

Ablation Path Saliency

Justus Sagemüller¹[0000-0003-1882-1096] and Olivier Verdier¹[0000-0003-3699-6244] {over, jsag}@hvl.no

Western Norway University of Applied Sciences

Abstract. Various types of saliency methods have been proposed for explaining black-box classification. In image applications, this means highlighting the part of the image that is most relevant for the current decision.

Unfortunately, the different methods may disagree and it can be hard to quantify how representative and faithful the explanation really is. We observe however that several of these methods can be seen as edge cases of a single, more general procedure based on finding a particular *path* through the classifier’s domain. This offers additional geometric interpretation to the existing methods.

We demonstrate furthermore that ablation paths can be directly used as a technique of its own right. This is able to compete with literature methods on existing benchmarks, while giving more fine-grained information and better opportunities for validation of the explanations’ faithfulness.

Keywords: Explainability, Classification, Saliency, Neural Networks, Visualisation, Gradient Descent

1 Introduction

The basic idea of *saliency* or *attribution* is to provide insights as to why a neural network produces a given output (for instance, a classification) for a given input (for instance, an image). There is no clear consensus in the literature as to what saliency should exactly be, but various properties that such a method should fulfill have been proposed. All the methods discussed here start out by contrasting the given input (also called *current target*) with another one, called *baseline*, which should be neutral in at least the sense of not displaying any of what causes the target image’s classification. The saliency problem then amounts to finding out what the features of the target are which cause it to be classified differently from the baseline.

In [14] the authors give axioms attempting to make it precise what that means. Of these, *sensitivity* captures most of the notion of saliency, namely, that the features on which the output is most sensitive should be given a higher saliency value. The authors give further axioms to narrow it down: implementation invariance, completeness, linearity and symmetry preservation. They obtain a corresponding method: the Integrated Gradient method. Despite the attempt to thus narrow down the choice of saliency method, Integrated Gradient has

not established itself as a default in the community. Indeed the axioms used to justify it are not altogether self-evident.

In [6], a method is provided whose construction is quite different. Instead of following axioms about the properties a method should have, they produce a result that has direct meaning associated to it, namely as a mask that preserves only certain parts of the input and removes others, optimised so that the classification is retained even at high degrees of ablation, i.e., when the mask only keeps small part of the target image. This method is highly appealing, but in practice the optimisation problem is ill-conditioned and can only be solved under help of regularization techniques. That prevents this technique, too, from being a definite saliency method or “the” saliency method.

Various other methods from the literature are in a broadly similar position, all with certain arguments for their use but also various practical limitations and no clear reason to favour them over the alternatives. In some cases there are evident mathematical relationships between the methods, but they have not been investigated thoroughly yet or exploited for a unifying generalization.

This is what our paper provides: it introduces *ablation paths*, which take up and extend the idea of integrating from the baseline to the target image. It combines this with the notion of ablation / masks, in that each step along the path can constitute a mask highlighting progressively smaller portions of the image. The main purpose of this is mathematical unification and better (meta-) understanding of the various methods, but ours can also be used as a saliency method by itself.

A summary how the method works: suppose first that images are defined over a domain Ω , which can be regarded as the set of pixels in the discrete case, or as a domain such as a square, for the image at infinite resolution. We define *ablation paths* as parameter dependent smooth masks $\varphi: [0, 1] \rightarrow \mathcal{C}(\Omega, \mathbb{R})$, with the further requirement that the mask at zero, $\varphi(0)$, should be zero over the domain Ω , and the mask at one, $\varphi(1)$, should be one over the domain Ω . We also impose that, at each pixel, the mask value increases over time (see Figure 1), and that this happens with a constant area speed: the area covered by the mask should increase linearly over time (see §3). Let F be the classifier, which outputs a probability between zero and one. We choose a current image of interest x_0 and a *baseline image* x_1 . The objective function P_{\uparrow} is then $P_{\uparrow}(\varphi) := \int_0^1 F(x_0 + \varphi(t)(x_1 - x_0)) dt$ (see §4). Assuming that $F(x_0) \simeq 1$ and $F(x_1) \simeq 0$, maximising the objective function means that we try to find an ablation path that stays as long as possible in the decision region of x_0 . Intuitively, we try to replace as many pixels of x_0 by pixels of x_1 while staying in the same class as x_0 .

2 Related Work

[12] defines a saliency map as the gradient of the network output at the given image. This would appear to be a sensible definition, but the resulting saliency is very noisy because the network output is roughly constant around any particular

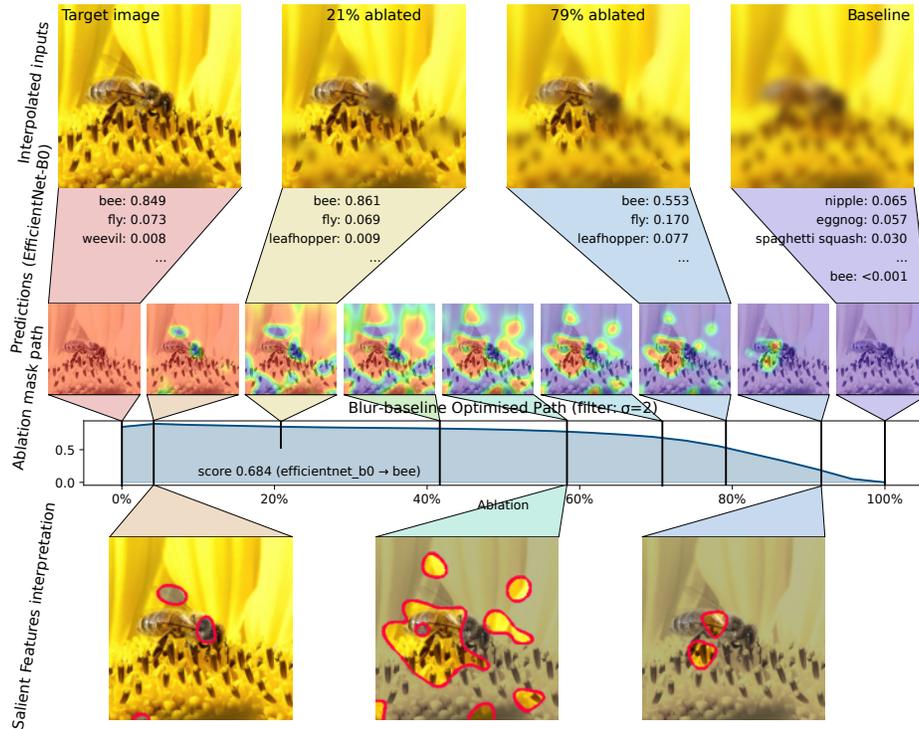


Fig. 1: Example of how an ablation path (sequence of masks, middle row) gives rise to a transition between a current target (a bee from ImageNet) and a baseline (blurred version of the same image).

image. [11] improves the situation by computing the gradient after each layer instead. This is, however, not a black-box method such as the one we propose. [8] computes an *influence function*, that is, a function that measures how the parameters would be changed by a change in the training data. Although it is a black-box method, it is not a saliency method per se. They use the gradient of the network output to find the pixel most likely to have a high saliency. The pixels that have most effect are given a higher saliency.

By contrast, [9] proposes to directly evaluate the saliency by finding out which pixels are most likely to affect the output, similarly to [6], but through statistical means instead of iterative optimisation. These methods can be seen as different ways of solving similar optimisation problems, the solution of which produces a mask (cf. § 3.1) highlighting features of importance.

There are also a number of meta-studies of saliency methods. [1] lists essential properties, for instance the requirement that the results should depend on the training data in a sense that perturbing *model parameters* should change the saliency. [7] proposes a number of properties that saliency methods should satisfy.

[2] compares several saliency methods and proposes a method to evaluate them (the sensitivity- n property).

These properties were not in the main focus of the design of our method, but we do fulfill the general criteria. For example, [7] is concerned with constant shifts in the inputs. If such a shift is applied to both the target and baseline, then the modification commutes through the interpolation and if we then assume a modified network that deals with such inputs in the same way as the original did with unshifted ones then all the gradients will be the same, therefore the optimised paths will also be the same.

Finally, [3] developed a method to recursively mask out ever finer “superpixels”. Their construction claims to be based on solid causality principles, and it does seem to largely avoid the adversarial issues that perturbation- / gradient- / statistics-based methods all have to content with (including ours), however the actual implementation approximates the principle in a nondeterministic, mathematically quite unclear way, and the technical details are highly pertaining to the image-classification application. (Still it is a black-box saliency method.)

3 Ablation Paths

3.1 Images and Masks

We assume data is represented by functions on a compact domain Ω . Examples for Ω may be $\Omega = \{1, \dots, n\} \times \{1, \dots, m\}$ (for pixel images), or a continuous domain $\Omega = [0, a] \times [0, b]$. What matters to us is that Ω is equipped with a positive measure. Without loss of generality, we assume the mass of that measure is one, i.e., $\int_{\Omega} 1 = 1$.

The data itself consists of functions on Ω with values in a vector space V (typically, the dimensions of V may be the *colour channels*). For the space of images, we choose

$$\mathcal{I} := \mathcal{C}(\Omega, V).$$

For our method to work, we need a space of *masks*, denoted \mathcal{M} , whose role is to select features between x_0 and x_1 . We associate to each mask $\theta \in \mathcal{M}$ an *interpolation operator* between two images x_0 and x_1 , denoted by $[x_0, x_1]_{\theta}$. This interpolator should have the property that $[x_0, x_1]_0 = x_0$ and $[x_0, x_1]_1 = x_1$. We will thus use the shorthand notation:

$$x_{\theta} := [x_0, x_1]_{\theta} \quad \theta \in \mathcal{M} \quad (x_0, x_1 \in \mathcal{I}).$$

The specific choice of masks and interpolation we make in this paper is

$$\mathcal{M} := \mathcal{C}(\Omega, \mathbb{R}), \quad [x_0, x_1]_{\theta} := (1 - \theta)x_0 + \theta x_1, \quad \theta \in \mathcal{M}, \quad x_0, x_1 \in \mathcal{I}$$

Another example of interpolation is what in [5] is called the *pyramid of blur* perturbations. We also explored this possibility for our method, but we obtain better results with the affine interpolation¹. (A possible reason that this blur

¹ Called “fade perturbation” in [5].

perturbation technique did not work well is the start state, which in that case contains none of the highest-frequency components at all. This makes it likely for the classifier to behave completely different from the target input.)

3.2 Ablation Paths

Definition 1. We define the set \mathcal{A} of ablation paths as the set of functions $\varphi: [0, 1] \rightarrow \mathcal{M}$ fulfilling the following properties:

Boundary conditions $\varphi(0) = 0$ and $\varphi(1) = 1$

Monotonicity $t_1 \leq t_2 \implies \varphi(t_1) \leq \varphi(t_2) \quad t_1, t_2 \in [0, 1]$

Constant speed $\int_{\Omega} \varphi(t) = t \quad t \in [0, 1].$

We will call monotone paths the paths that fulfill the first two conditions but not the third.

Note that the set \mathcal{A} of ablation paths is a *convex subset* inside the set of possible paths denoted by $\mathcal{P} := \mathcal{L}^{\infty}([0, 1], \mathcal{M})$.

Some comments on each of those requirements are in order. (i) 0 and 1 denote here the constant functions zero and one (which corresponds to the zero and one of the algebra \mathcal{M}) (ii) $\varphi(t_1) \leq \varphi(t_2)$ should be interpreted as usual as $\varphi(t_2) - \varphi(t_1)$ being pointwise nonnegative. (iii) If $t \mapsto \int_{\Omega} \varphi(t)$ is differentiable, this requirement can be rewritten as $\frac{d}{dt} \int_{\Omega} \varphi(t) = 1$, so it can be regarded as a *constant speed* requirement. This requirement is more a normalisation convention than a necessity, as is further detailed in Remark 2.

The simplest (requirement-fulfilling) ablation path is the *affine interpolation path*:

$$\ell(t) := t. \tag{1}$$

The mask is thus constant in space at each time t . This path is implicitly used in [14]: its image-application corresponds to affine interpolation between target- and baseline image.

Note that an ablation path without the constant-speed property can always be transformed into one that fulfils it. The proof is in § Appendix A.

Remark 1. In the sequel, we will abuse the notations and write φ as a function of one or two arguments depending on the context, that is, we will identify $\varphi(t) \equiv \varphi(t, \cdot)$. For instance, in the definition Definition 1 above, $\int_{\Omega} \varphi(t) \equiv \int_{\Omega} \varphi(t, \cdot) \equiv \int_{\Omega} \varphi(t, \mathbf{r}) \, d\mathbf{r}$.

Remark 2. If the ablation path φ is differentiable in time, the requirements in Definition 1 admit a remarkable reformulation. Define $\psi(t) := \frac{d}{dt} \varphi(t)$. All the requirements in Definition 1 are equivalent to the following requirements for a function $\psi: [0, 1] \times \Omega \rightarrow \mathbb{R}$:

$$\psi(t, \mathbf{r}) \geq 0, \quad \int_{\Omega} \psi(t, \mathbf{r}) \, d\mathbf{r} = 1, \quad \int_{[0,1]} \psi(t, \mathbf{r}) \, dt = 1 \quad t \in [0, 1], \mathbf{r} \in \Omega$$

The corresponding ablation path φ is then recovered by $\varphi(t) := \int_0^t \psi(s) \, ds$.

4 Score of an Ablation Path

We now define the *retaining score function* $P_{\uparrow}: \mathcal{P} \rightarrow \mathbb{R}$ from paths to real numbers by the integral

$$P_{\uparrow}(\varphi) := \int_0^1 F(x_{\varphi(t)}) dt. \quad (2)$$

Note that, as F is bounded between zero and one, so is $P_{\uparrow}(\varphi)$ for any ablation path $\varphi \in \mathcal{A}$. The main idea here is that $F(x_0) \simeq 1$ and $F(x_1) \simeq 0$, and $F(x) \leq 1$ holds always. So a high value of P_{\uparrow} means that the classification stays similar to that of x_0 over most of the ablation path, which is another way of saying that the characteristics of the original image are retained as best as possible whilst other features of the image are ablated away. See §6 for caveats.

Another score to consider is the *dissipating score*

$$P_{\downarrow}(\varphi) := 1 - \int_0^1 F(x_{\varphi(t)}) dt \quad (3)$$

which instead takes high values for paths that ablate crucial features for the current classification as quickly as possible. Optimisation of P_{\downarrow} corresponds roughly to what [6] call “deletion game”, whereas P_{\uparrow} corresponds to their “preservation game”, the difference to this work being that they optimise individual masks rather than constrained paths.

Intuitively, the first features to be deleted in a P_{\downarrow} -optimal path φ_{\downarrow} should correspond roughly to the ones longest preserved in a P_{\uparrow} -optimal path φ_{\uparrow} , meaning that a feature that is potent at retaining the classification should be removed early on if the objective is to change the classification. More generally, one would expect $\varphi_{\downarrow}(t)$ to be similar to $1 - \varphi_{\uparrow}(1 - t)$.

We observe this to be often *not* the case: specifically, there are many examples where either the classification is so predominant that it is almost indeterminate what features should be preserved longest (because any of them will be sufficient to retain the classification), or vice versa the classification is so brittle that it is indeterminate which ones should be removed first. It is however possible to *enforce* features to be considered simultaneously in a sense of their potency to preserve the classification when they are kept, and changing it when removed. This is achieved by optimising a path with the combined objective of retaining for the path itself and dissipating for its opposite: this is expressed by optimising the *contrastive score*

$$P_{\updownarrow}(\varphi) := P_{\uparrow}(\varphi) + P_{\downarrow}(1 - \varphi). \quad (4)$$

This too corresponds to ideas already used in previous work, called “hybrid game” or “symmetric preservation” [6][4].

A related possibility is to train both a retaining and a dissipating path in tandem, but with additional constraints to keep them in correspondence. Here, it is most useful to keep them not opposites of each other, but rather to keep

them as similar as possible. (Cf. Figure 4.) This is achieved by a score of the form

$$P_{\uparrow\downarrow}(\varphi_{\uparrow}, \varphi_{\downarrow}) := P_{\uparrow}(\varphi_{\uparrow}) + P_{\downarrow}(\varphi_{\downarrow}) + \lambda_{\pm} \|\varphi_{\uparrow} - \varphi_{\downarrow}\|, \quad (5)$$

where $\|\cdot\|$ could refer to various distance notions on the space of paths. We call the corresponding optimisation problem the *boundary-straddling method*, since (in the ideal of a classifier with exact decision boundaries) it rewards φ_{\uparrow} staying in the domain of x_0 as much as possible and φ_{\downarrow} in the domain of x_1 as much as possible, i.e. on the other side of the decision boundary but as close as possible. Thus, φ_{\uparrow} and φ_{\downarrow} effectively pinch the decision boundary between them.

For all the above score functions it is straightforward to compute the differential, e.g. dP_{\uparrow} , on the space of paths \mathcal{P} :

$$\langle dP_{\uparrow}, \delta\varphi \rangle = \int_0^1 \underbrace{\langle dF_{x_{\varphi(t)}}, \rangle}_{\in \mathcal{I}^*} \underbrace{(x_1 - x_0)}_{\in \mathcal{I}} \underbrace{\delta\varphi(t)}_{\in \mathcal{M}} dt \quad \delta\varphi \in \mathcal{P}.$$

So if we define the product of $D \in \mathcal{I}^*$ and $x \in \mathcal{I}$ producing an element in \mathcal{M}^* by $\langle xD, \varphi \rangle := \langle D, x\varphi \rangle$ as is customary², we can rewrite this differential as

$$\langle dP_{\uparrow}, \delta\varphi \rangle = \int_0^1 \langle (x_1 - x_0) dF_{x_{\varphi(t)}}, \delta\varphi(t) \rangle dt.$$

Note that we know that any ablation path is bounded, so $\varphi \in \mathcal{L}^{\infty}([0, 1], \mathcal{M})$, so the differential of P_{\uparrow} at φ can be identified with the function $dP_{\uparrow\varphi} = [t \mapsto (x_1 - x_0) dF_{x_{\varphi(t)}}]$ in $\mathcal{L}^1([0, 1], \mathcal{M}^*)$.

4.1 Relation with the Integrated Gradient Method

When this differential is computed on the interpolation path ℓ (1) and then *averaged*, then this is exactly the integrated average gradient [14]. More precisely, the integrated gradient is exactly $\int_0^1 dP_{\uparrow\ell(t)} dt$. Note that this is in fact an integrated *differential*, since we obtain an element in the dual space \mathcal{I}^* , and this differential should be appropriately smoothed along the lines of §5.1.

4.2 Relation to Pixel Ablation

Given any saliency function $\sigma \in \mathcal{M}$ (for example an integrated gradient, meaningful-perturbation, or grad-CAM result) we can define a path by

$$\tilde{\varphi}(t) := \mathbf{1}_{\sigma \leq \log(t/(1-t))} \text{ when } t \in (0, 1) \quad (6)$$

and $\tilde{\varphi}(0) := 0$, $\tilde{\varphi}(1) := 1$. This path is a monotone path, except in the module of images $\mathcal{I} = \mathcal{L}^2(\Omega, V)$, equipped with the ring of masks $\mathcal{M} = \mathcal{L}^{\infty}(\Omega)$. To be an ablation path, it still needs to be transformed into a constant speed path, which is always possible as explained in § Appendix A.

² For instance in the theory of distributions.

That results in a generalisation of the pixel-ablation scores used in [9] and [13]. In that case, the set Ω would be a discrete set of pixels, which are being sequentially switched from “on” to “off” by the (binary) mask.

Note that in the ranking, pixels with the same saliency would be ranked in an arbitrary way and added to the mask in that arbitrary order. In the method of Equation 6, we add such pixels all at once, which seems preferable because it does not incur an arbitrary bias between pixels. The time reparameterisation keeps the function constant longer to account for however *many* pixels were ranked the same. As long as the ranking is strict (no two pixels have the same saliency), the method is the same as discrete pixel ranking.

4.3 Relation to Meaningful Perturbations

In the saturated case, that is, if F only takes values zero and one (or in the limit towards this), our method reduces to finding the interpolation with the largest mask on the boundary, equivalent to the approach of [6]. This is a result of the following: suppose that the ablation path φ crosses the boundary at time t^* . It means that $F(x_{\varphi(t)})$ has value one until t^* and zero afterwards, so the score P_{\uparrow} defined in (2) is $P_{\uparrow}(\varphi) = t^*$. By the constant speed property, $t^* = \int_{\Omega} \varphi(t^*)$, so we end up maximising the mask area on the boundary.

4.4 Relation with RISE

The method used in [9] does not explicitly involve an optimisation problem like here, though they do use pixel ablation as some validation for the results. It does nevertheless resemble specifically the boundary-straddling method in the sense that it evaluates F for many different inputs on both sides of the decision boundary, and uses the classification results to weigh the features involved.

5 Optimisation Problem and Algorithm

We proceed to define the optimisation problem that we propose as a saliency method, and how to solve it numerically.

Conceptually we try to find the ablation path[s] (see Definition 1) that maximises one of the scores $P_{\uparrow}(\varphi)$, $P_{\downarrow}(\varphi)$, $P_{\uparrow\downarrow}(\varphi)$, or $P_{\uparrow\downarrow}(\varphi_{\uparrow}, \varphi_{\downarrow})$:

$$\max_{\varphi \in \mathcal{A}} P(\varphi).$$

Recall that the set \mathcal{A} of ablation paths is convex; however, since any of objective functions are not convex, this is not a convex optimisation problem.

The method we suggest is to follow a gradient direction. Such an approach is in general not guaranteed to approximate a global maximum, a common problem with many practical applications. However, empirical results (see § 7) suggests that gradient descent does often manage to approximate global maxima, particularly obvious in the unregularised near-perfect scores.

5.1 Gradient and Metric

In order to perform gradient descent, we need to be able to compute gradient vectors (sometimes called “natural gradients” in the literature). Strictly speaking the *differential* is an element of the dual space \mathcal{M}^* . In the Euclidean case that space is canonically isomorphic to \mathcal{M} , thus the common practice to use it directly as a *gradient* in \mathcal{M} , which is usable as contribution to a state update. In general, this requires first a map from that space to \mathcal{M} , and even in a discretised realisation it is prudent to consider this map explicitly, since the implied one depends on the (pixel) basis choice. A reasonable choice is the covariance operator associated to a smoothing operation. For a measure $\mu \in \mathcal{M}^*$, $\langle K\mu, \theta \rangle := \langle \mu, \int_{\Omega} k(\cdot - \mathbf{r})\theta(\mathbf{r}) \, d\mathbf{r} \rangle$, where k is a suitable smoothing function. We use here the same Gaussian blurring filter that is also applied between the optimisation steps for regularisation.

Since the optimisation problem is *constrained* (the ablation path φ being constrained by the requirements in Definition 1), following the gradient direction will in general leave the set \mathcal{A} . Because the constraints are convex, it is straightforward enough to project each gradient-updated version back to something that does fulfill them, and indeed that is the idea behind our algorithm (see § Appendix C for the technical details). However, in addition to the hard constraints there are also properties that are desirable but cannot directly be enforced. This is the subject of the next section.

5.2 Mask saturation

Recall that the masks we use in this paper are functions $\theta: \Omega \rightarrow [0, 1]$. The interpretation is that if $\theta(\mathbf{r}) = 0$, the pixel $\mathbf{r} \in \Omega$ of the reference image x_0 is used, whereas if $\theta(\mathbf{r}) = 1$, the pixel \mathbf{r} of the baseline image x_1 is used instead. Typically though, masks take value between zero and one. We notice that such intermediate values of masks produce blending of different images which lie far away from the natural distribution of images. What is problematic is that the classifier may put such blends of images in totally different classes. As analogy in human vision, an image of a person half-blended into an image of a hallway would not be seen as person present at 50%, but rather as something completely different; for instance, to a human, this half-present person would look more like a ghost than a person.

In order to alleviate this potential problem, our algorithm intersperses gradient descent in \mathcal{P} with both (hard) projections onto \mathcal{A} , as well as *soft projections* of the masks onto $\mathcal{M}_{\{0,1\}}$, the set of *saturated masks*, that is, masks taking value in the set $\{0, 1\}$. Concretely, this is done by tweaking the ablation path pointwise with a sigmoidal function³ that brings values lower than $\frac{1}{2}$ slightly closer to 0, and values greater than $\frac{1}{2}$ slightly closer to 1.

$$\varphi \leftarrow \Pi_{\text{sat}}(\varphi) := \frac{1}{2} \left(\frac{\tanh((\varphi \cdot 2 - 1) \cdot \zeta_{\text{sat}})}{\tanh(\zeta_{\text{sat}})} + 1 \right). \quad (7)$$

³ The exact definitions of Π_{sat} and Π_{pinch} are largely arbitrary, what matters are their attractive fixpoints; see Figure 2

The parameter ζ_{sat} determines how strongly this affects the path.

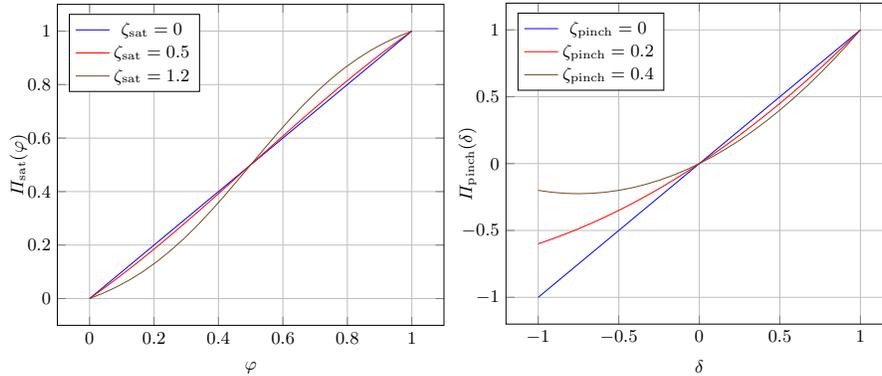


Fig. 2: The pointwise soft-projection functions for saturation and pinching.

5.3 Boundary-pinching

For the boundary-straddling method there is another requirement: making φ_{\uparrow} and φ_{\downarrow} similar to each other can be achieved by explicitly penalizing their distance in the score function, but in our implementation this too is done by a dedicated algorithm step that manipulates the masks pointwise to become more similar. For interpretability purposes it is particularly desirable for $\varphi_{\uparrow}(t)$ to contain only few features that $\varphi_{\downarrow}(t)$ does not, since that allows direct comparison between two images showing how inclusion of these features bring the classification into the target class. The exact difference in strength of features meanwhile is less relevant (even when the masks themselves are not boolean). Accordingly, we suggest a *pinching tweak* that diminishes specifically the smaller positive differences between φ_{\uparrow} and φ_{\downarrow} , in addition to any negative differences. The concrete form in our experiments is this: (recall that values close to 1 correspond to masked-away features)

$$\varphi_{\downarrow}(t, \mathbf{r}) \leftarrow \varphi_{\uparrow}(t, \mathbf{r}) + \Pi_{\text{pinch}}(\varphi_{\downarrow}(t, \mathbf{r}) - \varphi_{\uparrow}(t, \mathbf{r})) \quad (8)$$

where $\Pi_{\text{pinch}} : [-1, 1] \rightarrow [-1, 1], \delta \mapsto \Pi_{\text{pinch}}(\delta)$ is a continuous function with an attractive fixpoint at $\delta = 0$ (which is responsible for squelching unsubstantial contrasts between φ_{\uparrow} and φ_{\downarrow}), and a repulsive one at $\delta = 1$ (which allows the most salient features of φ_{\uparrow} to stay absent from φ_{\downarrow} , as necessary for a high $P_{\uparrow\downarrow}$).

The concrete definition of Π_{pinch} is uncritical, in our experiments we used

$$\Pi_{\text{pinch}}(\delta) := \delta(1 - \zeta_{\text{pinch}}) + \delta^2 \zeta_{\text{pinch}}.$$

Notice that in Equation 8, φ_{\uparrow} is not affected by φ_{\downarrow} , only vice versa. But conceptually, the update is performing a change to δ , i.e. the difference between the paths, rather than either of them individually.

6 Stability and adversarial effects

So far it was more or less taken for granted that a high-scoring ablation path owes its score to good highlighting of the features that were also responsible for the classification of x_0 . This assumption would be reasonable if any masked version of that image were classified either the same for the same reasons as the original, or else classified differently. For a black-box model, there is however no way of verifying this, and in fact it is simply not true in general.

It is well known [6] that sufficiently pathological masks can act as *adversarial attacks* [15] on an image, i.e. that masking out very minor parts of an image may affect the classification disproportionately and in ways that involve completely different neural activations (or whatever other concept is appropriate for the classifier architecture at hand). Gradient descent approaches tend to produce such examples easily. Although the study of adversarial effects is an important matter of its own right, they are hardly relevant for saliency purposes, because they do not necessarily involve the features that caused the classification of the original image.

A standard technique [6][9] employed to avoid that masks affect images adversarially is to regularize them in some sense of smoothness, e.g. by adding a total variation penalty. Intuitively, this at least prevents the masked image from featuring sharp edges or similar details not present in x_0 that the classifier might latch onto. We implemented this in terms of a simple Gaussian filter applied to each mask in the path after each optimisation step. Empirically, strong enough smoothing does largely avoid adversarial classification, but unfortunately there is no a-priori way of telling how smooth it needs to be. And too strong smoothing can also have detrimental effects. Not only does it prevent the exact localization of small, salient features, but it can even bias the outcome: in Figure 3,

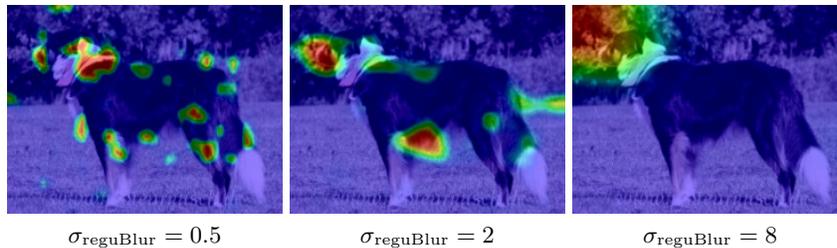


Fig. 3: Example of how both too little and too much regularisation can be detrimental for interpretability. Image from VOC2007 test set; saliency is class transition of a P_{\downarrow} -optimal path.

the strongly regularised saliency is not only condensed to a single location, but also specifically to a corner of the image. Our interpretation is that this happens because it reduces the total variation (since $\frac{3}{4}$ of the gradient of the mask lies

outside of the image). And although the mask still contains enough of the dog’s head to keep the classification, its maximum lies misleadingly in front of the nose.

Adversarialness, or generally instability of the mask-classification interaction, can also be viewed in terms of a decision boundary that is fractal-like crinkled in the high-dimensional image space, such that small exclaves of a class domain may reach far closer to x_0 than the bulk of that domain. This suggests that it would help to evaluate for many different masks rather than gradient-optimising a single one. Particularly the RISE method [9] benefits from this, by evaluating the classifier for a whole large random selection of masks. The reference implementation of [5] also optimises (individually) multiple masks of different sizes. This does not so much avoid adversarial examples as average out their contributions, whereas stable, faithfully-salient masks tend to agree. (This still requires also dedicated mask regularity, as otherwise the adversarial contributions overwhelm and the result appears as mere noise.)

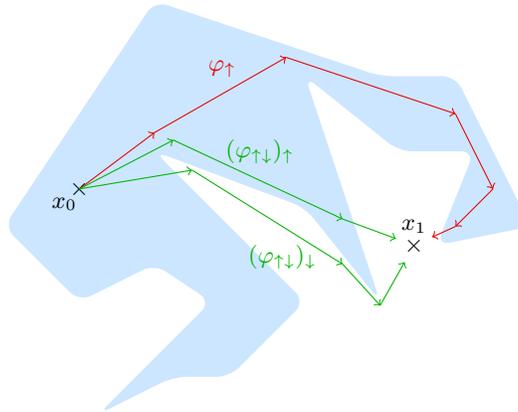


Fig. 4: Low-dimensional sketch of how irregularities in the decision boundary between two images can affect optimised paths. The P_{\uparrow} -optimised one takes a detour along a strong outlier (\approx adversarial) that allows it to stay almost completely in-domain, which causes the final section to come from a completely different direction in the end. The pair of $P_{\uparrow\downarrow}$ -optimised paths are still somewhat affected by outliers, but the pinching term causes them to mostly follow a more regular and consistent section of the boundary.

N.B.: this represents only very crudely the behaviour in real image classification applications, as inevitable with low-dimensional visualisations. In particular, the monotonicity condition is not represented at all here, and the pinching is here a symmetric L^2 -reduction, which is quite different from Equation 8.

It appears that the irregularity of the decision boundary is better described as thin outreaching folds rather than standalone islands (Figure 4), such that

even a monotone path can follow them up to an adversarial point before crossing the decision boundary very near x_0 . The scores P_{\uparrow} and P_{\downarrow} are particularly prone to procuring such paths, but we observe them also when optimising for P_{\downarrow} .

By comparison, $P_{\uparrow\downarrow}$ seems to be more reliable in practice. An intuitive reason is that the two paths have less possibility to simultaneously follow adversarial masks for the classes of both x_0 and x_1 , whilst also staying close to each other. That is certainly not inconceivable either, though.

7 Pointing game

We evaluate our saliency algorithm using the pointing game. This method was introduced in [18] and used, for instance, in [11][6]. It checks whether the maximum pixel of a saliency heatmap agrees with the location of a human-annotated⁴ object of the class of interest.

Assessments like the pointing game have their caveats for benchmarking saliency methods. One can for instance argue that the cases when the saliency points somewhere outside the bounding box are the most insightful ones, as they indicate that the classifier is using information from an unexpected part of the image (for instance, the background). Another caveat is that, if winning at the pointing game is the goal, a saliency method is only as good as its underlying classifier is. For these reasons, a pointing game score should not be considered as the predominant criterion for a good saliency method.

Nevertheless, it is reasonable to expect a useful saliency method to perform at least similarly well as the state of the art: if existing methods have proven capable of achieving high scores, this is after all indication that the classifier does to a significant degree base its decisions on spatially confined features of the real objects. Furthermore the pointing game provides a way of comparing the behaviour of different variations of a saliency method (such as different hyperparameters) somewhat more representatively than looking at individual image examples. Again, the best-scoring parameters are not necessarily the best for attribution purposes, but they are a reasonable starting point.

7.1 Heatmap reduction

The result of the methods introduced in this paper is one or multiple paths, whereas the pointing game expects a single heatmap. There are multiple ways of reducing to such a map:

⁴ The pointing game is generally used under the assumption that neither the classifier nor the saliency method have any direct training knowledge about the position annotations, i.e. it is not a test of how well a trained task generalizes but of an extrinsic notion of saliency.

Averaging. One can simply average over all the masks in a P_{\uparrow} -optimal path. This operation is (modulo a time renormalisation) left inverse to the pixel ablation of a saliency map (§ 4.2).

$$\overline{\varphi}_{\uparrow} := \int_0^1 \varphi(t) dt. \quad (9)$$

This works well in some cases, but the result can be disproportionately affected by low-discriminate contrasts of masks generated far from a decision boundary, which are unstable in a similar way to plain gradient methods.

Class transition. Taking the point of view that the decision boundary is what matters, one can seek the position where the path crosses it by tracking the classifier outputs along the path.

Empirically, this gives better results than averaging (both for the pointing game and, to our eyes, ease of interpretation), but it hinges on the assumption of there being a single boundary-crossing. In general, there may be multiple crossings, or the classifier might have a far more gradual transition, or (in case an explanation for a class different from the prediction for x_0 is sought) it might not cross a boundary at all. In our implementation, we therefore make a case distinction:

- If there exist t such that $F(\varphi(t))$ is dominated by the target class, then we select the largest of these t . In other words, we select the most confined mask that results in the classification of interest. Here (unlike the rest of the paper) we consider the full multi-class output of F , and by “dominate” we mean that the target class ranks higher than all others.
- If no such t exists, we select simply $\arg \max_t F(\varphi(t))$.

This may not be the best strategy in all applications, but it does guarantee always getting a result that can be compared in the pointing game. In critical applications it is likely better to discard paths that do not cross a boundary, and consult a different method in such a case.

Contrastive averaging. For the two paths optimizing $P_{\uparrow\downarrow}$, the property of interest is that they pinch the decision boundary between them. That means that for each t , the normal direction of the boundary is approximated by $\varphi_{\uparrow}(t) - \varphi_{\downarrow}(t)$ (at least coarsely, cf. Figure 4). This suggests averaging between these values, i.e.

$$\overline{\varphi}_{\uparrow\downarrow} := \int_0^1 (\varphi_{\uparrow}(t) - \varphi_{\downarrow}(t)) dt. \quad (10)$$

Indeed this appears to give comparatively good, stable results in practice. Our interpretation is that on any indiscriminate part of the path, the pinching tweak Equation 8 reduces $\varphi_{\uparrow}(t) - \varphi_{\downarrow}(t)$ so these parts do not contribute to the result like they would in Equation 9. The reason for this behaviour is that indiscriminate parts do not have a consistent F -gradient that would keep $\varphi_{\uparrow}(t)$ and $\varphi_{\downarrow}(t)$ apart during optimisation. On the other hand, stably-salient differences do keep them apart and therefore prevail in $\overline{\varphi}_{\uparrow\downarrow}$.

7.2 Results

We initially ran a custom implementation of the pointing game on individual classes (synsets) from the ImageNet dataset. In some cases even simple P_{\uparrow} -optimal paths perform well, e.g. 83% on the “Bee” synset with EfficientNet classifier, which is better than the result with saliency methods from the literature. These experiments turned however out to be not very representative: on larger and mixed datasets, we were not able to find hyperparameters that avoided high instability in the optimisation and consequently lower scores, especially in case of P_{\uparrow} .

For fair and representative comparison with the literature, we present here the result on a benchmark that was already used in [5]. Specifically, we used their TorchRay suite [16] to evaluate the saliency-result of our method (reduced to a heatmap), explaining the classifications by ResNet50 on the COCO14 validation dataset. We did not have the computational resources to run the whole set, so used the first 1000 images⁵ for a comparison to the literature state of the art of our best-scoring result (Table 1), which was in turn determined among variations of our method on the first 100 images (Table 2). For each image, a saliency is obtained for each of the annotated objects, so for example the 100 COCO images correspond to 310 optimised paths.

Method	VOC07 Test (All%/Diff%)	COCO14 Val (All%/Diff%)	Method	VOC07 Test (All%/Diff%)	COCO14 Val (All%/Diff%)
Ctr.	70.9/41.9	26.0/15.4	RISE	86.4/78.8	54.7/50.0
GCAM	90.5/80.4	57.1/49.2	GCAM	90.4/82.3	57.3/52.3
Ours	84.3/64.8	49.3/41.0	Extr	88.9/78.7	56.5/51.5

Table 1: *Left*: the highest-scoring results for the pointing game over 1000 images with ResNet50 classifier, for comparison with the state of the art. *Right*: excerpt from table 1 in [5], which contains the scores of more methods from the literature for the complete datasets.

The top scores are close to the state of the art, but do not quite reach the pointing accuracy of Grad-CAM, nor of extremal perturbation or RISE. It is not clear whether this is a result of fundamental limitations of our approach, of remaining stability problems that could be fixed with e.g. other regularisation

⁵ By “first 1000” we mean the 1000 images with the lowest ids. Note that the VOC and COCO sets are in random order, so that this should be a reasonably representative and reproducible selection. Comparing the score of Grad-CAM to the one on the full datasets confirms this.

The astute reader may notice that on the other hand, with only the first 100 images the results are systematically worse. This is less due to these images being more difficult, than artifact of the way the TorchRay benchmark gathers the results: specifically, it counts success rate for each class separately and averages in the end, but rates classes that are not even present in the smaller subset as 0% success.

approaches, or whether it is even a deficiency at all. Clearly the method does work in principle, and it is by construction faithful, so the somewhat higher mismatch rate could also be construed as higher sensitivity to aberrant or unstable behaviour of the classifier.

Method	opt.crit	intp.spc	ζ_{sat}	σ_{reguBlur}	postproc	<i>VOC07 Test</i> (All%/Diff%)	<i>COCO14 Val</i> (All%/Diff%)
Ctr.						71.4/36.6	26.5/11.2
GCAM						90.4/64.2	48.9/35.2
Abl.Path	P_{\uparrow}	blur-fade	0.8	8.0px		44.0/38.1	30.2/23.2
Abl.Path	P_{\downarrow}	blur-fade	0.8	2.0px		73.8/46.9	38.8/26.4
Abl.Path	P_{\downarrow}	blur-fade	0.8	7.0px		52.4/40.0	32.5/23.1
Abl.Path	P_{\downarrow}	blur-fade	0.8	7.0px	window	76.4/48.9	40.6/26.7
Abl.Path	P_{\downarrow}	blur-fade	0	7.0px	window	80.5/46.1	46.2/31.0
Abl.Path	$P_{\uparrow\downarrow}$	blur-fade	0.8	7.0px		72.2/47.5	48.3/34.8
Abl.Path	$P_{\uparrow\downarrow}$	pyramid	0.8	8.0px		75.2/47.1	39.5/26.5

Table 2: Some of our results for the pointing game over 100 images with ResNet50 classifier.

Different variations of our methods also perform quite differently. We cannot describe all the observations that could be made from these experiments here, nor is this necessarily useful (many of the trends here likely do not generalize to other data), but a selection that may be noteworthy:

- The simple single-path retaining (or dissipating) methods compete badly. See §6 for some possible explanations.
- The boundary-straddle method performs best on the COCO dataset, and also relatively good on the difficult subset of VOC. We propose that this is typically the best of our methods, though in particular on the simple subset of VOC its results are quite disappointing.
- The contrastive method performs well in particular on the simple subset of VOC, but only with very particular regularisation settings; see §7.3. With e.g. strong blurring and saturation but no windowing, it may perform worse than even the trivial center method, evidently an artifact of the effect shown in Figure 3.

7.3 Hyperparameter choice

Most of the saliency methods from the literature have some hyperparameters⁶, as does ours. The ideal choice of these parameters is little discussed. Unlike when training a machine learning model, there is no objective on which this choice should unambiguously be based, but it appears that several authors have

⁶ The authors in [5] emphasize that their method avoids hyperparameters, yet their examples rely on no fewer than 5 hard-coded number constants.

used the pointing game for this purpose. Apart from the aforementioned caveats, this also has the problem that the required position-annotations are simply not available for most applications.

In our case, there *is* an additional score with a clear meaning available: the ablation path score. And though it is evidently not the case that the hyperparameters leading to the highest ablation score give the best saliency results, we propose that studying only the ablation score can nevertheless provide some guidance for a good choice.

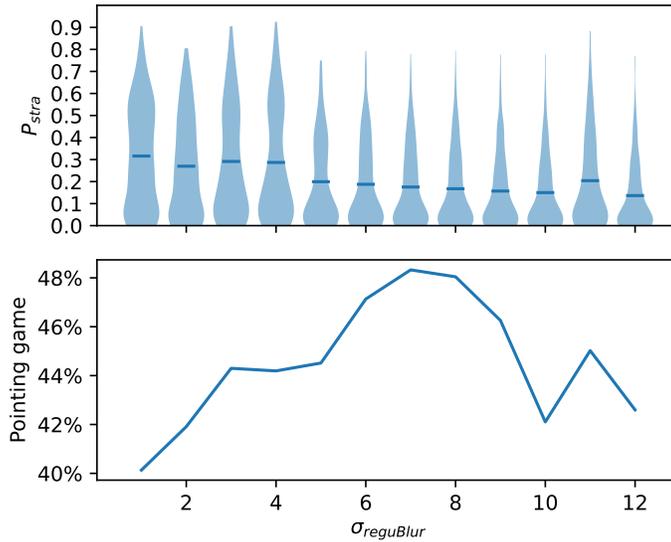


Fig. 5: Dependence on the size of the regularization filter, for both the distribution of boundary-straddling ablation-path scores and the pointing game score (evaluated with contrastive averaging as per §7.1). Based on ResNet50 and 100 images from the COCO14 dataset.

Specifically, the regularisation parameter is responsible for avoiding adversarial masks, which can be identified by a large population of paths with very high score. Observe in Figure 5 that the histograms at low σ_{reguBlur} have an upper bulge ($P_{\uparrow\downarrow} > 0.5$). After σ_{reguBlur} has been increased to a size of 5 pixels, the adversarial population vanishes, and accordingly the pointing game score rises towards its maximum at $\sigma_{\text{reguBlur}} = 7$.

Because the path score is a property without application-specific dimension, this phenomenon can also be expected in applications where very different regularisation is required compared the image classification examples here. This is therefore a possible criterion for hyperparameter choice when the pointing game or an analog is not available. Some care needs to be taken though: Figure 6 shows

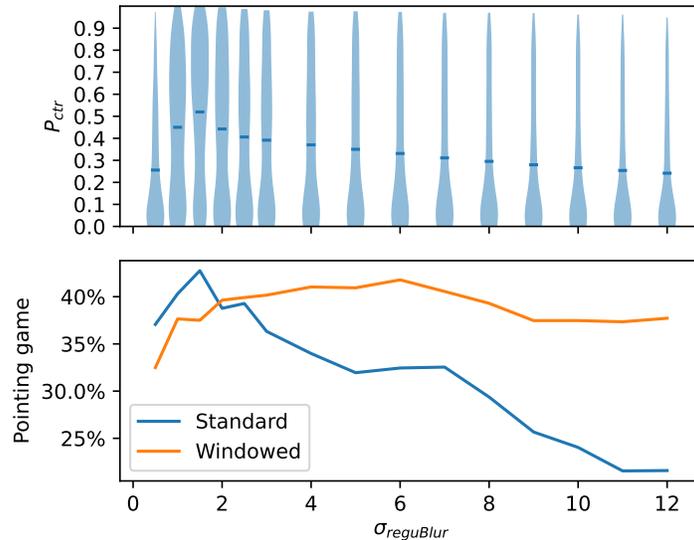


Fig. 6: Like Figure 5, but with paths optimised for the contrastive score. Pointing game evaluated on both standard class-transition masks (§ 7.1), and with a boundary-suppressing window applied.

an example where the pointing game with unmodified masks has a steadily decreasing score as the regularization is increased. We have identified the problem as stemming from the tendency of the regularization to push the argmax to the image boundary (cf. Figure 3). This particular effect can be prevented by post-processing the masks with a window that suppresses the boundaries⁷, in which case the rule suggested above is valid again.

8 Conclusion

We demonstrated that the ablation path formalism provides a usable saliency method that combines ideas from several previous methods within a single mathematical framework. The ablation path method can stand in for each of these methods to some extent, and is with suitable parameter choices also able to produce results that score similarly well in the pointing game.

This is a nontrivial result, because these methods appear quite different in their original formulations. And even though ours has strong similarity to [6] / [5], it was a priori not obvious that the restriction to a path instead of individual

⁷ This can be interpreted as applying prior knowledge of the location of objects in the dataset. However, there *are* also examples of objects close to the boundary. The window post-processing prevents these from being properly localised, which is why the top pointing-game score is still lower.

masks would still leave the optimisation problem solvable in practice. Indeed, for some inputs our method still struggles to converge on a human-reasonable explanation, even when other methods accomplish this. It is possible that in some cases there simply exist no paths that the classifier can follow in a well-behaved way. But for most of the examples we tested on, this does not seem to be a fundamental issue.

The main practical advantage of an ablation path, which is most evident when interactively browsing through it and tracking the exact classifier response, is the added information: unlike each of the previous methods, an ablation path offers a whole sequence of fine-grained changes to an input image. It thus offers a more thorough insight to the classification, while still ensuring the explanations form a consistent picture thanks to the monotonicity condition. Because each point in a path is associated with a concrete input to the network whose result can directly be inspected, we argue the method is *faithful* [17], whilst also being easy and intuitive to use. A caveat is that without suitable parameters (in particular regularisation), the explanations may highlight only the adversarial behaviour of a classifier, or even respond to spoofing by a classifier designed to detect the artificial inputs (“Volkswagening”). This possibility is to our knowledge common to all black-box saliency methods, so their use in critical applications should be considered carefully[10].

A disadvantage of our specific approach, in addition to the stability issues, is the rather high computational effort. A path requires many (ca. 50) optimisation steps, each of which require several classifier evaluations (ca. 20).

The method is perhaps best used in tandem with another one, for example Grad-CAM which is fast and stable but lacks the possibility of assessing faithfulness.

References

1. Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., Kim, B.: Sanity checks for saliency maps. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) *Advances in Neural Information Processing Systems* 31, pp. 9505–9515. Curran Associates, Inc. (2018), <http://papers.nips.cc/paper/8160-sanity-checks-for-saliency-maps.pdf>
2. Ancona, M., Ceolini, E., Öztireli, C., Gross, M.: Towards better understanding of gradient-based attribution methods for deep neural networks. *CoRR* (2017)
3. Chockler, H., Kroening, D., Sun, Y.: Explanations for occluded images. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. pp. 1234–1243 (October 2021)
4. Dabkowski, P., Gal, Y.: Real time image saliency for black box classifiers. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*. vol. 30. Curran Associates, Inc. (2017), <https://proceedings.neurips.cc/paper/2017/file/0060ef47b12160b9198302ebdb144dcf-Paper.pdf>
5. Fong, R., Patrick, M., Vedaldi, A.: Understanding deep networks via extremal perturbations and smooth masks pp. 2950–2958 (2019)
6. Fong, R.C., Vedaldi, A.: Interpretable explanations of black boxes by meaningful perturbation (Oct 2017)
7. Kindermans, P.J., Hooker, S., Adebayo, J., Alber, M., Schütt, K.T., Dähne, S., Erhan, D., Kim, B.: The (un)reliability of saliency methods pp. 267–280 (2019). https://doi.org/10.1007/978-3-030-28954-6_14, https://doi.org/10.1007/978-3-030-28954-6_14
8. Koh, P.W., Liang, P.: Understanding black-box predictions via influence functions **70**, 1885–1894 (06–11 Aug 2017), <https://proceedings.mlr.press/v70/koh17a.html>
9. Petsiuk, V., Das, A., Saenko, K.: Rise: Randomized input sampling for explanation of black-box models. *CoRR* (2018)
10. Rudin, C.: Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence* **1**(5), 206–215 (may 2019). <https://doi.org/10.1038/s42256-019-0048-x>, <https://doi.org/10.1038/s42256-019-0048-x>
11. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization (Oct 2017)
12. Simonyan, K., Vedaldi, A., Zisserman, A.: Deep inside convolutional networks: Visualising image classification models and saliency maps. *CoRR* (2013)
13. Sturmfels, P., Lundberg, S., Lee, S.I.: Visualizing the impact of feature attribution baselines. *Distill* (2020). <https://doi.org/10.23915/distill.00022>, <https://distill.pub/2020/attribution-baselines>
14. Sundararajan, M., Taly, A., Yan, Q.: Axiomatic attribution for deep networks **70**, 3319–3328 (06–11 Aug 2017), <https://proceedings.mlr.press/v70/sundararajan17a.html>
15. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks (2014)
16. Vedaldi, A.: Understanding deep networks via extremal perturbations and smooth masks. <https://github.com/facebookresearch/TorchRay> (2019)

17. Weller, A.: Transparency: Motivations and challenges. In: Samek, W., Montavon, G., Vedaldi, A., Hansen, L.K., Müller, K.R. (eds.) Explainable AI: interpreting, explaining and visualizing deep learning, vol. 11700, chap. 2. Springer Nature (2019)
18. Zhang, J., Bargal, S.A., Lin, Z., Brandt, J., Shen, X., Sclaroff, S.: Top-down neural attention by excitation backprop. International Journal of Computer Vision **126**(10), 1084–1102 (2017). <https://doi.org/10.1007/s11263-017-1059-x>

Appendix A Canonical Time Reparametrisation

Proof (Proof of ??). The function $m: [0, 1] \rightarrow \mathbb{R}$ defined by $m(t) := \int_{\Omega} \varphi(t)$ is increasing and goes from zero to one (since we assume that $\int_{\Omega} 1 = 1$).

Note first that if $m(t_1) = m(t_2)$, then $\varphi(t_1) = \varphi(t_2)$ from the monotonicity property. Indeed, supposing for instance that $t_1 \leq t_2$, and defining the element $\theta := \varphi(t_2) - \varphi(t_1)$ we see that on the one hand $\int_{\Omega} \theta = 0$, on the other hand, $\theta \geq 0$, so $\theta = 0$ and thus $\varphi(t_1) = \varphi(t_2)$.

Now, define $M := m([0, 1]) = \{s \in [0, 1] \mid \exists t \in [0, 1] m(t) = s\}$. Pick $s \in [0, 1]$.

If $s \in M$ we define $\psi(s) := \varphi(t)$ where $m(t) = s$ (and this does not depend on which t fulfills $m(t) = s$ from what we said above). We remark that $\int_{\Omega} \psi(s) = \int_{\Omega} \varphi(t) = m(t) = s$.

Now suppose that $s \notin M$. Define $s_1 := \sup(M \cap [0, s])$ and $s_2 := \inf(M \cap [s, 1])$ (neither set are empty since $0 \in M$ and $1 \in M$). Since $s_1 \in M$ and $s_2 \in M$, there are $t_1 \in [0, 1]$ and $t_2 \in [0, 1]$ such that $m(t_1) = s_1$ and $m(t_2) = s_2$. Finally define $\psi(s) := \varphi(t_1) + (s - s_1) \frac{\varphi(t_2) - \varphi(t_1)}{s_2 - s_1}$. In this case, $\int_{\Omega} \psi(s) = m(t_1) + (s - s_1) \frac{m(t_2) - m(t_1)}{s_2 - s_1} = s$. The path ψ constructed this way is still monotone, and it has the constant speed property, so it is an ablation path.

Appendix B \mathcal{L}^{∞} -optimal Monotonicity Projection

The algorithm proposed in § Appendix C for optimising monotone paths uses updates that can locally introduce nonmonotonicity in the candidate $\hat{\varphi}_1$, so that it is necessary to project back onto a monotone path φ_1 . The following routine⁸ performs such a projection in a way that is optimal in the sense of minimising the \mathcal{L}^{∞} -distance⁹, i.e.,

$$\sup_t |\varphi_1(t, \mathbf{r}) - \hat{\varphi}_1(t, \mathbf{r})| \leq \sup_t |\vartheta(t, \mathbf{r}) - \hat{\varphi}_1(t, \mathbf{r})|$$

for all $\mathbf{r} \in \Omega$ and any other monotone path ϑ .

The algorithm works separately for each \mathbf{r} , i.e., we express it as operating simply on continuous functions $p: [0, 1] \rightarrow \mathbb{R}$. The final step effectively *flattens out*, in a minimal way, any region in which the function was decreasing.

⁸ It is easy to come up with other algorithms for monotonising a (discretised) function.

One could simply *sort the array*, but that is not optimal with respect to any of the usual function norms; or clip the derivatives to be nonnegative and then rescale the entire function, but that is not robust against noise perturbations.

⁹ Note that the optimum is not necessarily unique.

Algorithm 1 Make a function $[0, 1] \rightarrow \mathbb{R}$ nondecreasing

```

 $\cup_i [l_i, r_i] \leftarrow \{t \in [0, 1] \mid p'(t) \leq 0\}$  ▷ Union of intervals where  $p$  decreases
for  $i$  do
   $m_i \leftarrow \frac{p(l_i) + p(r_i)}{2}$ 
   $l_i \leftarrow \max\{t \in [r_{i-1}, l_i] \mid p(t) \leq m_i\}$ 
   $r_i \leftarrow \min\{t \in [r_i, l_{i+1}] \mid p(t) \geq m_i\}$ 
end for
for  $i, j$  do
  if  $[l_i, r_i] \cap [l_j, r_j] \neq \emptyset$  then
    if  $m_j < m_i$ , merge the intervals and recompute  $m$  as the new center
  end if
end for
return  $t \mapsto \begin{cases} p(t) & \text{if } t \notin \cup_i [l_i, r_i] \\ m_i & \text{if } t \in [l_i, r_i] \end{cases}$ 

```

In practice, this algorithm is executed not on continuous functions but on a PCM-discretised representation; this changes nothing about the algorithm except that instead as real numbers, l, r and t are represented by integral indices.

Appendix C Path Optimisation Algorithm

As said in § 5, our optimisation algorithm is essentially gradient descent of a path φ : it repeatedly seeks the direction within the space of all paths that (first ignoring the monotonicity constraint) would affect the largest increase to $P(\varphi)$ as per § 4, for any of the defined score functions. Algorithm 2 shows the details of how this is done in presence of our constraints. In case of $P_{\uparrow\downarrow}$, the state φ is understood to consist of the two paths φ_{\uparrow} and φ_{\downarrow} .

As discussed before, the use of a gradient requires a metric to obtain a vector from the covector-differential, which could be either the implicit ℓ^2 metric on the discretised representation (pixels), or a more physical kernel/filter-based metric. In the present work, we base this on the regularisation filter.

Unlike with the monotonisation condition, the update can easily be made to preserve speed-constness by construction, by projecting for each t the gradient \mathbf{g} on the sub-tangent-space of zero change to $\int_{\Omega} \varphi(t)$, by subtracting the constant function times $\int_{\Omega} \mathbf{g}(t)$. Note this requires the measure of Ω to be normalised, or else considered at this point.

Then we apply these gradients time-wise as updates to the path, using a scalar product in the channel-space to obtain the best direction for φ itself (as opposed to the corresponding image composit $x_{\varphi,t}$).

The learning rate γ can be chosen in different ways. What we found to work best is to normalise the step size in a \mathcal{L}^{∞} sense, such that the strongest-affected pixel in the mask experiences a change of at most 0.7 per step. This is small

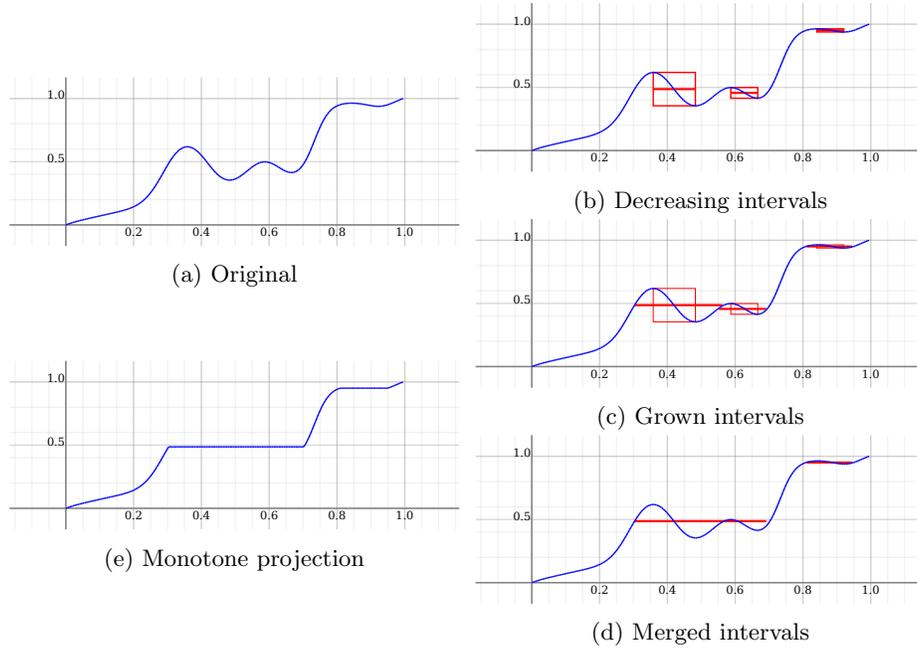


Fig. 7: Example view of the monotonicisation algorithm in practice. (a) contains decreasing intervals, which have been localised in (b). For each interval, the centerline is then extended to meet the original path non-decreasingly (c). In some cases, this will cause intervals overlapping; in this case merge them to a single interval and re-grow from the corresponding centerline (d). Finally, replace the path on the intervals with their centerline (e).

enough to avoid excessively violating the constraint, but not so small to make the algorithm unnecessarily slow.

Appendix D Baseline choice

The baseline image is prominently present in the input for much of the ablation path, and it is therefore evident that it will have a significant impact on the saliency. In line with previous work, we opted for a blurred baseline for the examples in the main paper, but even then there is still considerable freedom in the choice of blurring filter. Figure 8 shows two examples, where the result is not fundamentally, but still notably different.

Algorithm 2 Projected Gradient Descent

```

1:  $\varphi \leftarrow ((t, \mathbf{r}) \mapsto t)$  ▷ Start with linear-interpolation path
2: while  $\varphi$  is not sufficiently saturated do
3:   for  $t$  in  $[0, 1]$  do
4:      $x_{\varphi,t} := (1 - \varphi(t)) x_0 + \varphi(t) x_1$ 
5:     compute  $F(x_{\varphi,t})$  with gradient  $\mathbf{g} := \nabla F(x_{\varphi,t})$ 
6:     let  $\hat{\mathbf{g}} := \mathbf{g} - \int_{\Omega} \mathbf{g}$  ▷ ensure  $\hat{\mathbf{g}}$  does not affect mass of  $\varphi(t)$ 
7:     update  $\varphi(t, \mathbf{r}) \leftarrow \varphi(t, \mathbf{r}) - \gamma \langle \hat{\mathbf{g}}(\mathbf{r}) \mid |x_1 - x_0 \rangle$ , for  $\mathbf{r}$  in  $\Omega$ 
8:     (optional) apply a regularisation filter to  $\varphi(t)$ 
9:   end for
10:  (optional) adjust learning rate  $\gamma$  according to size of the actual step performed
11:  (optional) apply saturation to  $\varphi$  (§ 5.2)
12:  (optional) apply pinching to the paths  $\varphi_{\uparrow}, \varphi_{\downarrow}$  (§ 5.3)
13:  for  $\mathbf{r}$  in  $\Omega$  do
14:    re-monotonise  $t \mapsto \varphi(t, \mathbf{r})$ , using Algorithm 1
15:  end for
16:  clamp  $\varphi(t, \mathbf{r})$  to  $[0, 1]$  everywhere
17:  re-parametrise  $\varphi$ , such that  $\int_{\Omega} \varphi(t) = t$  for all  $t$  (using § Appendix A)
18: end while

```

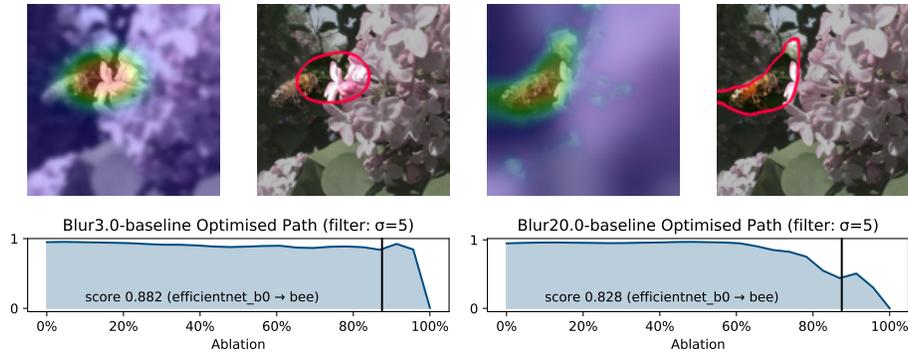


Fig. 8: An example of paths obtained with different-size blur baselines.