify this hypothesis,
 159 we substituted these specific operations in two representative CNNs, namely ResNet50 [26] (which
 160 incorporates strided convolutions and MaxPooling) and VGG16 [31] (details in appendix) with
 161 AveragePooling. This replacement ensured the preservation of information continuity in the gradient.
 162 The resulting plots for ResNet50 (VGG) are presented in Figures 5 - bottom curve (see appendix
 163 for VGG), displaying the power of high-frequency content using Kolmogorov image compression
 164 and Laplace-operator (see appendix) at each step. The depicted red shades represent the amount of
 165 high-frequency content in both models. We observe that prior to the initial dimension change, the
 166 quantity of high-frequency content remain comparable, suggesting that operations within a block of
 167 the same dimension does not significantly increase the power of high-frequency content. However,
 168 with the introduction of a downsampling layer, the curves for each model diverge, indicating a
 169 bigger contribution to the introduction of high frequencies by striding or MaxPooling compared to
 170 AveragePooling. Our findings corroborate the observations of [29], as the gradients (even averaged)
 171 following MaxPooling or strides exhibit checkerboard patterns, providing a plausible explanation for
 172 our quantitative observation of increased high-frequency content.

173 We employ the same pipeline to calculate the high-frequency content for both a random model and a
 174 trained model, using the identical set of models including ViT [27]. The resulting curves are depicted
 175 in Figures 5 - top curve, for ResNet50 (see appendix for VGG16 and for ViT), showcasing that there is
 176 no discrepancy in high-frequency content between the trained and random CNNs. Given our previous
 177 section’s demonstration that high-frequencies carry negligible information for the model, one would
 178 expect that training could potentially eliminate this content, leaving only relevant information to be
 179 processed. However, as our observations indicate the absence of such behavior despite the models
 180 accomplishing the task, we propose that the models were unable to adapt the gradient’s content,
 181 thereby suggesting it to be an inherent by-product of downsampling operations. In the case of the ViT,
 182 however, training appears to introduce some high frequencies from the initial operation, potentially
 183 arising from transformers’ pre-processing functions, such as image flattening via patches. These
 184 multiple findings suggest that high-frequency content emerges as a by-product of particular operations,
 185 predominantly observed in CNNs, which the models are unable to modulate during training. We
 186 therefore propose to consider most of the high-frequency content as noise. Consequently, when
 187 generating explanations for the models’ decisions, it is justifiable to disregard high frequencies as
 188 they offer limited or negligible information.

189 3.4 FORGrad: a simple strategy to remove noise

190 **An adapted σ^* per model** With **FORGrad**, we propose to remove high-frequency content,
 191 considered as noise, in order to obtain an optimal explanation related to the optimal frequency band
 192 from the gradient. We therefore propose to apply a low-pass filter on the Fourier spectrum of the
 193 gradient, employing multiple frequency cutoffs spaced evenly apart. For each filtered explanation,
 194 we compute the score from two different metrics. The first one Deletion – denoted $D(\varphi(\mathbf{x}))$ [9] – is
 195 a measure of the decrease in the likelihood of a particular class as the important pixels (identified
 196 by the saliency map) are systematically removed from the image. If the likelihood of the class
 197 experiences a rapid decrease, resulting in a small area under the probability curve, this is a strong
 198 indication of a good explanation. Complementary, Insertion, $I(\varphi(\mathbf{x}))$ [9] measures the significance
 199 of the pixels based on their capacity to create an image, and is calculated by measuring the increase
 200 in the probability of the class of interest as pixels are added in accordance with the generated
 201 importance map. Overall, we propose a heuristic to optimize our σ^* , representing the ideal bandwidth
 202 maximizing the difference $\sigma^* = \arg \max_{\sigma} \mathbb{E}_{\mathbf{x}} D(\varphi_{\sigma}(\mathbf{x})) - I(\varphi_{\sigma}(\mathbf{x}))$, combining the score of both
 203 metrics, on a subset of the validation set of ImageNet (1,000 images). Using both deletion and

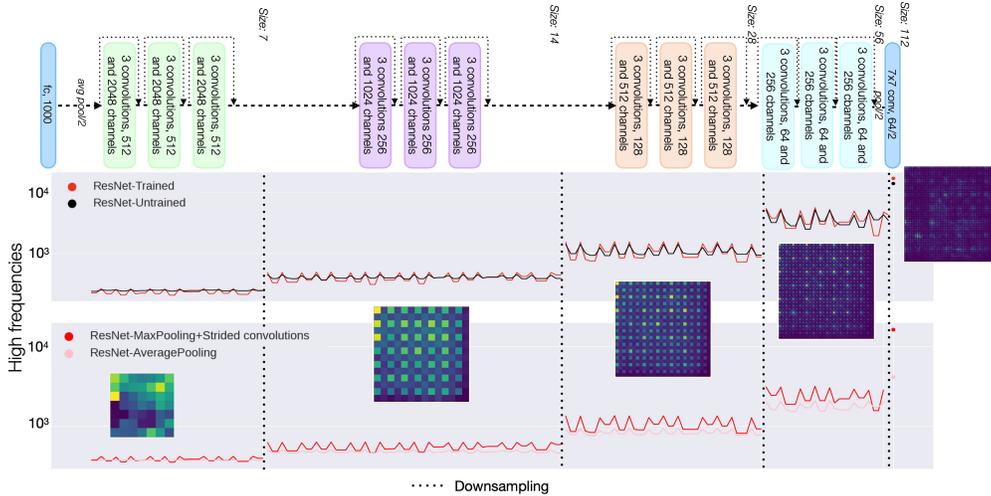


Figure 5: **Evolution of the high-frequency content in Resnet50.** We compute the high-frequency content along the depth of a ResNet50 varying either the weights or the pooling. The top curve represents the trained model, indicated by the red curve, while the untrained model is represented by the black curve. The bottom curve illustrates the impact of different poolings, with MaxPooling and stride shown in dark red and AveragePooling in pink. Each point on the graph corresponds to a layer within the models. In addition, we present visual examples of averaged gradients across 128 images after applying MaxPooling. Despite the averaging process, these examples exhibit checkerboard patterns, serving as a visual demonstration of the presence of high-frequencies.

204 insertion metrics can provide a more comprehensive evaluation of the quality of the attribution
 205 map or saliency map generated for a given model. The deletion metric is useful for identifying
 206 important regions of an image that contribute to a model’s decision, while the insertion metric is
 207 valuable for assessing the quality of the generated saliency map in terms of its ability to reconstruct
 208 the original image. By combining both metrics, we aim to assess the quality of the explanations
 209 generated by considering the impact of pixel removal and addition on the likelihood and significance
 210 of the target class, respectively. In the latter sections, we will consider the faithfulness metric to be
 211 the combination [Deletion-Insertion] Additionally, we also evaluate **FORGrad** on a third metric,
 212 MuFidelity, $F(\varphi(x))$, [32]. The fidelity correlation metric serves to verify the correlation between
 213 the attribution score and a random subset of pixels. To achieve this, a set of pixels is randomly chosen
 214 and set to a baseline state, after which a prediction score is obtained. The fidelity correlation metric
 215 evaluates the correlation between the decrease in the score and the significance of the explanation for
 216 each random subset created.

217 **Theoretical foundations** In this section, we build on the empirically demonstrated assumption that
 218 the gradient is noisy and prove, through a Fourier perspective, that **FORGrad** effectively recovers
 219 the true gradient. Moreover, assuming that the noise is originally Gaussian, we characterize the
 220 distribution of the noise in Fourier space. Finally, we propose a convergence bound for SmoothGrad,
 221 valid on finite samples, showing that it also recovers the true gradient at the cost of multiple samplings.

222 We denote $\|\cdot\|_F$ as the Frobenius norm and $\|\cdot\|_2$ as the spectral norm. Note that $\|\cdot\|_2 \leq \|\cdot\|_F$
 223 in order to interpret our results. Finally, we define $\mathbf{K}^\sigma \in \{0, 1\}^{W \times H}$ as the binary Fourier mask,
 224 parameterized by σ , that we used to filter high frequency, where each element $\mathbf{K}^\sigma_{(i,j)}$ is determined
 225 as $\mathbf{K}^\sigma_{(i,j)} = \mathbb{1}_{|i-\frac{W}{2}| \leq \sigma} \mathbb{1}_{|j-\frac{H}{2}| \leq \sigma}$, with $\mathbb{1}$ the indicator function and $\bar{\mathbf{K}}^\sigma = 1 - \mathbf{K}^\sigma$. As we have
 226 discussed above, the gradient of deep models is noisy, and in the following work, we consider that
 227 we only have access to $\nabla_{\mathbf{x}} \hat{\mathbf{f}}(\mathbf{x})$, a noisy estimator of $\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x})$ such that $\nabla_{\mathbf{x}} \hat{\mathbf{f}}(\mathbf{x}) = \nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}) + \varepsilon$
 228 with $\varepsilon \in \mathbb{R}^{W \times H}$. We do not assume any randomness for the noise so far. The following proposition
 229 develops the squared residual of the filtered noisy gradient as compared to the true one. Under the
 230 condition of finding the optimal filter, the gap between the two is naturally norm of the remaining
 231 noise post filtering.

232 **Proposition 3.1.** Let $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{Y}$ a predictor, and denote $\nabla \hat{\mathbf{f}} = \nabla \mathbf{f} + \varepsilon$ as the noisy gradient of \mathbf{f} ,
 233 with $\varepsilon \in \mathbb{R}^{W \times H}$. For $\sigma^* = \inf \{ \sigma : \|\mathcal{F}(\nabla \hat{\mathbf{f}}) \odot \mathbf{K}^\sigma\|_F^2 = 0 \}$, we have

$$\|\mathcal{F}^{-1}(\mathcal{F}(\nabla \hat{\mathbf{f}}) \odot \mathbf{K}^{\sigma^*}) - \nabla \mathbf{f}\|_F^2 = \|\mathcal{F}^{-1}(\mathcal{F}(\varepsilon) \odot \mathbf{K}^{\sigma^*})\|_F^2 \leq \|\varepsilon\|_F^2, \quad (1)$$

234 where \odot is the Hadamard product, \mathbf{K}^{σ^*} a binary mask for low-pass filtering of frequency σ , and
 235 $\bar{\mathbf{K}}^{\sigma^*}$ is the opposite mask.

236 *Remark 3.2.* This result holds as long as we find σ^* . There always exists a σ^* as the set always
 237 contains $\sigma = \max(H, W)$ which does not alter the Fourier spectrum of an image of size $W \times H$.
 238 However, finding σ^* poses a challenge, leading us to leverage XAI metrics as a heuristic.

239 With the information that the remaining gap between the filtered estimator and the true gradient is
 240 the remaining noise, of which the norm is upper bounded by the one of the original noise, we aim
 241 at measuring the reduction of the noise. In that way, we demonstrate the always-positive effect of
 242 **FORGrad** on gradient methods. In particular, under the assumption of Gaussian noise, we derive
 243 the distribution of the ratio $\|\varepsilon\|_F^2 / \|\mathcal{F}^{-1}(\mathcal{F}(\varepsilon) \odot \mathbf{K}^{\sigma^*})\|_F^2$.

244 **Proposition 3.3.** *Let the noise $\varepsilon \in \mathbb{R}^{W \times H}$ follow a normal distribution $\varepsilon \sim \mathcal{N}(0, \varsigma)^{\otimes N}$. Then the
 245 norm of the Fourier spectra of the noise $\|\mathcal{F}(\varepsilon)\|_F^2 \sim \Gamma(k = 2WH, \theta = \varsigma^2WH)$ and filtered noise
 246 $\|\mathcal{F}(\varepsilon) \odot \mathbf{K}^{\sigma^*}\|_F^2 \sim \Gamma(k = 8\sigma^2, \theta = 4\varsigma^2\sigma^2)$ follow Gamma distributions.*

247 *Therefore, the ratio of the two distributions $R = \|\mathcal{F}(\varepsilon)\|_F^2 / \|\mathcal{F}(\varepsilon) \odot \mathbf{K}^{\sigma^*}\|_F^2$ follows a Beta prime
 248 distribution $R \sim \beta'(2N, 8\sigma^2, 1, \frac{WH}{4\sigma^2})$.*

249 This result allows us to directly compute the distribution of the ratio of the norm of the original
 250 noise on the norm of the filtered noise (up to a scaling factor, by Parseval’s Theorem). Naturally,
 251 this distribution depends on the parameter σ of the filtering. From this, we can deduce probabilistic
 252 results, such as, for $\sigma = 10$ and $\varsigma = 1$, the ratio of the norms is larger than 70 with probability almost
 253 one.

254 Finally, in the following proposition, we obtain a non-asymptotic result on the concentration of the
 255 SmoothGrad procedure to its expected value based on the Matrix Bernstein inequality [33].

256 **Proposition 3.4.** *We recall that SmoothGrad is defined as $SG = \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x} + \delta_i)$
 257 with $\forall_{i=1, \dots, n} \delta_i \in \mathcal{N}(0, \varsigma)$. Here we compute SG on the noisy estimate of the gradient
 258 $\nabla_{\mathbf{x}} \hat{\mathbf{f}}(\mathbf{x} + \delta_i) \in \mathbb{R}^{W \times H}$, which is then a random matrix \widehat{SG} . Assuming our predictor
 259 $\mathbf{f} \in L\text{-Lip}(\mathcal{X})$ is L -Lipschitz. We denote $\|\cdot\|_2$ as the spectral norm, and define the variance
 260 as $\mathbb{V}(SG) = \max(\|\mathbb{E}((SG - \mathbb{E}SG) \cdot (SG - \mathbb{E}SG)^T)\|_2, \|\mathbb{E}((SG - \mathbb{E}SG)^T \cdot (SG - \mathbb{E}SG))\|_2)$.
 261 We then have, for $t > 0$,*

$$\mathbb{P}\left(\|\widehat{SG} - \mathbb{E}SG\|_2 \geq t\right) \leq (W + H) \cdot \exp\left(\frac{-t^2 n^2 / 2}{\mathbb{V}(\widehat{SG}) + 2Lt/3}\right). \quad (2)$$

262 Our results suggest that in order to effectively eliminate noise using the SmoothGrad method, several
 263 iterations are required as opposed to ours. For instance, to be at least $t = \frac{L}{10}$ away from its expected
 264 value, with probability at most 0.01, we need $n \approx 700$ iterations, for $\varsigma = 1$. Furthermore, the noisy
 265 SmoothGrad gradually approaches the expected outcome of the non-noisy SmoothGrad. Additionally,
 266 SmoothGrad alleviate the noise but at the cost of employing Monte-Carlo sampling.

267 4 Gradient-based methods perform better and are more efficient

268 **The new explanations are free from noise** Figure 6 presents qualitative examples of corrected
 269 gradients obtained using the Gradient Input method [34] combined with **FORGrad**. As we analyze
 270 the different images, we observe that gradually removing high frequencies from the gradients has a
 271 notable impact on the resulting explanation. The initially noisy patterns transform into larger patches
 272 until the saliency map effectively highlights the key features that represent the object for categorization.
 273 However, it is crucial to consider the optimal value of σ^* , as exceeding this threshold leads to the
 274 map spreading too widely and the explanation becoming less informative. This observation is further
 275 supported by the curve on the right, which demonstrates the evolution of the faithfulness score as σ
 276 changes. Prior to finding the optimal σ , the faithfulness score fluctuates around the initial value before
 277 gradually increasing to reach its optimal level. As expected, when all the information is removed
 278 (represented by the last point on the x-axis), the fidelity score drops to zero.

279 **A new ranking of attribution methods** We apply **FORGrad** on all the gradient-based attri-
 280 bution methods and report the scores in Table 1 for three models: ResNet50 [26], ViT [27] and
 281 ConvNeXT [28]. GradCAM methods can’t be tested on ViT because they are based on convolution so
 282 are limited to CNNs. We can observe two notable findings. Firstly, it is rare to encounter cases where

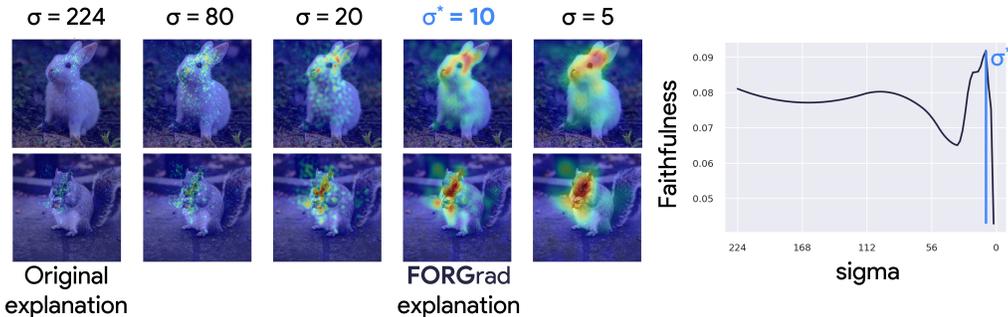


Figure 6: **FORGrad : selection of the optimal σ** . With **FORGrad**, we aim to derive the explanation from the gradient’s corrected version. To achieve this, we determine the optimal cutoff of high frequencies in order to maximize the faithfulness of the explanation. The images displayed on the left illustrate the progression for different cutoff values. On the right side, the curve represents the variation of the metric across 1,000 images from the validation set of ImageNet, as a function of σ , representing the cutoff value. σ^* represents the optimal value, selected as the one maximizing the faithfulness score, equivalent to [Deletion-Insertion].

	ResNet50				ViT				ConvNeXT				
	Del.(↓)	Ins.(↑)	Fid.(↑)	Comp.	Del.(↓)	Ins.(↑)	Fid.(↑)	Comp.	Del.(↓)	Ins.(↑)	Fid.(↑)	Comp.	
Gradient-based	Saliency[1]	0.77	0.85	0.07	$\Theta(2)$	0.77	0.82	0.01	$\Theta(2)$	0.85	0.86	0.05	$\Theta(2)$
	Saliency*	0.74	0.90	0.15	$\Theta(2)$	0.81	0.89	0.03	$\Theta(2)$	0.82	0.88	0.06	$\Theta(2)$
	GradInput[34]	0.76	0.87	0.05	$\Theta(2)$	0.78	0.88	0.01	$\Theta(2)$	0.83	0.89	0.05	$\Theta(2)$
	GradInput*	0.74	0.88	0.13	$\Theta(2)$	<u>0.74</u>	0.89	0.05	$\Theta(2)$	<u>0.81</u>	0.88	<u>0.07</u>	$\Theta(2)$
	SmoothGrad[5]	0.74	0.89	0.08	$\Theta(200)$	0.80	0.87	0.03	$\Theta(200)$	0.86	0.86	0.05	$\Theta(200)$
	SmoothGrad*	0.72	0.93	0.19	$\Theta(200)$	0.78	0.92	<u>0.04</u>	$\Theta(200)$	0.82	0.88	0.06	$\Theta(200)$
	VarGrad[35]	0.74	0.91	0.07	$\Theta(200)$	0.72	0.88	0.01	$\Theta(200)$	0.89	0.86	0.02	$\Theta(200)$
	VarGrad*	<u>0.73</u>	0.91	<u>0.18</u>	$\Theta(200)$	0.74	0.90	0.02	$\Theta(200)$	0.80	0.88	0.02	$\Theta(200)$
	Int.Grad[6]	0.75	0.88	0.06	$\Theta(200)$	0.78	0.86	0.01	$\Theta(200)$	0.82	0.90	0.05	$\Theta(200)$
	Int.Grad*	0.74	0.89	0.15	$\Theta(200)$	0.79	0.87	0.03	$\Theta(200)$	<u>0.81</u>	0.90	0.05	$\Theta(200)$
Prediction-based	GradCAM[8]	0.78	<u>0.92</u>	0.06	$\Theta(2)$	n.a	n.a	n.a	n.a	0.87	0.92	0.06	$\Theta(2)$
	GradCAM++[36]	0.75	0.93	0.08	$\Theta(2)$	n.a	n.a	n.a	n.a	0.90	0.92	0.02	$\Theta(2)$
	Occlusion[34]	0.75	0.85	0.06	$\Theta(1024)$	0.79	0.83	0.01	$\Theta(1024)$	0.83	0.88	<u>0.07</u>	$\Theta(1024)$
	HSIC[37]	0.72	<u>0.92</u>	0.05	$\Theta(2000)$	0.77	<u>0.91</u>	0.02	$\Theta(2000)$	0.80	0.92	0.05	$\Theta(2000)$
	Sobol[38]	0.74	<u>0.92</u>	0.06	$\Theta(4000)$	0.79	<u>0.91</u>	0.02	$\Theta(4000)$	0.82	<u>0.93</u>	0.08	$\Theta(4000)$
	RISE[9]	0.76	0.93	0.07	$\Theta(8000)$	0.80	0.92	0.01	$\Theta(8000)$	0.84	0.94	<u>0.07</u>	$\Theta(8000)$

Table 1: **Results on Faithfulness metrics**. Deletion, Insertion, and Fidelity scores obtained on 1,000 ImageNet validation set images, on an Nvidia V100 (For Deletion, lower is better and for Insertion and Fidelity, higher is better). Complexity Θ (Comp.) corresponds to the number of forward + backward passes required for computation, up to a factor that depends on the model. The first and second best results are in **bold** and underlined.

283 the scores after applying **FORGrad** are lower than the scores obtained before. In such instances, the
284 decrease in scores is typically observed in only one metric, either Deletion or Insertion. However,
285 since the other metric is optimized, the overall Faithfulness, as measured by [Deletion-Insertion],
286 remains at least as good as before. Secondly, even without explicitly optimizing the Fidelity metric,
287 we observe an improvement in this score across all methods and the three models analyzed. Fur-
288 thermore, after applying **FORGrad**, we observe that the scores of several gradient-based methods
289 surpass or at least match those of prediction-based methods. Notably, these gradient-based methods
290 offer the additional advantage of being significantly more computationally efficient, as evident from
291 the complexity column. In order to determine the best method for each model, we propose to
292 aggregate the scores from the three metrics to obtain a single global score for each method and model.
293 This resulting score, is denoted as $I(\varphi(\mathbf{x})) + F(\varphi(\mathbf{x})) - D(\varphi(\mathbf{x}))$, corresponding to the sum of
294 1-Deletion, Insertion and Fidelity score. Interestingly, in Table 2, we observe that the rankings change
295 when we incorporate **FORGrad** into the analysis. This shift leads to the inclusion of at least two
296 gradient-based methods among the top-5 for all three models. In the case of ResNet50, all five of the
297 top-performing methods are gradient-based, whereas only one of them occupied a position in the
298 previous ranking. Although some prediction-based methods, such as *Sobol* and *HSIC*, consistently

	ResNet50		ViT		ConvNeXT	
	Original	FORGrad	Original	FORGrad	Original	FORGrad
1	GradCAM++	SmoothGrad*	VarGrad	SmoothGrad*	Sobol	Sobol
2	HSIC	VarGrad*	HSIC	VarGrad*	RISE	RISE
3	RISE	Saliency*	Sobol	HSIC	HSIC	HSIC
4	Sobol	Int.Grad*	RISE	Sobol	Occlusion	GradInput*
5	VarGrad	GradInput*	GradInput	RISE	GradCAM	Int.Grad*

Table 2: **Global ranking before (original) and after FORGrad.** For each model, we show the 5 attribution methods with highest metrics, before and after applying **FORGrad**. The explanation maps were computed on 1000 images from the validation set of ImageNet, based on an aggregation of the three metrics computed by $I(\varphi(x)) + F(\varphi(x)) - D(\varphi(x))$.

299 demonstrate good performance, we demonstrate that gradient-based methods such as *SmoothGrad*
300 and *VarGrad* now perform nearly as well, with the added advantage of computational efficiency.

301 5 Limitations

302 In our study, we have proposed to find an optimal σ value representing an ideal cutoff to improve
303 explanations of gradient-based methods. However, we acknowledge that this optimal value is highly
304 dependent on the dataset, perhaps more so than on the model itself. Furthermore, while we have
305 chosen a single value that maximizes the scores across 1,000 images, it may be beneficial to use
306 different values for individual images, but would increase the computational costs. We also optimize
307 our value of σ only on 2 metrics, deletion and insertion. Even though it turns out to also increase the
308 fidelity score, we could potentially obtain even better results by optimizing on this metric as well. It’s
309 however, once again, a very resource-consuming method that we chose to avoid. Furthermore, in our
310 ranking computation, we combine metrics that do not precisely capture the same information. While
311 Deletion and Insertion can be aggregated, particularly since we optimize the difference between them,
312 it should be noted that Deletion, Insertion, and Fidelity are not directly comparable even if they range
313 between 0 and 1. We have proposed one approach to integrate these metrics and derive a ranking
314 based on the three scores. However, an alternative could involve producing separate rankings for each
315 individual score. If we had followed this approach, the **FORGrad** methods would have emerged as
316 the top-5 for both ResNet50 and ViT, according to MuFidelity.

317 6 Conclusion

318 This work started with an empirical observation: prediction-based and gradient-based methods
319 exhibit distinct power spectra in their attribution maps – with gradient-based methods exhibiting
320 higher power in the high frequencies compared to prediction-based methods. This led us to wonder
321 whether the frequency content of model gradients is merely noisy information. We demonstrate
322 that removing this content does not impair our ability to approximate the gradient and conclude that
323 high frequencies predominantly carry non-essential information. We further conducted an in-depth
324 analysis of gradient frequency content in CNNs across processing layers and found that downsampling
325 operations, such as max pooling and striding, contribute to the introduction of high frequencies.
326 This points to model aliasing as a likely cause of this high-frequency content. Interestingly, even
327 with training, CNNs are unable to prevent this phenomenon. These results hence raise the question:
328 Could high-frequencies be filtered out to improve the explanations derived from attribution methods?
329 We design an optimal filter, σ^* , and show that the filtering of attribution maps leads to significant
330 improvements in the quality of the explanations. These improvements were most pronounced for
331 gradient-based methods, which ended up approaching and sometimes even surpassing the much
332 more compute-intensive prediction-based methods. Overall, our work leads to a surprising result –
333 that the almost forgotten gradient-based methods turn out to contain all the information needed to
334 provide a faithful explanation of a model’s decision and that they can be as interpretable as the newest
335 methods. In future work, it would be worth exploring the influence of this noise on the model’s
336 performance and evaluating whether replacing certain operations that introduce noise could affect
337 both the accuracy and robustness of the models. Furthermore, considering that many adversarial
338 attacks are gradient-based and often exploit additive noise patterns, it is worth investigating whether
339 these attacks target the noisy high-frequency content in the gradients and whether they might be
340 prevented by using operations not introducing high-frequencies.

References

- 341
- 342 [1] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising
343 image classification models and saliency maps. In *Workshop, Proceedings of the International Conference
344 on Learning Representations (ICLR)*, 2013.
- 345 [2] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising
346 image classification models and saliency maps. In *Workshop Proceedings of the International Conference
347 on Learning Representations (ICLR)*, 2014.
- 348 [3] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Proceedings
349 of the IEEE European Conference on Computer Vision (ECCV)*, 2014.
- 350 [4] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for sim-
351 plicity: The all convolutional net. In *Workshop Proceedings of the International Conference on Learning
352 Representations (ICLR)*, 2014.
- 353 [5] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad:
354 removing noise by adding noise. In *Workshop on Visualization for Deep Learning, Proceedings of the
355 International Conference on Machine Learning (ICML)*, 2017.
- 356 [6] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings
357 of the International Conference on Machine Learning (ICML)*, 2017.
- 358 [7] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features
359 for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern
360 recognition*, pages 2921–2929, 2016.
- 361 [8] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and
362 Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In
363 *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- 364 [9] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box
365 models. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2018.
- 366 [10] Thomas Fel, Rémi Cadène, Mathieu Chalvidal, Matthieu Cord, David Vigouroux, and Thomas Serre. Look
367 at the variance! efficient black-box explanations with sobol-based sensitivity analysis. *Advances in Neural
368 Information Processing Systems*, 34:26005–26014, 2021.
- 369 [11] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the
370 predictions of any classifier. In *Knowledge Discovery and Data Mining (KDD)*, 2016.
- 371 [12] Beomsu Kim, Junghoon Seo, Seunghyeon Jeon, Jamyoungh Koo, Jeongyeol Choe, and Taegyun Jeon. Why
372 are saliency maps noisy? cause of and solution to noisy saliency maps. In *2019 IEEE/CVF International
373 Conference on Computer Vision Workshop (ICCVW)*, pages 4149–4157. IEEE, 2019.
- 374 [13] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and
375 Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In
376 *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- 377 [14] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert
378 Müller. How to explain individual classification decisions. *The Journal of Machine Learning Research*, 11:
379 1803–1831, 2010.
- 380 [15] Zhiqin John Xu. Understanding training and generalization in deep learning by fourier analysis. *arXiv
381 preprint arXiv:1808.04295*, 2018.
- 382 [16] Dong Yin, Raphael Gontijo Lopes, Jon Shlens, Ekin Dogus Cubuk, and Justin Gilmer. A fourier perspective
383 on model robustness in computer vision. *Advances in Neural Information Processing Systems*, 32, 2019.
- 384 [17] Yusuke Tsuzuku and Issei Sato. On the structural sensitivity of deep convolutional networks to the
385 directions of fourier basis functions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
386 Pattern Recognition*, pages 51–60, 2019.
- 387 [18] Gintare Karolina Dziugaite, Zoubin Ghahramani, and Daniel M Roy. A study of the effect of jpg
388 compression on adversarial images. *arXiv preprint arXiv:1608.00853*, 2016.
- 389 [19] Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens Van Der Maaten. Countering adversarial images
390 using input transformations. *arXiv preprint arXiv:1711.00117*, 2017.
- 391 [20] Sibong Song, Yueru Chen, Ngai-Man Cheung, and C-C Jay Kuo. Defense against adversarial attacks with
392 saak transform. *arXiv preprint arXiv:1808.01785*, 2018.
- 393 [21] Jason Jo and Yoshua Bengio. Measuring the tendency of cnns to learn surface statistical regularities. *arXiv
394 preprint arXiv:1711.11561*, 2017.
- 395 [22] Ramesh C. Jain, Rangachar Kasturi, and Brian G. Schunck. Machine vision. 1995.

- 396 [23] Said Pertuz, Domenec Puig, and Miguel Angel Garcia. Analysis of focus measure operators for shape-
397 from-focus. *Pattern Recognition*, 46(5):1415–1432, 2013.
- 398 [24] Hector Zenil, Jean-Paul Delahaye, and Cédric Gaucherel. Image characterization and classification by
399 physical complexity. *Complexity*, 17(3):26–42, 2012.
- 400 [25] Andrei N Kolmogorov. On tables of random numbers. *Sankhyā: The Indian Journal of Statistics, Series A*,
401 pages 369–376, 1963.
- 402 [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
403 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 404 [27] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
405 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth
406 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 407 [28] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A
408 convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
409 Recognition (CVPR)*, pages 11976–11986, June 2022.
- 410 [29] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2017. doi:
411 10.23915/distill.00007. <https://distill.pub/2017/feature-visualization>.
- 412 [30] Xueyan Zou, Fanyi Xiao, Zhiding Yu, Yuheng Li, and Yong Jae Lee. Delving deeper into anti-aliasing in
413 convnets. *International Journal of Computer Vision*, 131(1):67–81, 2023.
- 414 [31] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recogni-
415 tion. *arXiv preprint arXiv:1409.1556*, 2014.
- 416 [32] Umang Bhatt, Adrian Weller, and José MF Moura. Evaluating and aggregating feature-based model
417 explanations. *arXiv preprint arXiv:2005.00631*, 2020.
- 418 [33] Joel A Tropp et al. An introduction to matrix concentration inequalities. *Foundations and Trends® in
419 Machine Learning*, 8(1-2):1–230, 2015.
- 420 [34] Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Towards better understanding of gradient-
421 based attribution methods for deep neural networks. In *Proceedings of the International Conference on
422 Learning Representations (ICLR)*, 2018.
- 423 [35] Junghoon Seo, Jeongyeol Choe, Jamyounng Koo, Seunghyeon Jeon, Beomsu Kim, and Taegyun Jeon.
424 Noise-adding methods of saliency map as series of higher order partial derivative. In *Workshop on Human
425 Interpretability in Machine Learning, Proceedings of the International Conference on Machine Learning
426 (ICML)*, 2018.
- 427 [36] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++:
428 Generalized gradient-based visual explanations for deep convolutional networks. In *Proceedings of the
429 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2018.
- 430 [37] Paul Novello, Thomas Fel, and David Vigouroux. Making sense of dependence: Efficient black-box
431 explanations using dependence measure. In *Advances in Neural Information Processing Systems (NeurIPS)*,
432 2022.
- 433 [38] Thomas Fel, Remi Cadene, Mathieu Chalvidal, Matthieu Cord, David Vigouroux, and Thomas Serre. Look
434 at the variance! efficient black-box explanations with sobol-based sensitivity analysis. In *Advances in
435 Neural Information Processing Systems (NeurIPS)*, 2021.