

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Explaining Image Classifiers with Multiscale Directional Filters

Anonymous CVPR submission

Paper ID ****

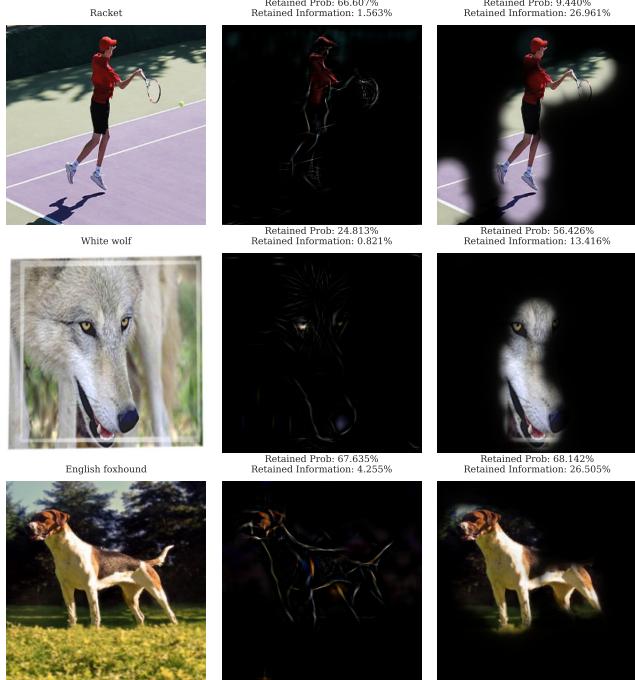


Figure 1. Left column: Imagenet images with VGG-19 [34] prediction. Middle column: Our new ShearletX explanation. Right column: Smooth pixel mask explanation from Fong et al. [11]. Our mask explanation method can delete irrelevant features such as background elements and texture when they are deemed irrelevant without producing explanation artifacts.

Abstract

We present a novel mask explanation method *ShearletX* that can localize and delete irrelevant features that current mask explanation methods cannot due to their inflexible smoothness constraints that protect against undesirable artifacts. Our method optimizes a sparse mask in the shearlet representation of an image to delete as much spatial information as possible while retaining the classifier’s prediction. We show theoretically and empirically that our method manages to guard against explanation artifacts, while retaining the ability localize the relevant features remarkably well.

1. Introduction

Modern image classifiers are known to be difficult to interpret and explain. Saliency maps comprise a well-established explainability tool that highlights important image regions for the classifier and helps interpret classification decisions. An important saliency approach frames saliency map computation as an optimization problem over masks [7, 9, 11, 12, 16, 19, 26]. The explanation mask is optimized to keep only parts of the image that suffice to retain the classification decision. The main advantage of the mask approach over other saliency map methods [6, 32, 36–39, 41, 42] is that the explanation mask is found by optimizing an objective that intuitively aligns with what is defined as relevance. However, Fong and Vedaldi [12] showed that an unregularized explanation mask is very susceptible to explanation artifacts and unreliable. Therefore, current practice [7, 11, 12] heavily regularizes and parameterizes the explanation masks to be smooth. Although smooth explanation masks can communicate useful explanatory information by roughly localizing objects, they are heavily constrained in their ability to delete irrelevant features in the image due to possible smoothness violations. As a result many details such as background elements, texture, and other spatially localized features cannot be deleted although they may be irrelevant to the classifier.

An optimal mask explanation method should be both resistant to artifacts and capable of deleting all irrelevant details. We present a novel method that achieves this by optimizing a sparse mask in the shearlet representation [21] of an image to delete as much spatial information as possible while retaining the classifier’s prediction. We theoretically and empirically show that sparsely masking in the shearlet representation instead of the pixel representation guards against artifacts, which we characterize as artificial edges. Due to shearlets’ ability to efficiently encode anisotropic features in images, we can localize the relevant image parts extremely well without producing explanation artifacts. Figure 1 illustrates the advantage of ShearletX over smooth pixel mask explanations with an example. For instance, ShearletX for the wolf image shows that the classifier has learned the core characteristics of a wolf: the wolf’s

108 eyes and distinctive snout. Such anisotropic details cannot
 109 be efficiently represented by a smooth explanation mask,
 110 which picks up details such as teeth and fur to obey the
 111 smoothness constraints. In the example, our novel method
 112 can tell us that 1) the classifier can correctly predict the wolf
 113 from the wolf’s eyes and snout shape. 2) Adding the fur and
 114 the teeth in the smooth pixel mask explanation double the
 115 class score. The examples in Figure 1 further reveal that
 116 a good ImageNet classifier only needs a rough sketch of
 117 the class instance to correctly classify many images, simi-
 118 lar to humans. Such new explanatory insights motivate
 119 why ShearletX will be an important new explanation tool
 120 for practitioners.
 121

2. Related Work

122 The explainability field has experienced a surge in re-
 123 search activity over the last decade, due to the obvious soci-
 124 etal need to explain machine learning models. We focus on
 125 explainability aspects of image classifiers, where saliency
 126 maps offer an important and useful way of understanding a
 127 classifier’s prediction. The community has also introduced
 128 other tools, such as concept-based methods [18] and inher-
 129 ently interpretable architectures [8]. Nevertheless, saliency
 130 map methods comprise the most common tools for explain-
 131 ing image classifiers. In the following, we review current
 132 saliency map methods.
 133

Pixel Attribution Methods

134 Many of the first saliency map methods assign a relevance
 135 score to each pixel indicating its relevance for the predic-
 136 tion. Such methods include Gradient Saliency [33], LRP
 137 [6], Guided Backprop [37] and Integrated Gradients [38].
 138 Although the fidelity of such explanations can be verified
 139 quantitatively, *e.g.* with deletion and insertion curves [6,29],
 140 they are to a large degree heuristic and not directly optimiz-
 141 ing for a well defined notion of relevance. Moreover, pixel
 142 attributions tend to be relatively jittery and have been even
 143 shown to be susceptible to adversarial attacks [10]. Other
 144 methods, such as LIME [31] and SHAP [25], form coalitions
 145 of pixels by first segmenting the image into superpix-
 146 els and assigning a relevance score to each superpixel. Not
 147 only has research exposed various vulnerabilities of LIME
 148 and SHAP [35] but superpixels are suboptimal due to their
 149 inflexibility. Within a superpixel one cannot say what is rel-
 150 evant and what not.
 151

Pixel Mask Explanations

152 Mask-based explanations do not attribute individual relevance
 153 scores to (super)pixels but rather optimize a mask to
 154 delete as much information of the image as possible while
 155 maximizing the probability score for the classification. The
 156

157 advantage of this approach is that one optimizes for a natu-
 158 ral interpretability objective that can be quickly validated in
 159 two steps: (1) Determining which and how much informa-
 160 tion was deleted by the mask (2) Computing the probability
 161 score after masking the image. The explanation mask is
 162 found as a solution to the optimization problem
 163

$$\max_{m \in \mathcal{M}} \mathbb{E}_{u \sim \nu} [\Phi(x \odot m + (1 - m) \odot u)] + \lambda \cdot \|m\|_1, \quad (1)$$

164 where
 165

Wavelet Mask Explanations

3. WaveletX

3.1. Wavelets for Images

3.2. Method

3.3. Theory

4. ShearletX

4.1. Shearlets for Images

4.2. Method

4.3. Theory

5. Experiments

5.1. Artificial Edges and Hallucinations

5.2. Conciseness and Preciseness

6. Limitations

7. Conclusion

8. Introduction

199 Modern image classifiers are known to be difficult to
 200 interpret and explain. Saliency maps comprise a well-
 201 established explainability tool that highlights important im-
 202 age regions for the classifier and helps interpret classifi-
 203 cation decisions. An important saliency approach [7, 9,
 204 11, 12, 16, 19, 26] frames saliency map computation as an
 205 optimization problem over masks. The explanation mask
 206 is optimized to keep only parts of the image that suffice
 207 to retain the classification decision. The main advantage
 208 of the mask approach over other saliency map methods
 209 [6, 32, 36–39, 41, 42] is that the explanation mask is found
 210 by optimizing an objective that intuitively aligns with what
 211 is defined as relevance. Mask explanations need to be ro-
 212 bust against corruptions and discriminate between unneces-
 213 sary and necessary features of the image. Fong and Vedaldi
 214 showed that pixel masks are susceptible to corruptions and
 215 propose to mitigate this with TV regularization to enforce

smooth masks. However, we identify a problem in TV regularized masks as the TV penalty compromises the ability to discriminate between unnecessary and necessary features of an image, such as texture and shape. A TV regularized mask can be thought of as a window that locates the relevant part fitting the window. However, features inside the window may not be important and difficult to delete due to the TV penalty. To address this, we revisit a recent technique proposed by Kolek et al. that gives explanations by masking in the wavelet domain (WaveletX). We first prove experimentally and theoretically that WaveletX is robust against explanation corruptions, similar to TV regularized pixel masks. Then we show that WaveletX is better suited to discriminate relevant and irrelevant features and thereby avoids the problems of TV regularized pixel masks. Moreover, we identify and address a significant problem of WaveletX. By adding a crucial spatial regularization term, we get rid of blurry backgrounds that are sometimes difficult to interpret, making WaveletX unambiguous and more interpretable. Moreover, we address the fact that Wavelets are not an optimal representation system for piece-wise smooth images and successfully apply the mask explanation approach in the optimal shearlet representation system. Results show that ShearletX can represent edges better in the explanation. Lastly, we also address an overlooked bias of in-distribution perturbations, realized as an inpainting GAN, that were claimed to be most appropriate for mask explanations [7, 16]. An explanation with inpainting GAN perturbations may reflect not what the classifier needs to see for the correct classification but rather what the inpainting GAN needs to infer the class relevant features. Overall, we are the first to identify all three biases of pixel mask explanations: 1) explanation corruption due to artificial edges, 2) inability to discriminate close relevant and irrelevant features due to TV penalty and 3) bias of an inpainting GAN perturbation. Our improved WaveletX and ShearletX are theoretically motivated and well-suited to mitigate all three biases.

9. Related Work

10. Mask Explanations

We model an image as a map $x : \Omega \rightarrow [0, 1]$, where Ω is the image domain. For a digital RGB image with $h \times w$ pixels, we have $\Omega = (\{1, \dots, h\} \times \{1, \dots, w\})^3$. For $K \in \mathbb{N}$ labels, we denote the K -simplex, the space of all probability distributions over K labels, as $\Delta^{(K)}$. An image classifier $\Phi : [0, 1]^\Omega \rightarrow \Delta^{(K)}$ is a map from image space to the label probability space. We denote a binary mask on an image as a map $m : \Omega \rightarrow \{0, 1\}$. Let $u : \Omega \rightarrow [0, 1]$ take random values from a probability distribution $u \sim \mathcal{V}$. We adopt the mathematical interpretation of explanation masks as solutions to a rate-distortion optimization problem as in

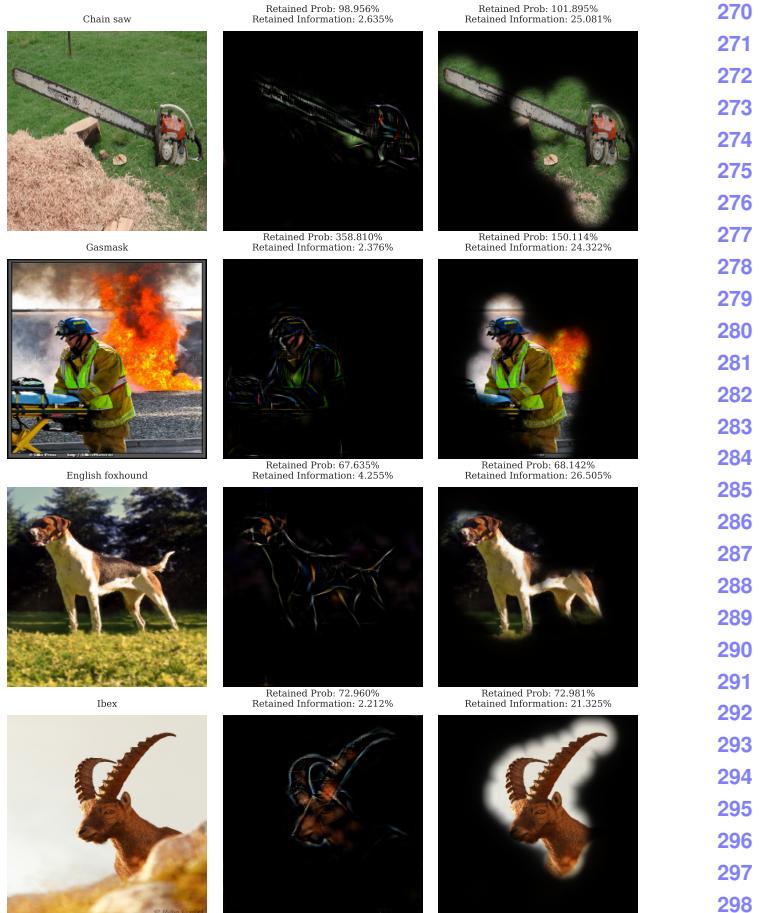


Figure 2. Caption

Kolek et al. in [20]. An optimal explanation mask m^* with sparsity level $s \in \mathbb{N}$ and prior $\mathcal{M} \subset [0, 1]^\Omega$ is defined as a solution to

$$\begin{aligned} \min_{m: \Omega \rightarrow \{0, 1\}} \quad & \mathbb{E}_{u \sim \mathcal{V}} \left[d(\Phi(x \odot m + (1 - m) \odot u), \Phi(x)) \right] \\ \text{subject to} \quad & \|m\|_0 \leq s \\ & m \in \mathcal{M}, \end{aligned} \quad (2)$$

where \odot denotes pointwise multiplication and $1 - m$ is pointwise subtraction by abuse of notation. The prior set \mathcal{M} can accommodate certain smoothness conditions such as low total variation and the $d(\cdot, \cdot)$ can be the ℓ_2 norm or set so that the minimization is equivalent to maximization of the predicted label. In practice, the discrete optimization problem is infeasible and the optimal explanation mask is approximated by solving a relaxed optimization problem

$$\min_{m: \Omega \rightarrow [0, 1]} \quad - \mathbb{E}_{u \sim \mathcal{V}} \left[\Phi(x \odot m + (1 - m) \odot u) \right] \quad (3)$$

$$+ \lambda_1 \cdot \|m\|_1 + \lambda_2 \cdot R_{\mathcal{M}}(m), \quad (4)$$

where $R_{\mathcal{M}}(m)$ is typically chosen as a TV penalty to enforce mask smoothness. The optimization is typically performed with projected SGD over the mask m [7, 11, 12, 20]. The perturbation distribution \mathcal{V} has priorly been chosen as Gaussian [12], blurring [12], constant [12], and inpainting GAN [7, 16].

10.1. CartoonX

Kolek et al. [19] applied the mask explanation framework in the wavelet domain of images instead of standard pixel space coined their method *CartoonX*. Since piece-wise smooth images are sparsely represented by wavelets [28], sparsely masking in the wavelet domain yields a piece-wise smooth (cartoon-like) images. In short, *CartoonX* was designed to extract the relevant piece-wise smooth part of an image.

Wavelets for Images

A wavelet $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a wave-like oscillation and the cornerstone of the *wavelet transform*

$$\mathcal{W}\{f\}(a, b) = \int_{\mathbb{R}^2} f(u) \frac{1}{a} \psi\left(\frac{u-b}{a}\right) du \quad (5)$$

of a function $f \in L^2(\mathbb{R}^2)$. The wavelet coefficient $\mathcal{W}\{f\}(a, b)$ extracts time-frequency information of the signal f at position $b \in \mathbb{R}^2$ and scale $a > 0$ (at a frequency proportional to $1/a > 0$). Three suitably chosen mother wavelets $\psi^1, \psi^2, \psi^3 \in L^2(\mathbb{R}^2)$ with dyadic dilations and translations yields an orthonormal basis

$$\left\{ \psi_{j,n}^k := \frac{1}{2^j} \psi^k \left(\frac{\cdot - 2^j n}{2^j} \right) \right\}_{j \in \mathbb{Z}, n \in \mathbb{Z}^2, 1 \leq k \leq 3} \quad (6)$$

of finite energy functions $L^2(\mathbb{R}^2)$. The three indices $k = 1, 2, 3$ correspond to directions (typically vertical, horizontal, and diagonal). A wavelet coefficient $\langle f, \psi_{j,n}^k \rangle$ has high amplitude if the function f has sharp transitions over the support of $\psi_{j,n}^k$. Suitable choices of a (mother) wavelets $\psi^1, \psi^2, \psi^3 \in L^2(\mathbb{R}^2)$, paired with an appropriate scaling function $\phi \in L^2(\mathbb{R}^2)$, define a multiresolution approximation. From the mother wavelet and scaling function one can construct the subspaces

$$V_j := \text{span}\{\phi_{j,n} \mid n \in \mathbb{Z}^2\}, \quad \phi_{j,n} := 2^{-j/2} \phi(2^{-j} \cdot -n) \\ W_j := \text{span}\{\psi_{j,n}^k \mid n \in \mathbb{Z}^2, 1 \leq k \leq 3\}, \quad (7)$$

which satisfy $V_j \subset V_{j-1}$, $V_j \oplus W_j = V_{j-1}$, and $L^2(\mathbb{R}^2) = \bigoplus_{j=-\infty}^J W_j \oplus V_J$. Therefore, for all $J \in \mathbb{Z}$, any finite energy signal decomposes into

$$f = \sum_{n \in \mathbb{Z}^2} a_n \phi_{J,n} + \sum_{1 \leq k \leq 3} \sum_{i \leq J} d_{j,n}^k \psi_{j,n}^k, \quad (8)$$

where $a_n = \langle f, \phi_{J,n} \rangle$ and $d_{j,n}^k = \langle f, \psi_{j,n}^k \rangle$ are the approximation coefficients at scale J and detail coefficients at scale $j-1$ respectively. In practice, images are discrete signals $x[n_1, n_2]$ with pixel values at discrete positions $n = (n_1, n_2) \in \mathbb{Z}^2$ but they can be associated with a function $f \in L^2(\mathbb{R}^2)$ approximated at a certain scale 2^L . The discrete wavelet transform of an image x then computes an invertible wavelet image representation:

$$\text{DWT}(x) = \left\{ a_{J,n} \right\}_n \cup \left\{ d_{j,n}^1, d_{j,n}^2, d_{j,n}^3 \right\}_{L < j \leq J, n} \quad (9)$$

corresponding to discretely sampled approximation and detail coefficients of f .

In CartoonX, the explanation mask is optimized on the wavelet image representation. A wavelet mask m is optimized with projected SGD to minimize

$$\begin{aligned} \min_m \quad & -\mathbb{E}_{u \sim \mathcal{V}} \left[\Phi(\text{DWT}^{-1}(m \odot \text{DWT}(x) + (1-m) \odot u)) \right] \\ & + \lambda_1 \cdot \|m\|_1, \end{aligned} \tag{10}$$

where $\lambda_1 > 0$ is a hyperparameter controlling the sparsity of the mask. The sparse wavelet mask is then applied to the wavelet representation and subsequently inverted to yield an approximately piece-wise smooth image $DWT^{-1}(m \odot DWT(x))$ that is sufficient to retain the classification decision. Kolek et al. [19] explain that fine details, that are still visible in the CartoonX explanation are relevant and fine details that are blurred or blacked out are irrelevant because it is expensive for the wavelet mask to retain fine details and cheap to keep smooth (very blurry) areas. The consequence is that the CartoonX explanation is harder to interpret when smooth image regions are not blacked out. Was the smooth region relevant or was it just inexpensive for the wavelet mask to keep the smooth region (see Figure Fig. 3 for an example)? In the following we propose to modify CartoonX to solve this problem and enforce regions to be black if they are irrelevant. This makes CartoonX easier to interpret and yields a novel kind of explanation method.

Shearlets for Images

The shearlet transform, which was introduced in [14], is based on applying translation, anisotropic dilation, and shearing to generator functions. To dilate and shear a function, we define the following three matrices:

$$A_a := \begin{pmatrix} a & 0 \\ 0 & \sqrt{a} \end{pmatrix}, \quad \tilde{A}_a := \begin{pmatrix} \sqrt{a} & 0 \\ 0 & a \end{pmatrix}, \quad S_s := \begin{pmatrix} 1 & s \\ 0 & 1 \end{pmatrix},$$

where $s, a \in \text{Real}$. Next, given $(a, s, t) \in \mathbb{R}_+ \times \mathbb{R} \times \mathbb{R}^2$, $\psi \in L^2(\mathbb{R}^2)$, and $x \in \mathbb{R}^2$, we define

$$\begin{aligned}\psi_{a,s,t,1}(x) &:= a^{-\frac{3}{4}} \psi \left(A_a^{-1} S_s^{-1}(x-t) \right), \\ \psi_{a,s,t,-1}(x) &:= a^{-\frac{3}{4}} \tilde{\psi} \left(\tilde{A}_a^{-1} S_s^{-T}(x-t) \right),\end{aligned}\tag{11}$$

432 where $\tilde{\psi}(x_1, x_2) := \psi(x_2, x_1)$ for all $x = (x_1, x_2) \in \mathbb{R}^2$.
 433 Following [13], we define the continuous shearlet transform
 434 as follows:

435 **Definition 10.1** (Continuous shearlet transform). Let $\psi \in L^2(\mathbb{R}^2)$. Then the family of functions $\psi_{a,s,t,\iota}: \mathbb{R}^2 \rightarrow \mathbb{R}$
 436 parametrized by $(a, s, t, \iota) \in \mathbb{R}^+ \times \mathbb{R} \times \mathbb{R}^2 \times \{-1, 1\}$ that
 437 are defined in (11) is called a *shearlet system*. The corresponding (*continuous*) *shearlet transform* is defined by

$$\mathcal{SH}_\psi : L^2(\mathbb{R}^2) \rightarrow L^\infty(\mathbb{R}^+ \times \mathbb{R} \times \mathbb{R}^2 \times \{-1, 1\}),$$

441 where
 442 $\mathcal{SH}_\psi(f)(a, s, t, \iota) := \langle f, \psi_{a,s,t,\iota} \rangle.$

443 As we shall see next, if the generator function ψ has di-
 444 rectional vanishing moments, then the asymptotic behav-
 445 ior as $a \rightarrow 0$ of the continuous shearlet transform of an
 446 L^2 -function f characterizes its wavefront set. The precise
 447 statement can be found in [13].

448 10.2. CartoonX with Spatial Regularization

449 To avoid unnecessary blurry regions in CartoonX we pe-
 450 nalyze unnecessary energy in pixel space with the ℓ_1 norm
 451 in pixel space. More precisely, we propose to optimize a
 452 wavelet mask m with projected SGD and the optimization
 453 objective

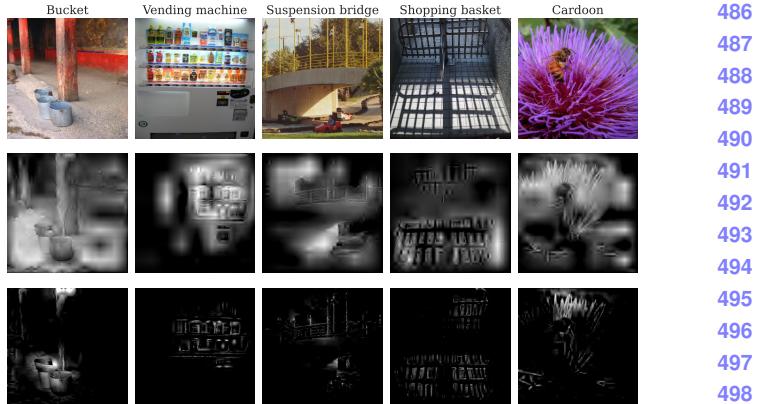
$$\begin{aligned} \min_m - \mathbb{E}_{u \sim \mathcal{V}} & \left[\Phi(\text{DWT}^{-1}(m \odot \text{DWT}(x) + (1 - m) \odot u)) \right] \\ & + \lambda_1 \cdot \|m\|_1 + \lambda_2 \cdot \|\text{DWT}^{-1}(m \odot \text{DWT}(x))\|_1. \end{aligned}$$

463 We illustrate the effect in Figure Fig. 3 by computing Car-
 464 toonX with and without spatial regularization for five ran-
 465 domly sampled images from Imagenet of different classes.
 466 The spatial regularization makes CartoonX much easier to
 467 interpret and resolves previous ambiguities.

468 10.3. CartoonX Theory

469 Fong and Vedaldi [12] first observed that explanation
 470 masks in pixel space are susceptible to explanation artifacts
 471 if the mask is not regularized to be smooth. In section
 472 Sec. 11.1, we demonstrate that many explanation artifacts
 473 of pixel explanation masks stem from artificially formed
 474 edges that make up artificial prototypical class features. It
 475 is therefore, desirable to have explanations that cannot cre-
 476 ate artificial edges. In this section, we show theoretically
 477 that CartoonX cannot create artificial edges when model-
 478 ing images and edges in a continuous space. We need to
 479 work in continuous space to obtain a rigorous definition of
 480 edges but the conclusion is verified later experimentally for
 481 discrete images.

482 In this theory section, we will model images as
 483 $L^2([0, 1]^2)$ functions $x: [0, 1]^2 \rightarrow [0, 1]$. Edges are sin-
 484 gularities, which we will model with the notion of Lipschitz



485 **Figure 3.** Fair qualitative comparison on images randomly sam-
 486 pled from Imagenet. First row: Images from Imagenet with clas-
 487 sification label. Second row: CartoonX without spatial regulariza-
 488 tion from Kolek et al. [19]. Third row: Our CartoonX with spatial
 489 regularization.

490 α regularity, which we introduce below. Our final theorem
 491 will show that masking an image $x: [0, 1]^2 \rightarrow [0, 1]$ in its
 492 wavelet domain cannot not create new (artificial) edges.

493 **Definition 10.2** (Lipschitz α regular functions). A function
 494 $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ is uniformly Lipschitz $\alpha \geq 0$ over a domain
 495 $\Omega \subset \mathbb{R}^2$ if there exists $K > 0$, such that for any $v \in \Omega$ one
 496 can find a polynomial p_v of degree $\lfloor \alpha \rfloor$ such that

$$\forall x \in \Omega, |f(x) - p_v(x)| \leq K|x - v|^\alpha. \quad (12)$$

500 The infimum of K , which satisfies the above equation, is the
 501 *homogenous* Hölder α norm $\|f\|_{\tilde{C}^\alpha}$. The Hölder α norm of
 502 $\|f\|_{C^\alpha} := \|f\|_{\tilde{C}^\alpha} + \|f\|_\infty$ also imposes that f is bounded.

503 For our results, we also need the defintion of *fast decaying*
 504 *derivatives*, which we state below.

505 **Definition 10.3** (Fast Derivative Decay). If $f: [0, 1]^2 \rightarrow \mathbb{R}$
 506 is $C^{\lfloor \alpha \rfloor}$ then we say the derivatives of f decay fast if and
 507 only if for all multi-indices $\nu = (\nu_1, \nu_2)$ with $\nu_1 + \nu_2 \leq \lfloor \alpha \rfloor$
 508 and for all $m \in \mathbb{N}$ there exists $C_m \in \mathbb{R}_+$ such that

$$\forall u \in [0, 1]^2 : |\partial^\nu \psi(t)| \leq \frac{C_m}{1 + |t|^m}. \quad (13)$$

509 Our result relies on the fact that function regularity as
 510 measured by Lipschitz α regularity is characterized by the
 511 decay of wavelet coefficients. This is shown in the follow-
 512 ing theorem, which is a small adaptation of Theorem 9.15
 513 in [27].

514 **Theorem 1.** Let $f: [0, 1]^2 \rightarrow \mathbb{R}$ be Lipschitz α on a do-
 515 main $\Omega_0 \subset [0, 1]^2$. Consider a mother-wavelet ψ with $\lfloor \alpha \rfloor$
 516 vanishing moments. Then, there exists a constant $B > 0$,

such that for all $1 \leq l \leq 3, j \in \mathbb{Z}, 0 \leq n < 2^{-j}$ with $\text{supp } \psi_{j,n}^l \subset \Omega_0$, we have

$$|\langle f, \psi_{j,n}^l \rangle| \leq B \cdot \|f\|_{\tilde{C}^\alpha} \cdot 2^{j(\alpha+1)}. \quad (14)$$

Moreover, if ψ is $C^{\lfloor \alpha \rfloor}$ and its derivatives are decaying fast (see Definition 10.3) then there exists a constant $A > 0$ such that for all $1 \leq l \leq 3, j \in \mathbb{Z}, 0 \leq n < 2^{-j}$ with $\text{supp } \psi_{j,n}^l \subset \Omega_0$, we have

$$A \cdot \|f\|_{\tilde{C}^\alpha} \cdot 2^{j(\alpha+1)} \leq |\langle f, \psi_{j,n}^l \rangle|. \quad (15)$$

In particular, inequality (14) implies f is Lipschitz α on Ω_0 .

The Proof is in the supplementary material. The above theorem essentially says that an image's Lipschitz α regularity in a domain $\Omega_0 \subset [0, 1]^2$, is completely determined by the decay of its wavelet coefficients, for wavelets $\psi_{j,n}^l$ that have spatial support in Ω_0 . Next we define, what we call the *admissible singularity set*, which contains 1d edges as a special case.

Definition 10.4 (Admissible Singularity Set). Let $0 \leq \alpha < 1 \leq \beta$. We say a set $\Gamma \subset [0, 1]^2$ is an admissible singularity set of order $(\alpha, \beta) \in \mathbb{R}_+^2$ for an image $x : [0, 1]^2 \rightarrow \mathbb{R}$ if the following three properties are satisfied:

1. The image x is Lipschitz α on Γ
2. The image x is Lipschitz β on $[0, 1]^2 \setminus \Gamma$
3. There exists two constants $0 < K_1 \leq K_2$ such that for all $j \in \mathbb{N}$

$$C_1 2^{-j} \leq \#\{I_{n,m}^{(j)} \mid \Gamma \cap I_{n,m}^{(j)} \neq \emptyset\} \leq C_2 2^{-j}, \quad (16)$$

where $I_{n,m}^{(j)} := [2^{-j} \cdot n, 2^{-j} \cdot (n+1)] \times [2^{-j} \cdot m, 2^{-j} \cdot (m+1)]$.

The most important example of an admissible singularity set $\Gamma \subset [0, 1]^2$ of order (α, β) for an image $x : [0, 1]^2 \rightarrow \mathbb{R}$, is an image that is differentiable outside of Γ (case $\beta = 1$) and has an edge in a 1d curve Γ (case $\alpha = 0$). The property in (16) is satisfied by 1d curves and sets of finitely many points but not by 2d regions. Next, we define the wavelet mask in our continuous image setting.

Definition 10.5 (Wavelet Mask in Continuum). Consider an image $x : [0, 1]^2 \rightarrow \mathbb{R}$ and let

$$m_\psi : \{(j, n) \in \mathbb{Z} \times \mathbb{Z}^2 \mid 2^j n \in [0, 1)\} \rightarrow \{0, 1\} \quad (17)$$

be a mask on the wavelet coefficients $\{\langle x, \psi_{j,n}^l \rangle\}_{1 \leq l \leq 3, j \leq J, 2^j n \in [0, 1)}$ and let

$$m_\phi : \{n \in \mathbb{Z}^2 \mid 2^j n \in [0, 1)\} \rightarrow \{0, 1\} \quad (18)$$

be a mask on the approximation coefficients $\{\langle x, \phi_{J,n}^l \rangle\}_{2^J n \in [0, 1)}$. We define the mean and standard deviation of wavelet coefficients at scale 2^j in x as $\mu_\psi[j]$ and $\sigma_\psi[j]$. They are computed as:

$$\mu_\psi[j] := \frac{1}{3 \cdot 2^{-2j}} \sum_{1 \leq l \leq 3} \sum_{n: 2^j n \in [0, 1)^2} \langle x, \psi_{j,n}^l \rangle \quad (19)$$

$$\sigma_\psi[j] := \left(\frac{1}{3 \cdot 2^{-2j}} \sum_{1 \leq l \leq 3} \sum_{n: 2^j n \in [0, 1)^2} |\langle x, \psi_{j,n}^l \rangle - \mu_\psi[j]|^2 \right)^{1/2} \quad (20)$$

Moreover, we define the wavelet obfuscation of x with the wavelet mask m_ψ, m_ϕ and perturbations $\{z_\psi[j, n], z_\phi[J, n]\}_{j \in \mathbb{Z}, n \in \mathbb{Z}^2}$ as

$$\hat{x} = \sum_{1 \leq l \leq 3} \sum_{j \leq J} \sum_{2^j n \in [0, 1)^2} \left(m_\psi[j, n] \langle x, \psi_{j,n}^l \rangle \right. \quad (21)$$

$$\left. + (1 - m_\psi[j, n]) z_\psi[j, n] \right) \psi_{j,n}^l$$

$$+ \sum_{2^J n \in [0, 1)} \left(m_\phi[J, n] \langle x, \phi_{J,n}^l \rangle \right. \quad (21)$$

$$\left. + (1 - m_\phi[J, n]) z_\phi[J, n] \right) \phi_{J,n}^l.$$

Note that the above definition of the wavelet obfuscation \hat{x} only differs from the obfuscation in the CartoonX method (see term in the expectation in Eq. (10)) by a) the fact that images are discrete in practice and by b) masks being only approximately binary. We finally give our final theorem, that shows the obfuscation of \hat{x} with appropriately chosen perturbations cannot create new edges.

Theorem 2 (CartoonX does not create artificial edges.). Consider an image $x : [0, 1]^2 \rightarrow \mathbb{R}$ in $L^2([0, 1]^2)$ with an admissible singularity set $\Gamma \subset [0, 1]^2$ of order (α, β) . Let ψ be a mother-wavelet with $\lfloor \alpha \rfloor$ vanishing moments, derivatives decaying fast and let ϕ be a corresponding scaling function for a multiresolution approximation. Then the following holds :

1. The mean wavelet coefficient at a fixed scale 2^j is zero:

$$\mu_\psi[j] = 0 \quad (22)$$

and the wavelet coefficient standard deviation at scale 2^j decays as

$$\sigma_\psi[j] = \mathcal{O}(2^{j(\min\{\alpha+3/2, \beta+1\})}). \quad (23)$$

2. Let $\mathcal{Q}_\psi \in \{\text{Unif}(\mu_\psi[j] - \sigma_\psi[j] 2^{j/2}, \mu_\psi[j] + \sigma_\psi[j] 2^{j/2}), \text{Const}(0)\}$ and let \mathcal{Q}_ϕ be an arbitrary distribution over \mathbb{R} . The obfuscation \hat{x} of x with any wavelet mask (m_ψ, m_ϕ) and perturbations $z_\psi[j, n] \sim \mathcal{Q}_\psi$ and $z_\phi[J, n] \sim \mathcal{Q}_\phi$ satisfies with probability one

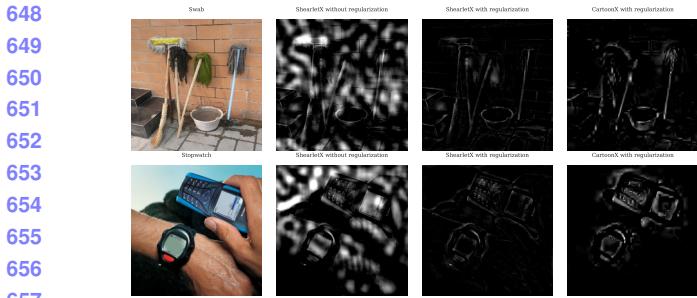


Figure 4. ShearletX

over randomness in $z_\phi[J, n], z_\psi[j, n]$ that \hat{x} is Lipschitz $\alpha + 1$ on $[0, 1]^2 \setminus \Gamma$, Lipschitz β on Γ , and \hat{x} has no admissible singularity set in $[0, 1]^2 \setminus \Gamma$. In particular, \hat{x} cannot contain edges (modeled as Lipschitz $\alpha < 1$ regions) that were not already present in x .

10.4. ShearletX

Wavelets sparsely represent piece-wise smooth (cartoon-like) images but not optimally. In fact, wavelets can only represent point-like singularities optimally. Piece-wise smooth images have edges as singularities and are optimally represented by a different representation system called *shearlets* [22]. Continuous shearlets constitute a 2-parameter dilation group, parametrized by parabolic scaling matrices and shear matrices. The continuous shearlets depend on three parameters, the scaling parameter $a > 0$, the shear parameter $s \in \mathbb{R}$ and the translation parameter $t \in \mathbb{R}^2$, and are defined as

$$\psi_{a,s,t}(u) = a^{-3/4}\psi(D_{a,s}^{-1}(u - t)), \quad (24)$$

where ψ is an appropriately chosen mother wavelet, $u \in \mathbb{R}^2$ and

$$D_{a,s} = \begin{bmatrix} a & -a^{1/2}s \\ 0 & a^{1/2} \end{bmatrix} \quad (25)$$

is a shearing and scaling matrix. The associated continuous shearlet transform is defined by

$$\text{SH}\{f\}(a, s, t) := \langle f, \psi_{a,s,t} \rangle. \quad (26)$$

There is of course also a digital shearlet transform, which we use to optimize and explanation mask in the shearlet domain of images. The explanation objective becomes:

$$\begin{aligned} \min_m & - \mathbb{E}_{u \sim \mathcal{V}} \left[\Phi(\text{SH}^{-1}(m \odot \text{SH}(x) + (1 - m) \odot u)) \right] \\ & + \lambda_1 \cdot \|m\|_1 + \lambda_2 \cdot \|\text{SH}^{-1}(m \odot \text{SH}(x))\|_1. \end{aligned}$$

We call the resulting explanation method *ShearletX*. In Figure Fig. 4, we compare ShearletX with and without spatial regularization to WaveletX (with spatial regularization).

We see that the spatial regularization term that we used to improve WaveletX, is again crucial for ShearletX. ShearletX produces sharper visualizations that are less blurry than WaveletX, which is directly the result of shearlets being better suited to represent piece-wise smooth images sparsely. Although ShearletX produces cleaner visualizations, it is much more computationally expensive due to the shearlet representation having more coefficients. In our implementation, an image of 226×226 pixels has $49 \times 226 \times 226$ shearlet parameters (49 shearing and scaling parameters and 226×226 locations) unlike the wavelet representation which has 226×226 wavelet coefficients. We leave a more efficient implementation for a faster ShearletX to future work.

10.5. Shearlets and wavefront set

The wavefront set is a concept that characterizes the oriented singularities of distributions, in particular, of $L^2(\mathbb{R}^2)$ functions.

Definition 10.6. [17, Section 8.1] Let $f \in L^2(\mathbb{R}^2)$ and $k \in \mathbb{N}$. A point $(x, \lambda) \in \mathbb{R}^2 \times \mathbb{S}^1$ is a *k-regular directed point* of f if there exist open neighbourhoods U_x and V_λ of x and λ , respectively and a smooth function $\phi \in C^\infty(\mathbb{R}^2)$ with $\text{supp } \phi \subset U_x$ and $\phi(x) = 1$ such that

$$|\widehat{\phi f}(\xi)| \leq C_k (1+|\xi|)^{-k} \quad \text{for all } \xi \in \mathbb{R}^2 \setminus \{0\} \text{ s.t. } \xi/|\xi| \in V_\lambda$$

holds for some $C_k > 0$. The *k-wavefront set* $\text{WF}_k(f)$ is the complement of the set of all *k*-regular directed points and the *wavefront set* $\text{WF}(f)$ is defined as

$$\text{WF}(f) := \bigcup_{k \in \mathbb{N}} \text{WF}_k(f),$$

One can make use of the shearlet transform to extract the wavefront set of a given function as follows.

Theorem 3. Let $f \in L^2(\mathbb{R}^2)$ and assume $(x_0, \lambda_0) \in \mathbb{R}^2 \times \mathbb{S}^1$ is a *k*-regular directed point of f for some $k \in \mathbb{N}$. Next, consider a continuous shearlet system with generator function $\psi \in H^l(\mathbb{R}^2)$, $l \in \mathbb{N}$, with Fourier transform $\widehat{\psi} \in L^1(\mathbb{R}^2)$ where ψ has $m \in \mathbb{N}$ vanishing moments in x_1 -direction, i.e.,

$$\int_{\mathbb{R}^2} \frac{|\widehat{\psi}(\xi_1, \xi_2)|^2}{|\xi_1|^{2m}} d\xi < \infty.$$

Finally, assume ψ displays the following asymptotic behavior for $p \in \mathbb{N}$:

$$|\psi(x)| = \mathcal{O}((1+|x|)^{-p}) \quad \text{for } |x| \rightarrow \infty.$$

Then, there exist a neighborhood $U_0 \subset \mathbb{R}^2$ of x_0 and a neighborhood $S_0 \subset \mathbb{S}^1$ of λ_0 such that

$$|\mathcal{SH}_\psi(f)(a, s, x, \lambda)| = \mathcal{O}\left(a^{\frac{p}{2}-\frac{3}{4}} + a^{\frac{m}{4}} + a^{\frac{3k}{4}-\frac{3}{4}} + a^{\frac{3l}{4}}\right) \quad \text{as } a \rightarrow 0$$

Method	Time
Integrated Gradients (Sundararajan et al. [38])	0.22s
Smooth Mask (Fong et al. [11])	20s
Wavelet Mask by Kolek et al. [19]	20s
WaveletX (ours)	20s
ShearletX (ours)	20s

Table 1. Computation time for explanation of Resnet18 [15] decision with 256×256 input image. It is well-known that mask explanations are much more computational expensive than more heuristic methods such as Integrated Gradients [38]. Our WaveletX implementation is faster than the implementation of Kolek et al. [19] due to better choices of learning rate. ShearletX is slower than WaveletX due to the mask on the shearlet representation being much larger.

for all $x \in U_0$ and all $s \in \mathbb{R}$ and $\iota \in \{-1, 1\}$ such that $\lambda(s, \iota) \in S_0$, where

$$\lambda(s, \iota) := \begin{cases} \left(\frac{1}{\sqrt{s^2 + 1}}, \frac{s}{\sqrt{s^2 + 1}} \right) & \text{if } \iota = 1, \\ \left(\frac{s}{\sqrt{s^2 + 1}}, \frac{1}{\sqrt{s^2 + 1}} \right) & \text{if } \iota = -1. \end{cases} \quad (27)$$

Remark. Under suitable assumptions on the shearlet generator ψ , the converse of Theorem 3 holds as well. More precisely, following [13], assume that ψ is sufficiently regular for any $k \in \mathbb{N}$. Next, let $(x_0, \lambda_0) \in \mathbb{R}^2 \times \mathbb{S}^1$ and assume there exist a neighborhood $U_0 \subset \mathbb{R}^2$ of x_0 and a neighborhood $S_0 \subset \mathbb{S}^1$ of λ_0 such that

$$|\mathcal{SH}_\psi(f)(a, s, x, \iota)| = \mathcal{O}(a^n) \quad \text{as } a \rightarrow 0,$$

holds for sufficiently large $n \in \mathbb{N}$ uniformly for $x \in U_0$ and all $s \in \mathbb{R}$, $\iota \in \{-1, 1\}$ such that $\lambda(s, \iota) \in S_0$. Then, $(x_0, \lambda_0) \notin \text{WF}_k(f)$.

Theorem 3 and Remark 10.5 demonstrate that the wavefront set is completely determined by the decay properties of the shearlet transform. This implies that in the continuous setting, one can compute the wavefront set of a function by first computing its continuous shearlet transform, then analyzing the pairs of point and direction where this shearlet transform exhibits rapid decay as $a \rightarrow 0$.

11. Explanation Bias of Mask Explanations

Mask based explanations can be corrupted by unintended bias. Fong and Vedaldi [12] first showed that pixel space masks without total variation regularization can be corrupt. In light of, our understanding of adversarial vulnerabilities of neural networks, one may think that the explanation corruption comes from uninterpretable adversarial artifacts. However, we empirically find that mostly the corruptions come from artificial edges that make up class prototypes.

11.1. Artificial Edges & Hallucinations

For most instances pixel masks with no or low total variation regularization will highlight a sensible region of the image that activates the class label. When the explanation is corrupted we find it is due to artificial edges forming semantically meaningful concepts that activate the predicted label but do not explain what the classifier saw in the image (see Figure for examples). This can be mitigated with a high TV regularization on the pixel mask. As we show in Section Sec. 11.2, high TV regularization has its own drawback. WaveletX presents itself as a method that is not vulnerable to artificial edges that corrupt the explanation. In Section Sec. 10.3, we proved that CartoonX in continuum cannot produce artificial edges. Our experiments confirm that this translates to the digital case. We do not observe examples as in Figure for WaveletX and we ran experiments that estimate and quantify the amount of artificial edges for pixel masks, WaveletX, each with and without regularization. For that, we compute each explanation on 1000 randomly drawn images from the Imagenet validation set. We quantify the amount of artificial edges in an explanation E for image x with the *hallucination score*:

$$\text{Hallucination Score} := \frac{\#\text{[Edges}(E)\setminus\text{Edges}(x)]}{\#\text{Edges}(x)} \quad (28)$$

The hallucination score counts the number of edges that are present in the explanation E but not in the original image x normalized by the total number of edges in x . To extract the edges we use a shearlet-based edge extractor [5, 30, 40]. To compensate for estimation errors in the edge extractor, we slightly thicken the edges in $\text{Edges}(x)$ for the term $\#\text{[Edges}(E)\setminus\text{Edges}(x)]$. Figure illustrates the edge extractor and artificial edges computed in pixel masks, WaveletX, and ShearletX. In Figure ??, we present a scatterplot comparing the expected distortion and hallucination score of the 1000 explanations for the distinct mask method. The result shows that pixel masks without TV regularization produce more artificial edges than WaveletX (with and without spatial regularization). In particular, pixel masks without TV regularization have a very fat tail for very high artifact explanations. This reflects what we observed empirically: the majority of explanations are to corrupted but there is a significant amount of outliers that are corrupted by a large amount of artificial edges.

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

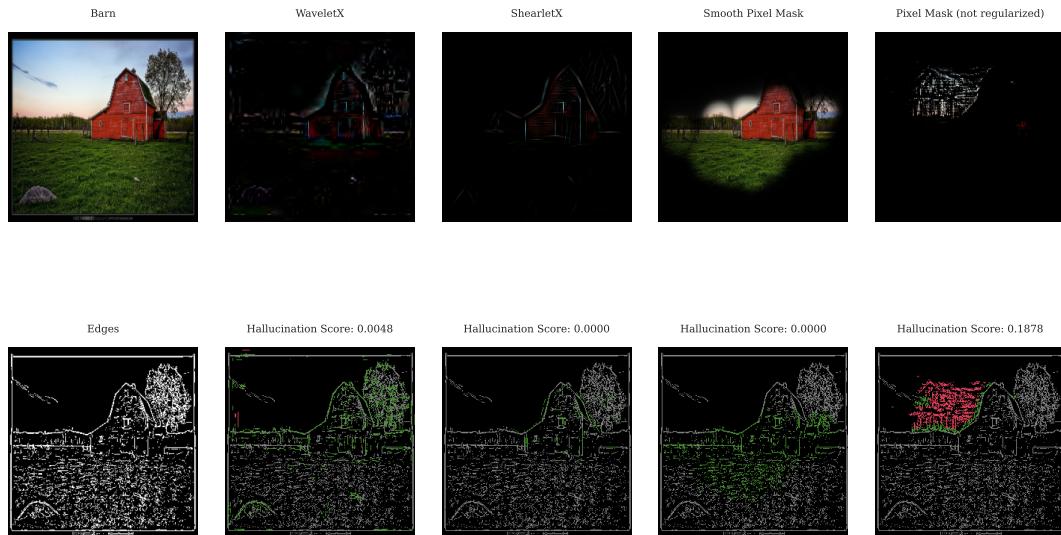


Figure 5. Visualization of Artificial Edges in Explanations

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

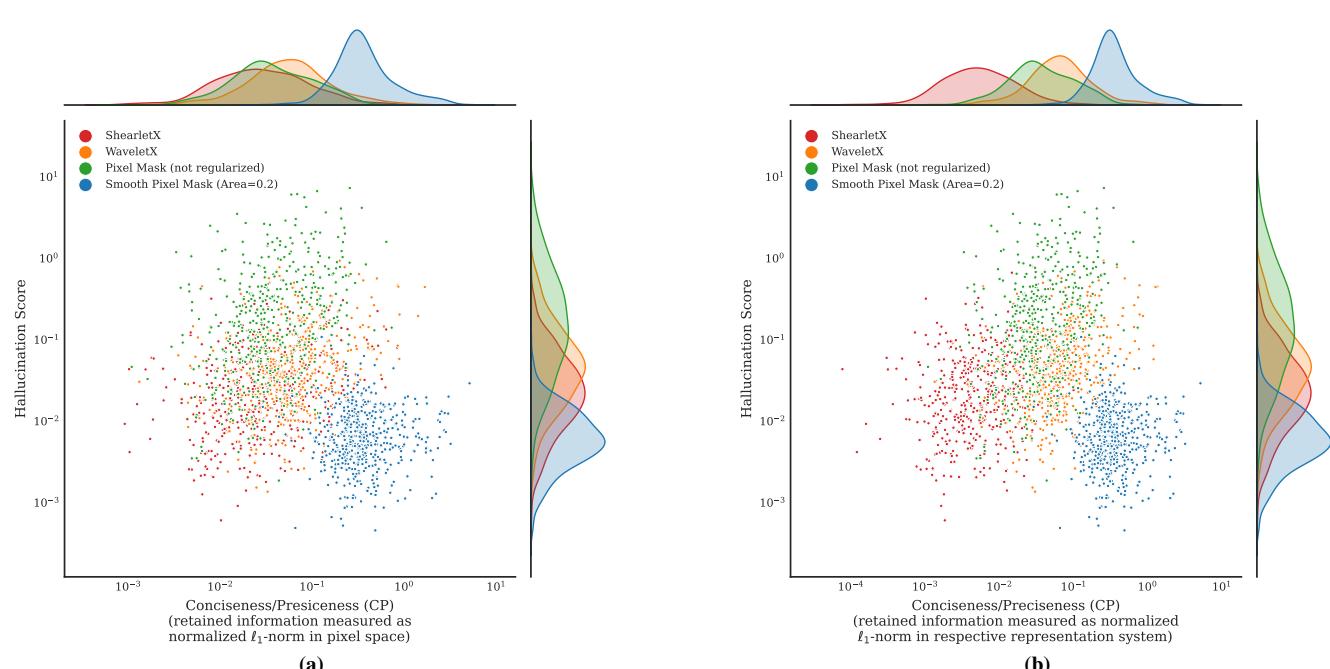
913

914

915

916

917

Figure 6. (a) Retained information for conciseness measured as normalized ℓ_1 -norm in pixel domain. (b) Retained information for conciseness measured as normalized ℓ_1 -norm in respective representation system. Evaluated on Resnet18 [15] on 100 random validation Imagenet samples.

11.2. Discriminating Spatially Close Features

11.3. The Inpainting GAN Bias

12. Sanity Checks

Layer Randomization

Class Localization

Shearletx is not an edge detector

13. Limitations

14. Conclusion

14.1. Language

972	Method	Retained Probability	Retained Info (Pixel)	Retained Info (Representation)	1026
973	Smooth Pixel Mask [11] (area = 0.2)	73.76%	%27.55	27.55%	1027
974	Smooth Pixel Mask [11] (area = 0.1)	51.63%	14.27%	14.72%	1028
975	Smooth Pixel Mask [11] (area = 0.05)	44.90%	7.41%	7.41%	1029
976	Pixel Mask (not regularized)	117.66%	3.97%	3.97%	1030
977	WaveletX (our second best)	66.63%	6.63%	7.21%	1031
978	ShearletX (our best)	115.81%	2.77%	0.74%	1032
979					1033

Table 2. Results. Ours is better.

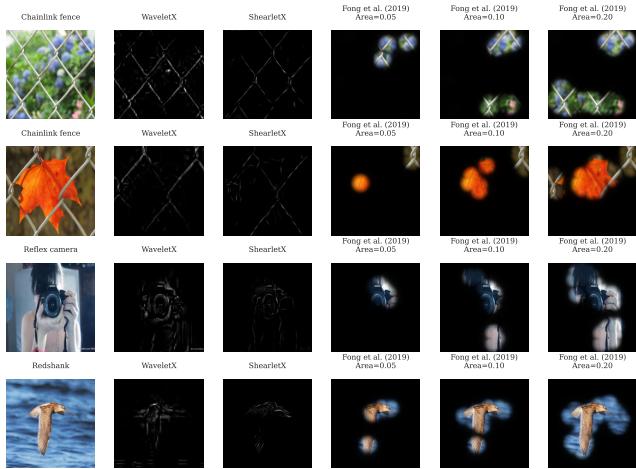


Figure 7. We compare our WaveletX and ShearletX to the smooth masks by Fong et al. (2019) [11] for various area sizes of the mask. Smooth masks fail to unambiguously highlight the chain as relevant.

14.2. Dual submission

Please refer to the author guidelines on the CVPR 2023 web page for a discussion of the policy on dual submissions.

14.3. Paper length

Papers, excluding the references section, must be no longer than eight pages in length. The references section will not be included in the page count, and there is no limit on the length of the references section. For example, a paper of eight pages with two pages of references would have a total length of 10 pages. **There will be no extra page charges for CVPR 2023.**

Overlength papers will simply not be reviewed. This includes papers where the margins and formatting are deemed to have been significantly altered from those laid down by this style guide. Note that this L^AT_EX guide already sets figure captions and references in a smaller font. The reason such papers will not be reviewed is that there is no provision for supervised revisions of manuscripts. The reviewing process cannot determine the suitability of the paper for presentation in eight pages if it is reviewed in eleven.

14.4. The ruler

The L^AT_EX style defines a printed ruler which should be present in the version submitted for review. The ruler is provided in order that reviewers may comment on particular lines in the paper without circumlocution. If you are preparing a document using a non-L^AT_EX document preparation system, please arrange for an equivalent ruler to appear on the final output pages. The presence or absence of the ruler should not change the appearance of any other content on the page. The camera-ready copy should not contain a ruler. (L^AT_EX users may use options of cvpr.sty to switch between different versions.)

Reviewers: note that the ruler measurements do not align well with lines in the paper — this turns out to be very difficult to do well when the paper contains many figures and equations, and, when done, looks ugly. Just use fractional references (*e.g.*, this line is 087.5), although in most cases one would expect that the approximate location will be adequate.

14.5. Paper ID

Make sure that the Paper ID from the submission system is visible in the version submitted for review (replacing the “*****” you see in this document). If you are using the L^AT_EX template, **make sure to update paper ID in the appropriate place in the tex file.**

14.6. Mathematics

Please number all of your sections and displayed equations as in these examples:

$$E = m \cdot c^2 \quad (29)$$

and

$$v = a \cdot t. \quad (30)$$

It is important for readers to be able to refer to any particular equation. Just because you did not refer to it in the text does not mean some future reader might not need to refer to it. It is cumbersome to have to use circumlocutions like “the equation second from the top of page 3 column 1”. (Note that the ruler will not be present in the final copy, so is not

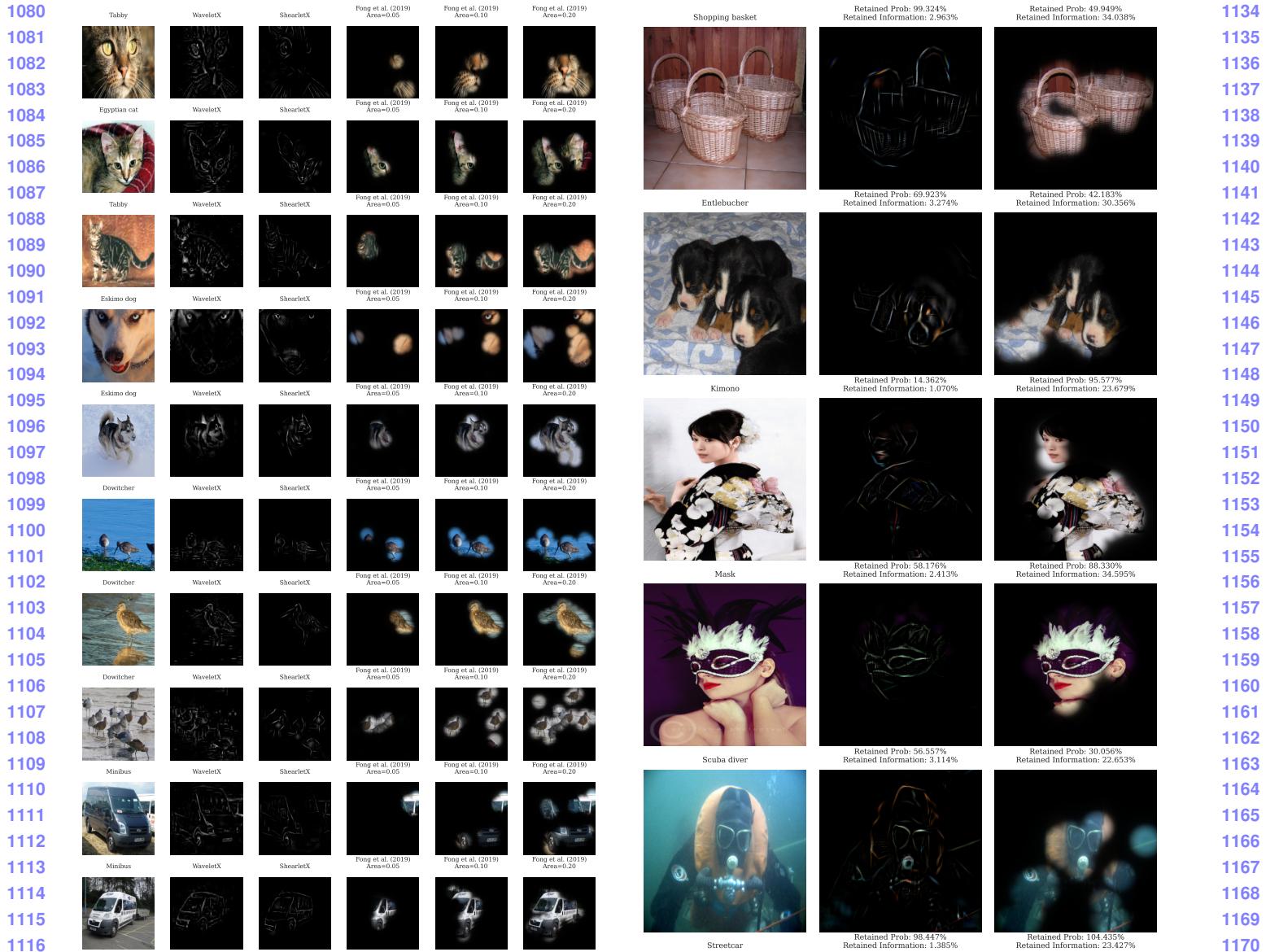


Figure 8. We compare our WaveletX and ShearletX to the smooth masks by Fong et al. (2019) [11] for various area sizes of the mask. Smooth masks fail to unambiguously highlight the chain as relevant.

an alternative to equation numbers). All authors will benefit from reading Mermin’s description of how to write mathematics: <http://www.pamitc.org/documents/mermin.pdf>.

14.7. Blind review

Many authors misunderstand the concept of anonymizing for blind review. Blind review does not mean that one must remove citations to one’s own work—in fact it is often impossible to review a paper unless the previous citations are known and available.

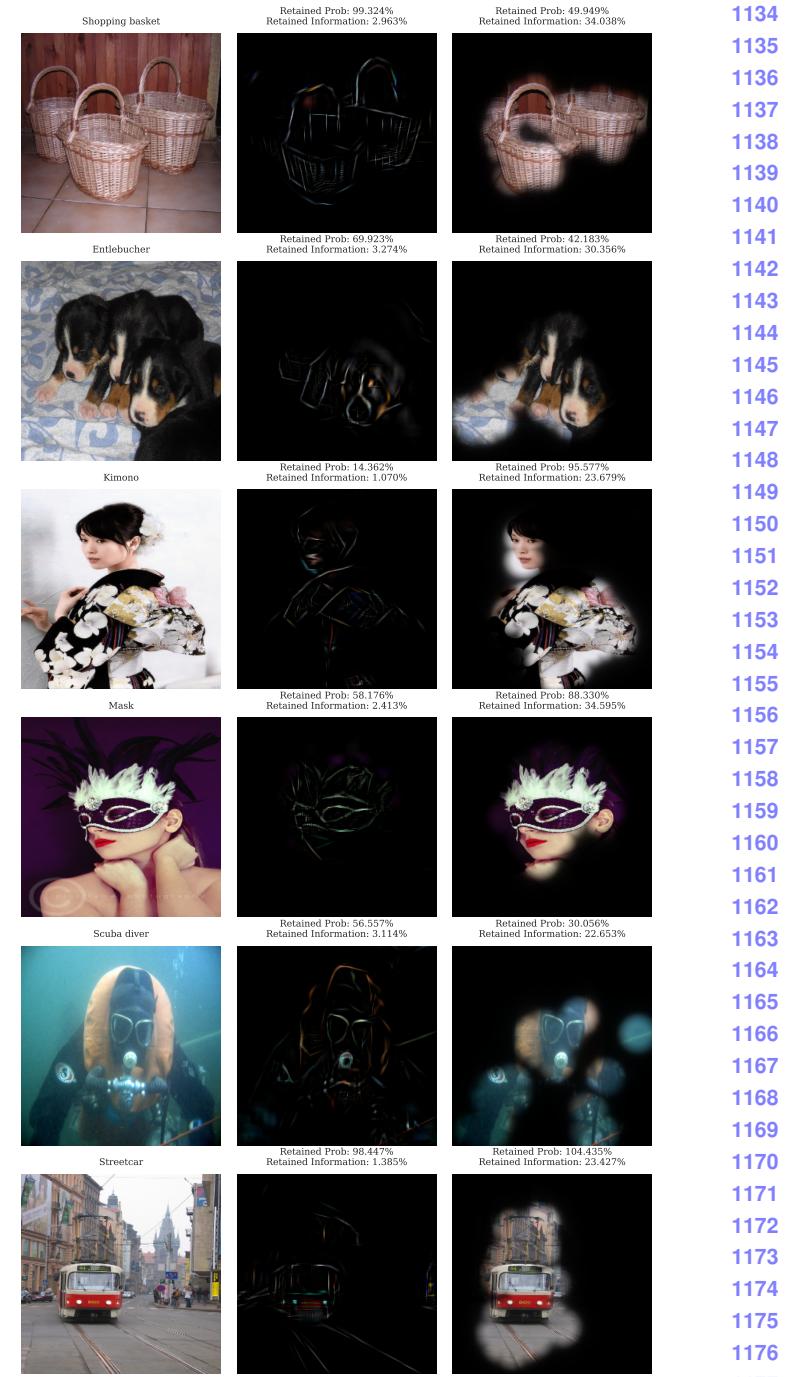
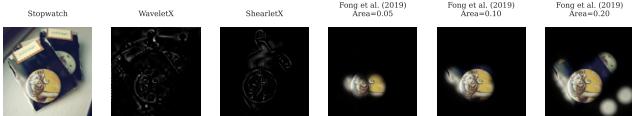


Figure 9. Caption

Blind review means that you do not use the words “my” or “our” when citing previous work. That is all. (But see below for tech reports.)

Saying “this builds on the work of Lucy Smith [1]” does not say that you are Lucy Smith; it says that you are building on her work. If you are Smith and Jones, do not say “as we show in [7]”, say “as Smith and Jones show in [7]” and at

1188



1189

1190

1191

1192

1193

1194

1195

1196

1197

the end of the paper, include reference 7 as you would any other cited work.

An example of a bad paper just asking to be rejected:

An analysis of the frobnicatable foo filter.

In this paper we present a performance analysis of our previous paper [1], and show it to be inferior to all previously known methods. Why the previous paper was accepted without this analysis is beyond me.

[1] Removed for blind review

An example of an acceptable paper:

An analysis of the frobnicatable foo filter.

In this paper we present a performance analysis of the paper of Smith *et al.* [1], and show it to be inferior to all previously known methods. Why the previous paper was accepted without this analysis is beyond me.

[1] Smith, L and Jones, C. “The frobnicatable foo filter, a fundamental contribution to human knowledge”. Nature 381(12), 1-213.

If you are making a submission to another conference at the same time, which covers similar or overlapping material, you may need to refer to that submission in order to explain the differences, just as you would if you had previously published related work. In such cases, include the anonymized parallel submission [23] as supplemental material and cite it as

[1] Authors. “The frobnicatable foo filter”, F&G 2014 Submission ID 324, Supplied as supplemental material fg324.pdf.

Finally, you may feel you need to tell the reader that more details can be found elsewhere, and refer them to a technical report. For conference submissions, the paper must stand on its own, and not *require* the reviewer to go to a tech report for further details. Thus, you may say in the body of the paper “further details may be found in [24]”. Then submit the tech report as supplemental material. Again, you may not assume the reviewers will read this material.

Sometimes your paper is about a problem which you tested using a tool that is widely known to be restricted to a single institution. For example, let’s say it’s 1969, you have solved a key problem on the Apollo lander, and you believe that the CVPR70 audience would like to hear about your solution. The work is a development of your celebrated 1968 paper entitled “Zero-g frobnication: How being the only people in the world with access to the Apollo lander source code makes us a wow at parties”, by Zeus *et al.*

You can handle this paper like any other. Do not write “We show how to improve our previous work [Anonymous, 1968]. This time we tested the algorithm on a lunar lander [name of lander removed for blind review]”. That would be silly, and would immediately identify the authors. Instead write the following:

We describe a system for zero-g frobnication. This system is new because it handles the following cases: A, B. Previous systems [Zeus *et al.* 1968] did not handle case B properly. Ours handles it by including a foo term in the bar integral.

...

The proposed system was integrated with the Apollo lunar lander, and went all the way to the moon, don’t you know. It displayed the following behaviours, which show how well we solved cases A and B: ...

As you can see, the above text follows standard scientific convention, reads better than the first version, and does not explicitly name you as the authors. A reviewer might think it likely that the new paper was written by Zeus *et al.*, but cannot make any decision based on that guess. He or she would have to be sure that no other authors could have been contracted to solve problem B.

FAQ

Q: Are acknowledgements OK?

A: No. Leave them for the final copy.

Q: How do I cite my results reported in open challenges?

A: To conform with the double-blind review policy, you can report results of other challenge participants together with your results in your paper. For your results, however, you should not identify yourself and should not mention your participation in the challenge. Instead present your results referring to the method proposed in your paper and draw conclusions based on the experimental comparison to other results.

14.8. Miscellaneous

Compare the following:

$\$conf_a\$$	$conf_a$
$\$\mathit{conf}_a\$$	$conf_a$

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

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



Figure 11. Example of caption. It is set in Roman so that mathematics (always set in Roman: $B \sin A = A \sin B$) may be included without an ugly clash.

See The TeXbook, p165.

The space after *e.g.*, meaning “for example”, should not be a sentence-ending space. So *e.g.* is correct, *e.g.* is not. The provided \eg macro takes care of this.

When citing a multi-author paper, you may save space by using “et alia”, shortened to “*et al.*” (not “*et. al.*” as “*et*” is a complete word). If you use the \etal macro provided, then you need not worry about double periods when used at the end of a sentence as in Alpher *et al.* However, use it only when there are three or more authors. Thus, the following is correct: “Frobnication has been trendy lately. It was introduced by Alpher [1], and subsequently developed by Alpher and Fotheringham-Smythe [2], and Alpher *et al.* [3].”

This is incorrect: “... subsequently developed by Alpher *et al.* [2] ...” because reference [2] has just two authors.

15. Formatting your paper

All text must be in a two-column format. The total allowable size of the text area is $6\frac{7}{8}$ inches (17.46 cm) wide by $8\frac{7}{8}$ inches (22.54 cm) high. Columns are to be $3\frac{1}{4}$ inches (8.25 cm) wide, with a $\frac{5}{16}$ inch (0.8 cm) space between them. The main title (on the first page) should begin 1 inch (2.54 cm) from the top edge of the page. The second and following pages should begin 1 inch (2.54 cm) from the top edge. On all pages, the bottom margin should be $1\frac{1}{8}$ inches (2.86 cm) from the bottom edge of the page for 8.5×11 -inch paper; for A4 paper, approximately $1\frac{5}{8}$ inches (4.13 cm) from the bottom edge of the page.

15.1. Margins and page numbering

All printed material, including text, illustrations, and charts, must be kept within a print area $6\frac{7}{8}$ inches (17.46 cm) wide by $8\frac{7}{8}$ inches (22.54 cm) high. Page numbers should be in the footer, centered and $\frac{3}{4}$ inches from the bottom of the page. The review version should have page num-

bers, yet the final version submitted as camera ready should not show any page numbers. The L^AT_EX template takes care of this when used properly.

15.2. Type style and fonts

Wherever Times is specified, Times Roman may also be used. If neither is available on your word processor, please use the font closest in appearance to Times to which you have access.

MAIN TITLE. Center the title $1\frac{3}{8}$ inches (3.49 cm) from the top edge of the first page. The title should be in Times 14-point, boldface type. Capitalize the first letter of nouns, pronouns, verbs, adjectives, and adverbs; do not capitalize articles, coordinate conjunctions, or prepositions (unless the title begins with such a word). Leave two blank lines after the title.

AUTHOR NAME(s) and **AFFILIATION(s)** are to be centered beneath the title and printed in Times 12-point, non-boldface type. This information is to be followed by two blank lines.

The **ABSTRACT** and **MAIN TEXT** are to be in a two-column format.

MAIN TEXT. Type main text in 10-point Times, single-spaced. Do NOT use double-spacing. All paragraphs should be indented 1 pica (approx. $\frac{1}{6}$ inch or 0.422 cm). Make sure your text is fully justified—that is, flush left and flush right. Please do not place any additional blank lines between paragraphs.

Figure and table captions should be 9-point Roman type as in Figs. 11 and 12. Short captions should be centred.

Callouts should be 9-point Helvetica, non-boldface type. Initially capitalize only the first word of section titles and first-, second-, and third-order headings.

FIRST-ORDER HEADINGS. (For example, **1. Introduction**) should be Times 12-point boldface, initially capitalized, flush left, with one blank line before, and one blank line after.

SECOND-ORDER HEADINGS. (For example, **1.1. Database elements**) should be Times 11-point boldface, initially capitalized, flush left, with one blank line before, and one after. If you require a third-order heading (we discourage it), use 10-point Times, boldface, initially capitalized, flush left, preceded by one blank line, followed by a period and your text on the same line.

15.3. Footnotes

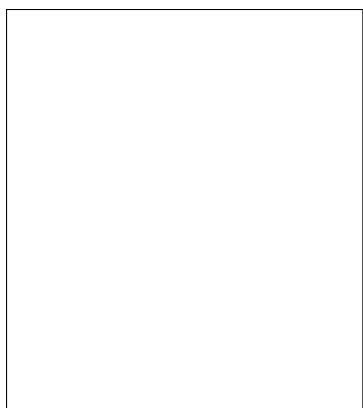
Please use footnotes¹ sparingly. Indeed, try to avoid footnotes altogether and include necessary peripheral observations in the text (within parentheses, if you prefer, as in this sentence). If you wish to use a footnote, place it at the

¹This is what a footnote looks like. It often distracts the reader from the main flow of the argument.

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

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417

(a) An example of a subfigure.



(b) Another example of a subfigure.

1418
1419**Figure 12.** Example of a short caption, which should be centered.

1420

1421 bottom of the column on the page on which it is referenced.
1422 Use Times 8-point type, single-spaced.
1423

15.4. Cross-references

1424 For the benefit of author(s) and readers, please use the

1425 `\cref{...}`1426 command for cross-referencing to figures, tables, equations,
1427 or sections. This will automatically insert the appropriate label alongside the cross-reference as in this example:1428 To see how our method outperforms previous
1429 work, please see Fig. 11 and Tab. 3. It is also
1430 possible to refer to multiple targets as once, *e.g.* to
1431 Figs. 11 and 12a. You may also return to Sec. 15
1432 or look at Eq. (30).1433 If you do not wish to abbreviate the label, for example at the
1434 beginning of the sentence, you can use the1435 `\Cref{...}`

1436 command. Here is an example:

1437 Figure 11 is also quite important.

1438 15.5. References

1439 List and number all bibliographical references in 9-point
1440 Times, single-spaced, at the end of your paper. When refe-
1441 rented in the text, enclose the citation number in square
1442 brackets, for example [23]. Where appropriate, include
1443 page numbers and the name(s) of editors of referenced
1444 books. When you cite multiple papers at once, please
1445 make sure that you cite them in numerical order like this
1446 [1, 2, 4, 23, 24]. If you use the template as advised, this will
1447 be taken care of automatically.

Method	Frobnability
Theirs	Frumpy
Yours	Frobby
Ours	Makes one's heart Frob

Table 3. Results. Ours is better.

1448 15.6. Illustrations, graphs, and photographs

1449 All graphics should be centered. In L^AT_EX, avoid using
1450 the `center` environment for this purpose, as this adds po-
1451 tentially unwanted whitespace. Instead use1452 `\centering`1453 at the beginning of your figure. Please ensure that any
1454 point you wish to make is resolvable in a printed copy of the
1455 paper. Resize fonts in figures to match the font in the body
1456 text, and choose line widths that render effectively in print.
1457 Readers (and reviewers), even of an electronic copy, may
1458 choose to print your paper in order to read it. You cannot
1459 insist that they do otherwise, and therefore must not assume
1460 that they can zoom in to see tiny details on a graphic.1461 When placing figures in L^AT_EX, it's almost always best to
1462 use `\includegraphics`, and to specify the figure width
1463 as a multiple of the line width as in the example below1464

```
\usepackage{graphicx} ...
1465 \includegraphics[width=0.8\linewidth]
1466   {myfile.pdf}
```

1467 15.7. Color

1468 Please refer to the author guidelines on the CVPR 2023
1469 web page for a discussion of the use of color in your doc-
1470 ument.1471 If you use color in your plots, please keep in mind that a
1472 significant subset of reviewers and readers may have a color

1512 vision deficiency; red-green blindness is the most frequent
 1513 kind. Hence avoid relying only on color as the discriminative
 1514 feature in plots (such as red *vs.* green lines), but add a
 1515 second discriminative feature to ease disambiguation.
 1516

1517 1518 16. Final copy

1519 You must include your signed IEEE copyright release
 1520 form when you submit your finished paper. We MUST have
 1521 this form before your paper can be published in the proceedings.
 1522

1523 Please direct any questions to the production editor in
 1524 charge of these proceedings at the IEEE Computer Society
 1525 Press: <https://www.computer.org/about/contact>.
 1526

1527 1528 References

- 1530 [1] FirstName Alpher. Frobnication. *IEEE TPAMI*, 12(1):234–
 1531 778, 2002. 13, 14
- 1532 [2] FirstName Alpher and FirstName Fotheringham-Smythe.
 1533 Frobnication revisited. *Journal of Foo*, 13(1):234–778, 2003.
 1534 13, 14
- 1535 [3] FirstName Alpher, FirstName Fotheringham-Smythe, and
 1536 FirstName Gamow. Can a machine frobnicate? *Journal
 1537 of Foo*, 14(1):234–778, 2004. 13
- 1538 [4] FirstName Alpher and FirstName Gamow. Can a computer
 1539 frobnicate? In *CVPR*, pages 234–778, 2005. 14
- 1540 [5] Andrade-Loarca, Kutyniok, and Öktem. Shearlets as feature
 1541 extractor for semantic edge detection: the model-based and
 1542 data-driven realm. *R. Soc. A*, Volume 476, Issue 2243, 2020.
 1543 8
- 1544 [6] Sebastian Bach, Alexander Binder, Grégoire Montavon,
 1545 Frederick Klauschen, Klaus-Robert Müller, and Wojciech
 1546 Samek. On pixel-wise explanations for non-linear classifier
 1547 decisions by layer-wise relevance propagation. *PLoS ONE*,
 1548 10(7):e0130140, 2015. 1, 2
- 1549 [7] Chun-Hao Chang, Elliot Creager, Anna Goldenberg, and
 1550 David Duvenaud. Explaining image classifiers by counter-
 1551 factual generation. In *Proceedings of the 7th International
 1552 Conference on Learning Representations, ICLR*, 2019. 1, 2,
 1553 3, 4
- 1554 [8] Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia
 1555 Rudin, and Jonathan K Su. This looks like that: deep learning
 1556 for interpretable image recognition. *Advances in neural
 1557 information processing systems*, 32, 2019. 2
- 1558 [9] Piotr Dabkowski and Yarin Gal. Real time image saliency for
 1559 black box classifiers. In I. Guyon, U. Von Luxburg, S. Ben-
 1560 gio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett,
 1561 editors, *Advances in Neural Information Processing Systems*,
 1562 volume 30. Curran Associates, Inc., 2017. 1, 2
- 1563 [10] Ann-Kathrin Dombrowski, Maximilian Alber, Christopher
 1564 Anders, Marcel Ackermann, Klaus-Robert Müller, and Pan
 1565 Kessel. Explanations can be manipulated and geometry is
 1566 to blame. In H. Wallach, H. Larochelle, A. Beygelzimer, F.
 1567 d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in*

- 1566 *Neural Information Processing Systems*, volume 32. Curran
 1567 Associates, Inc., 2019. 2
- 1568 [11] Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Under-
 1569 standing deep networks via extremal perturbations and
 1570 smooth masks. In *Proceedings of the IEEE/CVF Interna-
 1571 tional Conference on Computer Vision (ICCV)*, October
 1572 2019. 1, 2, 4, 8, 10, 11, 12
- 1573 [12] Ruth C. Fong and Andrea Vedaldi. Interpretable explana-
 1574 tions of black boxes by meaningful perturbation. In *2017 IEEE
 1575 International Conference on Computer Vision (ICCV)*, pages
 1576 3449–3457, 2017. 1, 2, 4, 5, 8
- 1577 [13] Philipp Grohs. Continuous shearlet frames and resolution of
 1578 the wavefront set. *Monatsh. Math.*, 164(4):393–426, 2011.
 1579 5, 8
- 1580 [14] Kanghui Guo, Gitta Kutyniok, and Demetrio Labate. Sparse
 1581 multidimensional representations using anisotropic dilation
 1582 and shear operators. In *Wavelets and Splines*, pages 189–
 1583 201, Nashville, TN, 2005. Nashboro Press,. 4
- 1584 [15] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep
 1585 residual learning for image recognition. *2016 IEEE Confer-
 1586 ence on Computer Vision and Pattern Recognition (CVPR)*,
 1587 pages 770–778, 2016. 8, 9
- 1588 [16] Cosmas Heiß, Ron Levie, Cinjon Resnick, Gitta Kutyniok,
 1589 and Joan Bruna. In-distribution interpretability for challeng-
 1590 ing modalities. *Preprint arXiv:2007.00758*, 2020. 1, 2, 3,
 1591 4
- 1592 [17] Lars Hormander. *The Analysis of Linear Partial Differential
 1593 Operators. I, Distribution Theory and Fourier Analysis*.
 1594 Grundlehren Der Mathematischen Wissenschaften. Springer,
 1595 1990. 7
- 1596 [18] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie J. Cai,
 1597 James Wexler, Fernanda B. Viégas, and Rory Sayres. Inter-
 1598 pretability beyond feature attribution: Quantitative testing
 1599 with concept activation vectors (tcav). In *ICML*, 2018. 2
- 1600 [19] Stefan Kolek, Duc Anh Nguyen, Ron Levie, Joan Bruna, and
 1601 Gitta Kutyniok. Cartoon explanations of image classifiers. In
 1602 *European Conference of Computer Vision (ECCV)*, 2022. 1,
 2, 4, 5, 8
- 1603 [20] Stefan Kolek, Duc Anh Nguyen, Ron Levie, Joan Bruna, and
 1604 Gitta Kutyniok. A rate-distortion framework for explaining
 1605 black-box model decisions. In *xxAI - Beyond Explainable
 1606 AI*, page 91–115, 2022. 3, 4
- 1607 [21] G. Kutyniok and D. Labate. *Shearlets: Multiscale Analysis
 1608 for Multivariate Data*. Applied and Numerical Harmonic
 1609 Analysis. Birkhäuser Boston, 2012. 1
- 1610 [22] Gitta Kutyniok and Wang-Q Lim. Compactly supported
 1611 shearlets are optimally sparse. *Journal of Approximation
 1612 Theory*, 163(11):1564–1589, 2011. 7
- 1613 [23] FirstName LastName. The frobnicatable foo filter, 2014.
 1614 Face and Gesture submission ID 324. Supplied as supple-
 1615 mental material fg324.pdf. 12, 14
- 1616 [24] FirstName LastName. Frobnication tutorial, 2014. Supplied
 1617 as supplemental material tr.pdf. 12, 14
- 1618 [25] Scott M Lundberg and Su-In Lee. A unified approach to in-
 1619 terpreting model predictions. In I. Guyon, U. V. Luxburg,
 1620 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R.

- 1620 Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 1621 2017. 2
- 1622 [26] Jan Macdonald, Stephan Wäldchen, Sascha Hauch, and Gitta Kutyniok. A rate-distortion framework for explaining neural network decisions. *Preprint arXiv:1905.11092*, 2019. 1, 2
- 1623 [27] Stéphane Mallat. *A wavelet tour of signal processing (2. ed.)*. Academic Press, 1999. 5
- 1624 [28] Stéphane Mallat. *A Wavelet Tour of Signal Processing (Third Edition)*, chapter 11.3. Academic Press, third edition edition, 1625 2009. 4
- 1626 [29] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. In 1627 *BMVC*, 2018. 2
- 1628 [30] Rafael Reisenhofer, Johannes Kiefer, and Emily King. 1629 Shearlet-based detection of flame fronts. *Experiments in Fluids*, 57, 02 2016. 8
- 1630 [31] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 1631 “why should i trust you?”: Explaining the predictions of any 1632 classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1633 2016. 2
- 1634 [32] Ramprasaath R. Selvaraju, Abhishek Das, Ramakrishna 1635 Vedantam, Michael Cogswell, Devi Parikh, and Dhruv Batra. 1636 Grad-cam: Visual explanations from deep networks via 1637 gradient-based localization. *International Journal of Computer Vision*, 128:336–359, 2019. 1, 2
- 1638 [33] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 1639 Deep inside convolutional networks: Visualising 1640 image classification models and saliency maps. *Preprint arXiv:1312.6034*, 2014. 2
- 1641 [34] Karen Simonyan and Andrew Zisserman. Very deep 1642 convolutional networks for large-scale image recognition. *ICLR*, 1643 2015. 1
- 1644 [35] Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and 1645 Himabindu Lakkaraju. Fooling lime and shap: Adversarial 1646 attacks on post hoc explanation methods. In *AAAI/ACM Conference on Artificial Intelligence, Ethics, and Society (AIES)*, 1647 2020. 2
- 1648 [36] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, 1649 and Martin Wattenberg. Smoothgrad: removing noise by 1650 adding noise. In *Workshop on Visualization for Deep Learning, ICML*, 2017. 1, 2
- 1651 [37] J.T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller. Striving for simplicity: The all convolutional net. In 1652 *ICLR (workshop track)*, 2015. 1, 2
- 1653 [38] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic 1654 attribution for deep networks. In *Proceedings of the 34th 1655 International Conference on Machine Learning, ICML*, volume 70, 1656 page 3319–3328, 2017. 1, 2, 8
- 1657 [39] Jacopo Teneggi, Alexandre Luster, and Jeremias Sulam. Fast 1658 hierarchical games for image explanations. *IEEE Transactions on 1659 Pattern Analysis and Machine Intelligence*, pages 1–11, 2022. 1, 2
- 1660 [40] Sheng Yi, Demetrio Labate, Glenn R. Easley, and Hamid 1661 Krim. A shearlet approach to edge analysis and detection. 1662 *IEEE Transactions on Image Processing*, 18(5):929–941, 1663 2009. 8
- 1664 [41] Matthew D. Zeiler and R. Fergus. Visualizing and understanding 1665 convolutional networks. In *ECCV*, 2014. 1, 2
- 1666 [42] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. 1667 In *2016 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2921–2929, 2016. 1, 2
- 1668 The proof of Theorem
- 1669
- 1670
- 1671
- 1672
- 1673