
CRAFT: Concept Recursive Activation FacTORIZATION for Explainability

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Despite their considerable potential, concept-based explainability methods have
2 received relatively little attention, and explaining *what's* driving models' decisions
3 and *where* it's located in the input is still an open problem. To tackle this, we revisit
4 unsupervised concept extraction techniques for explaining the decisions of deep
5 neural networks and present CRAFT – a framework to generate concept-based
6 explanations for understanding individual predictions and the model's high-level
7 logic for whole classes. CRAFT takes advantage of a novel method for recursively
8 decomposing higher-level concepts into more elementary ones, combined with a
9 novel approach for better estimating the importance of identified concepts with
10 Sobol indices. Furthermore, we show how implicit differentiation can be used to
11 generate concept-wise attribution explanations for individual images. We further
12 demonstrate through fidelity metrics that our proposed concept importance estimation
13 technique is more faithful to the model than previous methods, and, through
14 human psychophysics experiments, we confirm that our recursive decomposition
15 can generate meaningful and accurate concepts. Finally, we illustrate CRAFT's
16 potential to enable the understanding of predictions of trained models on multiple
17 use-cases by producing meaningful concept-based explanations.*

18 1 Introduction

19 Interpreting the decisions of modern machine learning models such as neural networks remains a
20 major challenge. The need for robust and reliable explainability methods has never been more urgent
21 as machine learning is being applied to an ever increasing range of domains, including safety critical
22 ones. The application of the General Data Protection Regulation law (GDPR) [1] in the European
23 Union has drawn the attention of the general public to the rights they should have on their data.
24 This kickstarted a race for other needs, with more and more regulation agencies asking for the right
25 for AI decisions to be explainable to users – e.g. European AI act [2], EASA concepts for design
26 assurance [3].

27 In order to try to meet this need, an array of explainability methods have already been proposed. Most
28 of these methods aim at explaining what inputs (or pixels in an image) are driving the model's decision.
29 These so-called attribution methods yield heatmaps that indicate the importance of individual pixels.
30 Among the most notable ones is LIME [4], which was initially developed to try to locally – that is,
31 at an instance level – understand models' predictions to identify possible biases in vision models.
32 Multiple improvements have since been introduced – either by better harnessing the information
33 provided by gradients to estimate the importance of individual pixels [5, 6, 7, 8, 9, 10, 11, 12],
34 leveraging image perturbations to evaluate the sensitivity of a model's output [13, 14] or, more
35 recently, via the use of formal methods to generate explanations [15].

*Our code is available at anonymous.4open.science/r/craft-concept-explanation-4351.

36 However, all the aforementioned methods focus on one side of explainability – answering the question
37 of *where* – i.e., where in an image are the pixels that are critical to the decision located. They leave
38 the question of *what* – i.e., what visual features are actually driving decisions – entirely open. We
39 argue that this limitation is one of the main reasons why these methods fail in some cases to help
40 users, for instance, identify the source of a system’s bias or its failure cases as shown in [16]. Feature
41 visualization methods [17, 18] characterize the selectivity of individual neurons (or neural channels or
42 arbitrary directions in the neural activation space) via the synthesis of input stimuli which maximize
43 their responses and can partially answer this question. Still in this vein, [19, 20, 21] proposed to use
44 the training dataset to identify the samples that contribute the most to the model’s decision. Finally,
45 closer to our work, a new line of research has recently been initiated [22] based on high-level concepts.
46 The goal of this branch is to find humanly interpretable concepts in the activation space of a layer in a
47 neural network. This approach can give positive results, but in its original formulation, it requires
48 prior knowledge on the relevant concepts, and more importantly, the labeling of a dataset for each of
49 the concepts we want to extract. Hence, several works have proposed to automate the concept search
50 based only on the training dataset and without explicit human supervision. The most prominent
51 of these techniques, ACE [23], uses a combination of segmentation and clustering techniques, but
52 requires heuristics to remove outliers. This method unlocks the possibility of large scale concept
53 extraction without additional labeling or human supervision. Nevertheless, it suffers from several
54 problems: each segment can only belong to one cluster, the choice of the layer from which to retrieve
55 the concepts is not clear, and the amount of information lost during the outlier rejection phase can
56 be a cause of concern. More recently, [24] proposes to leverage matrix decompositions on internal
57 feature maps to discover concepts.

58 It is important to note that current work does not offer a link between their global and local ex-
59 planations, nor do they offer an answer to the question of which layer to choose to perform the
60 decomposition. Building up on these conclusions, we revisit these concept extraction techniques by
61 using Non-Negative Matrix Factorization (NMF) and propose 3 different ingredients to answer these
62 questions simultaneously, thereby introducing CRAFT, a new automatic concept extraction method.
63 We can summarize our main contributions as follows:

- 64 • A novel approach for the automated extraction of high-level concepts learned by deep neural
65 networks.
- 66 • A recursive procedure to automatically decompose concepts into sub-concepts, starting
67 with the last layer of the model and working our way inwards. We validate the benefit of
68 this recursivity – i.e. decomposing concepts into sub-concepts – with human psychophysical
69 experiments which show that (i) that the decomposition of a concept yields more coherent
70 sub-concepts (ii) the groups of points formed by these sub-concepts are more refined and
71 appear meaningful to humans (expert or non-expert).
- 72 • A novel technique to quantify the importance of individual concepts on a model’s predictions
73 using Sobol indices coming from the field of Sensitivity Analysis.
- 74 • A novel *Concept Attribution Map* (CAM) method to backpropagate each of the concept
75 values independently into the pixel space by leveraging the implicit function theorem,
76 allowing us to locate the concept in a given input image. This effectively unlocks the ability
77 to apply all the white-box [5, 6, 7, 8, 9, 12, 25] and black-box [4, 26, 13, 14] explainability
78 techniques in the literature to obtain concept-wise attribution maps.
- 79 • A demonstration of the approach combining local and global explanations to accurately
80 explain predictions and understand complex failure cases.

81 2 Related Work

82 **Explaining *where*** The widespread use of black-box machine learning methods including deep
83 convolutional neural networks in myriads of computer vision tasks prompted a need to understand
84 where in the input image the model looked to make predictions. These explanatory heatmaps can
85 be generated through completely different approaches depending on whether access to gradients is
86 provided. If it is indeed the case, there’s a plethora of different methods that harnesses intermediary
87 information inside the neural network to create these explanations [5, 8, 7, 28, 6, 29, 9, 12]. However,
88 they have been found to induce confirmation bias [30] and to be vulnerable to adversarial attacks [31].
89 Somewhat differently, there are other methods [10, 11] that harness gradients to optimize masks
90 to maximize the impact on the predictions, and thus determine the most important parts of the
91 input for the model. However, if only the input and its corresponding output are available, other



Figure 1: **CRAFT Results for the prediction ‘Chain Saw’**. First, our method uses NMF to extract from the train set (ILSVRC2012 [27]) the most relevant concepts used by the network (ResNet50V2). Then, the global influence of these concepts on predictions is measured using Sobol indices (right panel). Finally, the method provides local explanations through *Concept Attribution Maps* (heatmap associated to a concept, and computed using grad-CAM by backpropagating through the NMF concept values with implicit differentiation). Besides, concepts can be interpreted by looking at crops that maximize the NMF coefficients. For the class ‘Chain Saw’, the detected concepts seem to be: \mathcal{C}_0 for the chainsaw engine, \mathcal{C}_2 for the saw blade, \mathcal{C}_4 for the human head, \mathcal{C}_{18} for the vegetation, \mathcal{C}_{21} for the jeans and \mathcal{C}_{22} for the tree trunk.

92 techniques exist that enable the generation of attribution maps by locally estimating the importance
 93 of each input pixel: LIME [4], RISE [13], and more recently, an attribution method based on Sobol
 94 indices [14, 32]. Crucially, they propose to input perturbed versions of the example one wishes to
 95 explain and either construct a linear model to determine the importance of each region of the input,
 96 leverage Monte-Carlo methods to this end, or compute the Sobol indices [32] associated to them as a
 97 measure of their influence on the model. Concretely, we will be exploiting all this literature to locate
 98 the important parts of images with respect to what we will call “high-level” concepts by generating
 99 concept-wise attribution maps.

100 **Explaining «what»** There have been studies [33, 34] that indicate that CNNs trained on the
 101 ImageNet dataset [27] rely heavily on textures to classify, and largely disregard the shapes. For
 102 this reason, some researchers suggest that attribution maps might not be enough to explain models’
 103 predictions [17], and that explainability methods revealing the role of the textures are a must. Namely,
 104 in [18] and [17], explanations are generated as the inputs that would maximize the neural activation of
 105 a given layer with respect to a given class. However, these explanations may not be easily interpretable
 106 by humans. Finally, other approaches suggest to modify the structure of the neural network, either by
 107 constraining the convolutional layers to naturally provide visual explanations [35], or by forcing it to
 108 generate prototypes for the classes [36], but our main focus are post-hoc methods that can be applied
 109 to pre-trained neural networks and don’t need further training.

110 **Concept discovery.** In [22], Kim et al. proposed an alternative to explaining the *what*: they built a
 111 database with different concepts (such as “stripes”) to extract a concept vector in the latent space of a
 112 given layer. Then they proposed to estimate the importance of this concept vector using the directional
 113 derivative of the model’s predictions with respect to this concept vector. However, it is a supervised
 114 approach, and thus, only applicable when we have prior knowledge of the concepts in play. The
 115 natural extension of this idea is automatic discovery of concepts in an unsupervised fashion, without
 116 the need for prior knowledge or labelled concept datasets. As such, in [23], a technique is proposed to
 117 discover these “high-level” concepts: they perform segmentation at different resolutions on patches
 118 of images, cluster them and select the most significant based on perception and Testing with Concept
 119 Activation Vectors (TCAV) [22] scores. However, the quality of the result is highly dependent on the
 120 segmentation scheme and on the layer used for perception scores. Building up on this technique, [24]
 121 propose to generate a bank of concepts for each class by performing dimensionality reductions on the
 122 activation maps flattened over the channel dimension. Once the factorization done, the reconstruction
 123 of the activation of the image can then be interpreted as a combination of a set of concepts and a
 124 coefficient associated to these concepts. Not all factorization-based methods are equal though. Their
 125 large-scale human experiments show an interesting trend: Non-negative Matrix Factorization (NMF)
 126 is widely preferred over Principal Component Analysis (PCA) or ACE for generating meaningful

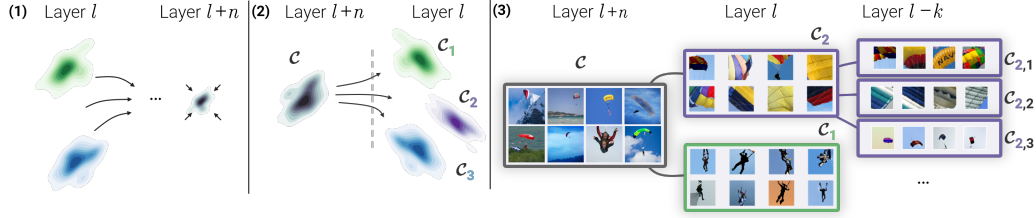


Figure 2: **(1) Neural Collapse (Amalgamation)** classifiers need to be able to linearly separate each class at the last layer, and to do this, the activations of the same class must merge during the forward pass until they all converge to the one-hot vector of the class in the logits layer. This may result in activations that are too concentrated to be broken down into meaningful concepts. **(2) Recursive process** When a concept is not understood (e.g., \mathcal{C}), we propose to decompose it into multiple sub-concepts (e.g., $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$) using the activations from an earlier layer to overcome the aforementioned neural collapse issue. **(3) Example of concept recursive decomposition** using CRAFT on the class ‘Parachute’ of ILSVRC2012 [27].

127 concepts for humans. Finally, [37] defines the notion of *completeness* of a concept bank and proposes
 128 a method to learn a complete set of concepts using Shapley values [26].

129 3 Overview of the method

130 In this section, we first describe our Concept Activations Factorization method by pointing out
 131 the differences that set our technique apart from previous work. We then proceed to introduce
 132 the three new ingredients that make up CRAFT: (1) a method to recursively decompose concepts
 133 into sub-concepts, (2) a new approach to better estimate the importance of extracted concepts
 134 and (3) how we unlock any attribution method to create *Concept Attribution Maps*, using implicit
 135 differentiation [38, 39, 40].

136 **Notation** In this work, we consider a general supervised learning setting, where $(x_1, \dots, x_n) \in \mathcal{X}^n$
 137 are n points and $(y_1, \dots, y_n) \in \mathcal{Y}^n$ their associated labels. Unless specified, all points are assumed
 138 to have the same labels. We are given a (machine-learnt) black-box predictor $f : \mathcal{X} \rightarrow \mathcal{Y}$, which at
 139 some test input x predicts the output $f(x)$. Without loss of generality, we assume that f is a neural
 140 network composed of k layers, and we denote $f(x) = h_k \circ h_{k-1} \circ \dots \circ h_1(x)$ with $h_l(x) \subseteq \mathbb{R}^p$
 141 being the intermediate activations for the layer l and $h_l(x)_i$ an activation for the same layer. Further,
 142 we require non-negative activations: $h_l(x)_i \geq 0 : \forall i \in \{1, \dots, p\}$, which amounts to choosing a layer
 143 whose activation function $\sigma_l(x) \geq 0$. In particular, this assumption is verified by any architecture that
 144 utilizes *ReLU*, but any non-negative activation function works. Finally, we denote $h_{l,k}$ the function
 145 going from the layer l to the output of the model f .

146 3.1 Concept Activations Factorization

147 As illustrated in Fig.3, we propose to use Non-negative matrix factorization activations to find a
 148 basis of concepts. Inspired by ACE [23], we will use sub-regions of images to attempt to identify
 149 coherent concepts. Instead of using segmentation – which naturally introduces artifacts due to the
 150 inpainting required by a baseline value –, we start by taking random crops of each image in our
 151 dataset (e.g, a set of points that the model predicts as belonging to the same class) to form an auxiliary
 152 dataset $\mathbf{X} \in \mathbb{R}^{n \times d}$ such that $\mathbf{X}_i = \tau(x_i)$ with τ a crop function. Given a layer l , we obtain the
 153 activations for the random crops $\mathbf{A} = h_l(\mathbf{X}) \in \mathbb{R}^{n \times p}$. In the case where f is a convolutional neural
 154 network, a global average pooling is applied on the activations. We recall that all the elements of \mathbf{A}
 155 are non-negative real numbers.

156 We are now ready to apply Non Negative Matrix Factorization (NMF) to decompose the positive
 157 activations \mathbf{A} , into a product of non-negative, low rank matrices $\mathbf{U}(\mathbf{A}) \in \mathbb{R}^{n \times r}$ and $\mathbf{W} \in \mathbb{R}^{p \times r}$,
 158 with:

$$\min_{\mathbf{U} \geq 0, \mathbf{W} \geq 0} \frac{1}{2} \|\mathbf{A} - \mathbf{U}\mathbf{W}^T\|_F^2 \quad (1)$$

159 Where $\|\cdot\|_F$ denotes the Frobenius norm. One of the appealing properties of NMF is the low
 160 rank constraint $r \ll \min(n, p)$. Simply put, NMF can be understood as the joint learning of \mathbf{W} ,
 161 a dictionary of CAVs – “concept bank” in Figure 3 – that maps a \mathbb{R}^p basis onto \mathbb{R}^r , and \mathbf{U} the
 162 coefficients of vectors \mathbf{A} expressed in this new basis. The minimization of the reconstruction error

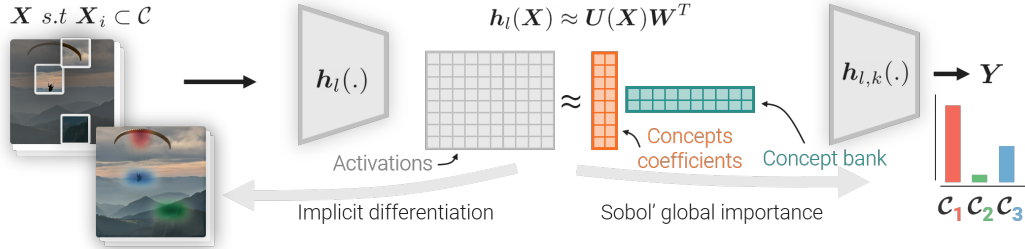


Figure 3: **Overview of CRAFT.** Starting from a set of crops X containing a concept \mathcal{C} (e.g., crops images of the class Parachute), we send random crops to a layer l to get activations $h_l(X)$. We then factorize the activation into two lower rank matrices, U and W . W is what we call a concept bank (a base of concepts), while U corresponds to the coefficients in this new basis. We then extend the method with 3 new ingredients: (1) the recursivity by proposing to re-decompose a concept (e.g., take a new set of point containing \mathcal{C}_1) at an earlier layer $l' < l$, (2) a better importance estimation using Sobol indices and (3) leveraging implicit differentiation to generate *Concept Attribution Maps* allowing to localize concepts in an image.

163 $\frac{1}{2}\|A - UW\|_F^2$ ensures that the new basis contains (mostly) relevant concepts. Intuitively, the
 164 non-negativity constraints $U \geq 0, W \geq 0$ encourage (i) the sparsity of W (useful for creating
 165 disentangled concepts), (ii) the sparsity of U (convenient for selecting a minimal set of useful
 166 concepts) and (iii) the imputation of missing data [41], which corresponds to the sparsity pattern
 167 of *post-ReLU* activations A . We shall also note that each original activation A_i coming from the
 168 input x_i can be approximated by its reconstruction $h_l(\tau(x_i)) = U_i W^T = \sum_{j=1}^r U_{i,j} W_j^T$. This
 169 approach is attractive as each activation can be understood as a composition of concepts.

170 While other methods in the literature solve a similar problem (such as low rank factorization using
 171 SVD or ICA), the NMF has stepped up as both fast, effective and is known to yield meaningful
 172 concepts to humans [42, 43, 24]. Finally, once the concept bank W is precomputed, we can associate
 173 the concept coefficients $U(x)$ to any new input x (e.g a full image) by solving the underlying
 174 Non-Negative Least Squares (NNLS) problem $\min_{U \geq 0} \frac{1}{2}\|h_l(x) - U(x)W^T\|_F^2$, and therefore
 175 have its decomposition in the concept base.

176 In essence, the core of our method can be summarized as follows: using a set of images, we re-interpret
 177 their embedding at a given layer l as a composition of concepts that can be easily understood by
 178 humans. In the next section we show how we can recursively apply concept activation factorizations
 179 on a layer $l' < l$ for an image containing a previously computed concept.

180 3.2 Ingredient 1: A Recursive Flavor

181 One of the most apparent issues in previous works [23, 24] is the choice of the layer at which the
 182 activation maps are extracted. Depending on this, certain concepts start getting amalgamated [44]
 183 into one, resulting in incoherent and indecipherable clusters, as illustrated in Fig 2. We posit that
 184 this can be solved by iteratively applying our decomposition at different layer-depths, and for the
 185 concepts that remain difficult to understand, look for their sub-concepts at earlier layers by isolating
 186 the images that contain them. This allows us to build hierarchies of concepts for each class.

187 We offer a simple solution consisting of reapplying our method to a concept by performing a second
 188 step of Concept Activation Factorization on a set of points that contain the concept \mathcal{C} in order to
 189 refine it and create sub-concepts (e.g., decompose \mathcal{C} into $\{\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3\}$) see see Fig.2 for an illustrative
 190 example. Note that we generalize current methods in the sense that taking points (x_1, \dots, x_n) that
 191 are clustered in the logits layer (belonging to the same class) and decomposing them in a previous
 192 layer – as done in [23, 24] – is a valid recursive step. For a more general case, let us assume that
 193 a set of points that contain a common concept is obtained using a first step of Concept Activation
 194 Factorization. We then look for a set of points with a high coefficient for the concept of our choice
 195 to perform the next factorization. Formally, with a factorization for a layer l UW^T and a concept
 196 index i , this set of points is defined as $\mathcal{C} = \{\tau(x_j) : U(A_j)_i \geq \lambda\}$ In practice, we assume λ to be
 197 equal to the 90th percentile of the values of U_i . Given this new set of points, we can then re-apply the
 198 Concept Matrix Factorization method to a earlier layer l' – with $l' < l$ – to obtain the sub-concept's
 199 decomposition from the initial concept.

200 **3.3 Ingredient 2: Sobol indices for enhanced concept importance estimation**

201 A common concern with concept extraction methods is that what makes sense to humans is not
 202 necessarily what is being used by the model to predict. To avoid this kind of confirmation bias during
 203 our concept analysis phase, we can estimate the global importance of the extracted concepts. To do
 204 so, [22] proposed an estimator based on directional derivatives: the partial derivative of the model
 205 output with respect to the concept vector. While this measure is theoretically well founded, it relies
 206 on the same principle as gradient-based methods, and thus, suffers from the same pitfalls: neural
 207 network models have noisy gradients [5, 7]. Hence, the farther the chosen layer is from the output,
 208 the noisier the directional derivative score will be.

209 Since we essentially want to know which concept has the greatest effect on the output of the model,
 210 it is natural to consider the field of Sensitivity Analysis [45, 46, 32, 47]. In this section, we briefly
 211 recall the classical total Sobol indices and how to apply it to our problem. The complete derivation of
 212 the Sobol-Hoeffding decomposition is presented in the appendix D.

213 Formally, we place ourselves at layer l and perform our Concept Activation Factorization, providing us
 214 with U, W . A natural way to estimate the importance of a concept U_i is to measure the fluctuation of
 215 the model’s output $h_{l,k}(UW^T)$ in response to meaningful perturbations of the concept coefficient U_i .
 216 Concretely, with $M = (M_1, \dots, M_r) \in [0, 1]^r$, here an i.i.d sequence of real-valued random variables,
 217 we introduce a concept fluctuation to reconstruct a perturbed activation $\tilde{A} = (U \odot M)W^T$ (e.g.,
 218 the masks can be used to put a concept value to zero). We can then propagate this perturbed
 219 activation to the model output $Y = h_{l,k}(\tilde{A})$. Thus, an important concept will have a large variance
 220 on the model output while an unused concept will barely change it.

221 Finally, we can capture the importance that a concept might have as a main effect – along with its
 222 interactions with other concepts – on the model’s output by calculating the expected variance that
 223 would remain if all the indices of the masks except the M_i were to be fixed. This yields the general
 224 definition of the Total Sobol indices.

225 **Definition 3.1 (Total Sobol indices).** *The total Sobol index S_{T_i} , which measures the contribution*
 226 *of a concept U_i as well as its interactions of any order with any other concepts to the model output*
 227 *variance, is given by:*

$$S_{T_i} = \frac{\mathbb{E}_{M_{\sim i}}(\mathbb{V}_{M_i}(Y|M_{\sim i}))}{\mathbb{V}(Y)} = \frac{\mathbb{E}_{M_{\sim i}}(\mathbb{V}_{M_i}(h_{l,k}((U \odot M)W^T)|M_{\sim i}))}{\mathbb{V}((U \odot M)W^T)} \quad (2)$$

228 In a practical way, this index can be calculated efficiently [48, 49, 50, 51, 52], more details on the
 229 sampling (Quasi-Monte Carlo) and the estimator used are left in appendix D.

230 **3.4 Ingredient 3: Unlocking Concept Attribution Map**

231 Attribution methods are useful for determining the regions deemed important by the model for
 232 the decision, but they lack the information about what exactly triggered it. We have seen that we
 233 can already extract this information from the matrices U and W , but as it is, we cannot know to
 234 which part of the image the model associates each concept, and thus, better comprehend the model’s
 235 decisions. In this section, we will show how we can unlock the set of attribution methods (forward
 236 and backward mode) to find where a concept is located in the input image (see Fig.1). Forward
 237 attribution methods don’t rely on gradients and only use inference information, whereas backward
 238 methods require to back-propagate through the network’s layers. By application of the chain rule,
 239 computing $\frac{\partial U}{\partial x}$ requires access to $\frac{\partial U}{\partial A}$.

240 To do so, it could be tempting to solve the linear system $UW^T = A$. However, this problem is
 241 ill-posed since W^T is low rank. A standard approach is to calculate the Moore-Penrose pseudo-
 242 inverse $(W^T)^+$, which solves rank deficient systems by looking at the minimum norm solution [53].
 243 In practice $(W^T)^+$ is computed with the Singular Value Decomposition (SVD) of W^T . Unfor-
 244 tunately, SVD is also the solution to the *unstructured minimization* of $\frac{1}{2}\|A - UW^T\|_F^2$ by the
 245 Eckart–Young–Mirsky theorem [54]. Hence, the non negativity constraints – i.e $U \geq 0, W \geq 0$ – of
 246 the NMF are ignored, which prevents approaches based on solving $U^T W = A^T$ from succeeding.
 247 Other issues stem from the fact that the U, W decomposition is generally not unique.

248 Our third contribution consists on tackling this problem to allow the use of attribution methods – i.e.
 249 *Concept Attribution Maps* – by proposing a strategy to differentiate through the NMF layer.

250 **Implicit differentiation of NMF layers** The NMF problem 1 is NP-hard [55], and it is not convex
 251 with respect to the input pair (U, W) . However, fixing the value of one of the two factors and

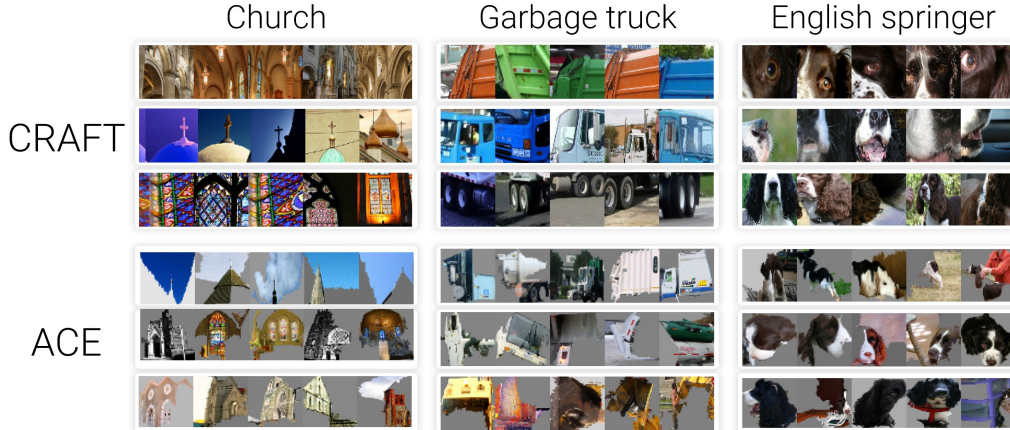


Figure 4: **Qualitative comparison.** We compare concepts found by our method (top) to those extracted with ACE [23] (bottom) for the classes *Church*, *Garbage truck* and *English springer* from ILSVRC2012 [27].

252 optimizing the other turns the NMF formulation into a pair of Non Negative Least Squares (NNLS)
 253 problems (see Equation 3), which are convex. This ensures that alternating minimization (a standard
 254 approach for NMF) of (\mathbf{U}, \mathbf{W}) factors will eventually reach a local (and global) minimum:

$$\mathbf{U}_{t+1} = \arg \min_{\mathbf{U} \geq 0} \frac{1}{2} \|\mathbf{A} - \mathbf{U}\mathbf{W}_t^T\|_F^2 \quad \mathbf{W}_{t+1} = \arg \min_{\mathbf{W} \geq 0} \frac{1}{2} \|\mathbf{A} - \mathbf{U}_t\mathbf{W}^T\|_F^2 \quad (3)$$

255 Each of the NNLS problem fulfills the KKT conditions[56, 57], which can be encoded in the so-called
 256 *optimality function* \mathbf{F} , see Equation 10 Appendix C.2. The implicit function theorem [39] allows us
 257 to use implicit differentiation [38, 39, 58] to efficiently compute the Jacobians $\frac{\partial \mathbf{U}}{\partial \mathbf{A}}$ and $\frac{\partial \mathbf{W}}{\partial \mathbf{A}}$ without
 258 requiring to back-propagate through each of the iterations of the NMF solver:

$$\frac{\partial(\mathbf{U}, \mathbf{W}, \bar{\mathbf{U}}, \bar{\mathbf{W}})}{\partial \mathbf{A}} = -(\partial_1 \mathbf{F})^{-1} \partial_2 \mathbf{F} \quad (4)$$

259 However, this requires the dual variables $\bar{\mathbf{U}}$ and $\bar{\mathbf{W}}$, which are not computed by Scikit-learn’s [59]
 260 popular implementation[†]. Consequently, we leverage the work of [62] and we re-implement our own
 261 solver with Jaxopt [40] based on ADMM [63], a GPU friendly algorithm (see Appendix C.2).

262 We start by performing Concept Activations Factorization – i.e we precompute the concept bank \mathbf{W}
 263 by solving the NMF. Concept Attribution Maps of a new input \mathbf{x} are calculated by solving the NNLS
 264 problem $\min_{\mathbf{U} \geq 0} \frac{1}{2} \|\mathbf{h}_l(\mathbf{x}) - \mathbf{U}\mathbf{W}^T\|_F^2$. The implicit differentiation of NMF layer $\frac{\partial \mathbf{U}}{\partial \mathbf{A}}$ is integrated
 265 into classical back-propagation to obtain $\frac{\partial \mathbf{U}}{\partial \mathbf{x}}$. Most interestingly, this technical advance unlocks all
 266 white-box explainability methods [5, 6, 7, 8, 9, 12] to generate concept-wise attribution maps and
 267 trace the part of the image that triggered the detection of the concept. Additionally, it is even possible
 268 to employ black-box methods [4, 13, 26, 14] since it only amounts to solving an NNLS problem.

269 4 Experimental evaluation

270 We used CRAFT to explain a ResNet50V2 trained on the ILSVRC2012 [27] data set (ImageNet). We
 271 selected a subset of 10 classes, each containing 1000 images (those recommended by ImageNette[‡]).
 272 In all of our experiments, $r = 25$, like in [23] and the cropping function τ consists on randomly
 273 choosing 10 square 64×64 patches for each image. We start by qualitatively validating CRAFT by
 274 showing that: (1) the method yields concepts that are easy to interpret (see Fig. 4), (2) the combination
 275 of local and global explanations allows to explain complex failure cases otherwise unexplainable
 276 with only the attribution methods (see Fig. 5). Then, we validate independently the new ingredients
 277 brought by the method by showing quantitatively that (3) recursivity allows us to refine concepts,
 278 making them more meaningful to humans with the help of two psychophysics experiments, and (4)
 279 Sobol indices allow for a better estimation of concept importance. Additional experiments, including

[†]Scikit-learn uses a Block coordinate descent algorithm [60, 61], with a randomized SVD initialization.

[‡]<https://github.com/fastai/imagenette>

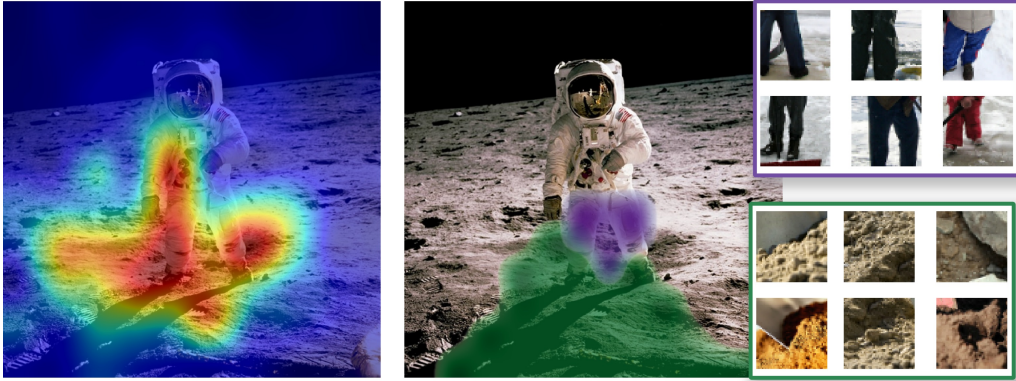


Figure 5: **This is a Shovel.** We compare a heatmap generated by RISE [13] (left) with the *Concept Attribution Maps* generated with our implicit differentiation pipeline and Grad-CAM (right) on the explanations of the two most influential concepts that drove the ResNet50’s decision. We found a first concept that seems to be associated with textures of dirt commonly found in the images of the class *Shovel*. The second concept elucidated by CRAFT is located on the astronaut’s pants, which he confuses with the ski suits of people clearing snow from their driveway with a shovel.

280 a sanity check and an example of activation maximization (Deep dream) on the concept bank, as well
 281 as many other examples of local explanations for randomly picked images from ILSVRC2012, are
 282 included in appendix B.

283 We leave a discussion on the limitations of this method and on the broader impact in appendix A.

284 4.1 Example of CRAFT concepts

285 Figure 4 compares the examples of concepts found by CRAFT against those found by ACE [23] for 3
 286 classes of Imagenet. For each class the concepts are ordered by importance (the highest being the
 287 most important). ACE uses a clustering technique and TCAV to estimate importance, while CRAFT
 288 uses the method introduced in 3 and Sobol to estimate importance. These examples illustrate one
 289 of the weaknesses of ACE: the segmentation used can introduce biases through the baseline value
 290 used [64, 10]. The concepts found by CRAFT seem distinct: (vault, cross, stained glass) for the
 291 Church class, (dumpster, truck door, two-wheeler) for the garbage truck, and (eyes, nose, fluffy ears)
 292 for the English Springer. More examples can be found in the appendix.

293 4.2 Explaining complex failure cases

294 One of the goals of explainability is to
 295 explain the failure cases of the models
 296 studied. Figure 5 shows an example
 297 of an incorrect prediction: the model
 298 in question – here still a ResNet50 –
 299 predicts ‘shovel’. Moreover, the at-
 300 tribution method on the left – here
 301 RISE [13] – does not tell us much
 302 except that the evidence for shovel
 303 seems to be located at the level of the
 304 ground and the lower torso and legs
 305 of the astronaut. With CRAFT, we can however study the concepts found by the model at these
 306 locations. There are two of them: the first concept in green, aims at the lunar ground and refers to the
 307 rocks often seen next to shovels in the dataset. The second concept in purple is aimed at the legs of
 308 the astronaut and refers to the legs of a person, often in a ski suit, which he takes for the astronaut’s.

309 4.3 Validation of Recursivity

310 To evaluate the meaningfulness of the extracted high-level concepts, we performed psychophysical
 311 experiments with human subjects, to whom we requested to answer a survey in two phases. Further-
 312 more, we distinguished two different audiences: on the one hand, experts in machine learning, and on

	Experts ($n = 36$)	Laymen ($n = 37$)
<i>Intruder</i>		
Acc. Concept	70.19%	61.08%
Acc. Sub-Concept	74.81% ($p = 0.18$)	67.03% ($p = 0.043$)
<i>Binary choice</i>		
Sub-Concept	76.1% ($p < 0.001$)	74.95% ($p < 0.001$)
Odds Ratios	3.53	2.99

Table 1: **Results from the psychophysics experiments.**

313 the other hand, people with no particular knowledge in computer vision. Both groups of participants
 314 were volunteers and didn't receive any monetary compensation. Some examples of the developed
 315 interface are available the appendix E.

316 **Intruder detection experiment,** we make users identify the intruder out of a series of five segments
 317 belonging to a certain class, with the odd one being taken from a different concept but from the same
 318 class. Now, we compare the results of this intruder detection with a concept (e.g., C_1) coming from a
 319 layer l and one of its sub-concepts (e.g., C_{12} in Fig.2) extracted using our recursive method. If the
 320 concept (or sub-concept) is meaningful, then it should be easy for the users to find the intruder. Table 1
 321 summarizes our results, showing that indeed both concepts and sub-concepts are meaningful, and
 322 that recursivity can lead to a slightly higher understanding of the generated concepts (significant for
 323 non-experts, not significant for experts) and might suggest a way to make concepts more interpretable.

324 **Binary choice experiment,** In order to test the improvement of the meaningfulness of the sub-
 325 concept generated with recursivity with respect to the larger parent concept, we showed participants a
 326 segment belonging to a subcluster and to the parent cluster (e.g., $\tau(x) \subset C_{11} \subset C_1$) without specifying
 327 why those images are grouped together. We then we asked which of the two clusters (i.e., C_{11} or C_1)
 328 seemed to accommodate the image the best. If our hypothesis is correct, then the concept refinement
 329 brought by recursivity should help form more coherent clusters. The results in Table 1 are satisfying,
 330 since in both the expert and non-expert groups, the participants chose the sub-cluster by more than 74%
 331 of the times. We measure the significance of our results by fitting a binomial logistic regression to our
 332 data, and we find that both groups are more likely to choose the sub-concept cluster (at a $p < 0.001$).
 333

334 4.4 Fidelity analysis

335 We propose to simultaneously verify that
 336 the concepts are faithful to the model and
 337 that the concept importance estimator per-
 338 forms better than TCAV [22] by using the
 339 fidelity metrics introduced in [23, 24].
 340 These metrics are similar to the one used
 341 for attribution methods, which consist on
 342 studying the change of the logit score when
 343 removing/adding pixels considered impor-
 344 tant. Nevertheless, we do not make these
 345 modifications in the pixel space but in the
 346 concept space: once U, W are computed,
 347 we reconstruct the matrix $A \approx UW^T$ us-
 348 ing only the most important concept (or
 349 removing the most important concept for deletion), and study the score in output of the model. As
 350 can be seen from Fig. 6, ranking the extracted concepts using Sobol's importance score results in
 351 much steeper curves than when they are sorted by their TCAV scores. We confirm these results with
 352 other matrix factorization techniques (PCA, ICA, RCA) in the Appendix F.

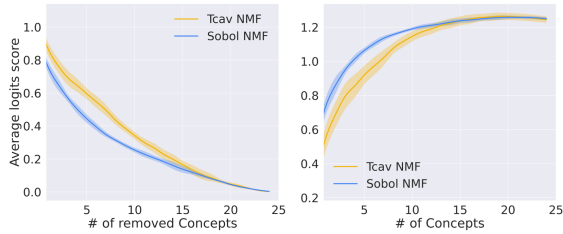


Figure 6: **(Left)** Deletion curve (lower is better). **(Right)** Insertion curves (higher is better). Whether in deletion or insertion, the score – calculated on more than 100,000 images – shows that using Sobol indices yield to better estimates of important concepts.

353 5 Conclusion

354 In this paper, we introduced a method for automatically extracting human-scrutable concepts from
 355 Deep Neural Network: CRAFT. Our method allows to explain a pre-trained model both in a per-class
 356 and per-instance basis by highlighting both *what* the model saw when predicting the class label
 357 and *where* it is located, which, as we have shown, exhibits complementary benefits. The approach
 358 relies on three novel ingredients: 1) exploiting the recursive nature of the feature extraction chains in
 359 CNNs to find decompositions where each concept is clearly understandable; 2) measuring concept
 360 importance through Sobol indices to more accurately identify which concepts influence a model's
 361 decision for a given class; and 3) harnessing implicit differentiation to backpropagate through NMF
 362 blocks, thus enabling the use of any attribution method to create concept-wise local explanations that
 363 we call *Concept Attribution Maps*. Human experiments confirmed the validity of the approach and
 364 that concepts identified by CRAFT are meaningful. We hope that this work will guide further efforts
 365 in the search for concept-based explainability methods and that further connections between local
 366 and global explanations will be made.

367 References

- 368 [1] Margot E Kaminski and Jennifer M Urban. The right to contest ai. *Columbia Law Review*,
369 121(7):1957–2048, 2021. 1
- 370 [2] Mauritz Kop. Eu artificial intelligence act: The european approach to ai. Stanford-Vienna
371 Transatlantic Technology Law Forum, Transatlantic Antitrust . . . , 2021. 1
- 372 [3] Christoph Torens, Umut Durak, and Johann C Dauer. Guidelines and regulatory framework for
373 machine learning in aviation. In *AIAA Scitech 2022 Forum*, page 1132, 2022. 1
- 374 [4] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining
375 the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international
376 conference on knowledge discovery and data mining*, pages 1135–1144, 2016. 1, 2, 3, 7, 24
- 377 [5] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smooth-
378 grad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017. 1, 2, 6, 7,
379 28
- 380 [6] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:
381 Visualising image classification models and saliency maps. In *In Workshop at International
382 Conference on Learning Representations*. Citeseer, 2014. 1, 2, 7
- 383 [7] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In
384 *International conference on machine learning*, pages 3319–3328. PMLR, 2017. 1, 2, 6, 7
- 385 [8] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi
386 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based
387 localization. In *Proceedings of the IEEE international conference on computer vision*, pages
388 618–626, 2017. 1, 2, 7
- 389 [9] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving
390 for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014. 1, 2, 7
- 391 [10] Ruth C Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful
392 perturbation. In *Proceedings of the IEEE international conference on computer vision*, pages
393 3429–3437, 2017. 1, 2, 8
- 394 [11] Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Understanding deep networks via extremal
395 perturbations and smooth masks. In *Proceedings of the IEEE/CVF international conference on
396 computer vision*, pages 2950–2958, 2019. 1, 2
- 397 [12] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep
398 features for discriminative localization. In *Proceedings of the IEEE conference on computer
399 vision and pattern recognition*, pages 2921–2929, 2016. 1, 2, 7
- 400 [13] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation
401 of black-box models. *arXiv preprint arXiv:1806.07421*, 2018. 1, 2, 3, 7, 8, 24, 26
- 402 [14] Thomas Fel, Rémi Cadène, Mathieu Chalvidal, Matthieu Cord, David Vigouroux, and Thomas
403 Serre. Look at the variance! efficient black-box explanations with sobol-based sensitivity
404 analysis. *Advances in Neural Information Processing Systems*, 34, 2021. 1, 2, 3, 7, 24
- 405 [15] Thomas Fel, Melanie Ducoffe, David Vigouroux, Remi Cadene, Mikael Capelle, Claire
406 Nicodeme, and Thomas Serre. Don't lie to me! robust and efficient explainability with
407 verified perturbation analysis. *arXiv preprint arXiv:2202.07728*, 2022. 1
- 408 [16] Thomas Fel, Julien Colin, Rémi Cadène, and Thomas Serre. What i cannot predict, i do not
409 understand: A human-centered evaluation framework for explainability methods. *arXiv preprint
410 arXiv:2112.04417*, 2021. 2
- 411 [17] Chris Olah, Arvind Satyanarayan, Ian Johnson, Shan Carter, Ludwig Schubert, Katherine
412 Ye, and Alexander Mordvintsev. The building blocks of interpretability. *Distill*, 2018.
413 <https://distill.pub/2018/building-blocks>. 2, 3

- 414 [18] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2017.
415 <https://distill.pub/2017/feature-visualization>. 2, 3, 16, 21
- 416 [19] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions.
417 In *International conference on machine learning*, pages 1885–1894. PMLR, 2017. 2
- 418 [20] Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training
419 data influence by tracing gradient descent. *Advances in Neural Information Processing Systems*,
420 33:19920–19930, 2020. 2
- 421 [21] Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point
422 selection for explaining deep neural networks. *Advances in neural information processing*
423 *systems*, 31, 2018. 2
- 424 [22] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al.
425 Interpretability beyond feature attribution: Quantitative testing with concept activation vectors
426 (tcav). In *International conference on machine learning*, pages 2668–2677. PMLR, 2018. 2, 3,
427 6, 9, 21, 26
- 428 [23] Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-
429 based explanations. *Advances in Neural Information Processing Systems*, 32, 2019. 2, 3, 4, 5, 7,
430 8, 9, 26
- 431 [24] Ruihan Zhang, Prashan Madumal, Tim Miller, Krista A Ehinger, and Benjamin IP Rubinstein.
432 Invertible concept-based explanations for cnn models with non-negative concept activation
433 vectors. *arXiv preprint arXiv:2006.15417*, 2020. 2, 3, 5, 9, 26
- 434 [25] Thomas Fel, Mélanie Ducoffe, David Vigouroux, Rémi Cadène, Mikael Capelle, Claire
435 Nicodème, and Thomas Serre. Don’t lie to me! robust and efficient explainability with
436 verified perturbation analysis. *arXiv preprint arXiv:2202.07728*, 2022. 2
- 437 [26] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions.
438 *Advances in neural information processing systems*, 30, 2017. 2, 4, 7
- 439 [27] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
440 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
441 *recognition*, pages 248–255. Ieee, 2009. 3, 4, 7, 26, 27
- 442 [28] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert
443 Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by
444 layer-wise relevance propagation. *PLOS ONE*, 10(7):1–46, 07 2015. 2
- 445 [29] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In
446 *European conference on computer vision*, pages 818–833. Springer, 2014. 2
- 447 [30] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim.
448 Sanity checks for saliency maps. *Advances in neural information processing systems*, 31, 2018.
449 2, 28
- 450 [31] Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. Fooling
451 lime and shap: Adversarial attacks on post hoc explanation methods. In *Proceedings of the*
452 *AAAI/ACM Conference on AI, Ethics, and Society*, pages 180–186, 2020. 2
- 453 [32] Ilya M Sobol. Global sensitivity indices for nonlinear mathematical models and their monte
454 carlo estimates. *Mathematics and computers in simulation*, 55(1-3):271–280, 2001. 3, 6
- 455 [33] Katherine Hermann, Ting Chen, and Simon Kornblith. The origins and prevalence of texture
456 bias in convolutional neural networks. *Advances in Neural Information Processing Systems*,
457 33:19000–19015, 2020. 3
- 458 [34] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann,
459 and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias
460 improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. 3

- 461 [35] Quanshi Zhang, Ying Nian Wu, and Song-Chun Zhu. Interpretable convolutional neural
462 networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
463 pages 8827–8836, 2018. 3
- 464 [36] Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K Su.
465 This looks like that: deep learning for interpretable image recognition. *Advances in neural*
466 *information processing systems*, 32, 2019. 3, 16
- 467 [37] Chih-Kuan Yeh, Been Kim, Sercan Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar.
468 On completeness-aware concept-based explanations in deep neural networks. *Advances in*
469 *Neural Information Processing Systems*, 33:20554–20565, 2020. 4
- 470 [38] Steven George Krantz and Harold R Parks. *The implicit function theorem: history, theory, and*
471 *applications*. Springer Science & Business Media, 2002. 4, 7, 24
- 472 [39] Andreas Griewank and Andrea Walther. *Evaluating derivatives: principles and techniques of*
473 *algorithmic differentiation*. SIAM, 2008. 4, 7, 24
- 474 [40] Mathieu Blondel, Quentin Berthet, Marco Cuturi, Roy Frostig, Stephan Hoyer, Felipe Llinares-
475 López, Fabian Pedregosa, and Jean-Philippe Vert. Efficient and modular implicit differentiation.
476 *arXiv preprint arXiv:2105.15183*, 2021. 4, 7, 24
- 477 [41] Bin Ren, Laurent Pueyo, Christine Chen, Élodie Choquet, John H Debes, Gaspard Duchêne,
478 François Ménard, and Marshall D Perrin. Using data imputation for signal separation in
479 high-contrast imaging. *The Astrophysical Journal*, 892(2):74, 2020. 5
- 480 [42] Yu-Xiong Wang and Yu-Jin Zhang. Nonnegative matrix factorization: A comprehensive review.
481 *IEEE Transactions on Knowledge and Data Engineering*, 25(6):1336–1353, 2013. 5
- 482 [43] Xiao Fu, Kejun Huang, Nicholas D Sidiropoulos, and Wing-Kin Ma. Nonnegative matrix
483 factorization for signal and data analytics: Identifiability, algorithms, and applications. *IEEE*
484 *Signal Process. Mag.*, 36(2):59–80, 2019. 5
- 485 [44] Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the
486 terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*,
487 117(40):24652–24663, 2020. 5
- 488 [45] Bertrand Iooss and Paul Lemaître. A review on global sensitivity analysis methods. In *Uncer-*
489 *tainty management in simulation-optimization of complex systems*, pages 101–122. Springer,
490 2015. 6
- 491 [46] Ilya M Sobol. Sensitivity analysis for non-linear mathematical models. *Mathematical modelling*
492 *and computational experiment*, 1:407–414, 1993. 6, 25
- 493 [47] RI Cukier, CM Fortuin, Kurt E Shuler, AG Petschek, and J Ho Schaibly. Study of the sensitivity
494 of coupled reaction systems to uncertainties in rate coefficients. i theory. *The Journal of chemical*
495 *physics*, 59(8):3873–3878, 1973. 6
- 496 [48] Andrea Saltelli, Paola Annoni, Ivano Azzini, Francesca Campolongo, Marco Ratto, and Stefano
497 Tarantola. Variance based sensitivity analysis of model output. design and estimator for the total
498 sensitivity index. *Computer physics communications*, 181(2):259–270, 2010. 6
- 499 [49] Amandine Marrel, Bertrand Iooss, Beatrice Laurent, and Olivier Roustant. Calculations of
500 sobol indices for the gaussian process metamodel. *Reliability Engineering & System Safety*,
501 94(3):742–751, 2009. 6
- 502 [50] Alexandre Janon, Thierry Klein, Agnes Lagnoux, Maëlle Nodet, and Clémentine Prieur. Asymp-
503 totic normality and efficiency of two sobol index estimators. *ESAIM: Probability and Statistics*,
504 18:342–364, 2014. 6, 25
- 505 [51] Art B Owen. Better estimation of small sobol’sensitivity indices. *ACM Transactions on*
506 *Modeling and Computer Simulation (TOMACS)*, 23(2):1–17, 2013. 6

- 507 [52] Stefano Tarantola, Debora Gatelli, and Thierry Alex Mara. Random balance designs for the
508 estimation of first order global sensitivity indices. *Reliability Engineering & System Safety*,
509 91(6):717–727, 2006. 6
- 510 [53] João Carlos Alves Barata and Mahir Saleh Hussein. The moore–penrose pseudoinverse: A
511 tutorial review of the theory. *Brazilian Journal of Physics*, 42(1):146–165, 2012. 6
- 512 [54] Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank.
513 *Psychometrika*, 1(3):211–218, 1936. 6
- 514 [55] Stephen A Vavasis. On the complexity of nonnegative matrix factorization. *SIAM Journal on*
515 *Optimization*, 20(3):1364–1377, 2010. 6
- 516 [56] William Karush. Minima of functions of several variables with inequalities as side constraints.
517 *M. Sc. Dissertation. Dept. of Mathematics, Univ. of Chicago*, 1939. 7, 23
- 518 [57] Harold W Kuhn and Albert W Tucker. Nonlinear programming proceedings of the second
519 berkeley symposium on mathematical statistics and probability. *Neyman*, pages 481–492, 1951.
520 7, 23
- 521 [58] Bradley M Bell and James V Burke. Algorithmic differentiation of implicit functions and
522 optimal values. In *Advances in Automatic Differentiation*, pages 67–77. Springer, 2008. 7, 24
- 523 [59] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
524 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-
525 learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830,
526 2011. 7
- 527 [60] Andrzej Cichocki and Anh-Huy Phan. Fast local algorithms for large scale nonnegative matrix
528 and tensor factorizations. *IEICE transactions on fundamentals of electronics, communications*
529 *and computer sciences*, 92(3):708–721, 2009. 7, 24
- 530 [61] Cédric Févotte and Jérôme Idier. Algorithms for nonnegative matrix factorization with the
531 β -divergence. *Neural computation*, 23(9):2421–2456, 2011. 7, 24
- 532 [62] Kejun Huang, Nicholas D Sidiropoulos, and Athanasios P Liavas. A flexible and efficient
533 algorithmic framework for constrained matrix and tensor factorization. *IEEE Transactions on*
534 *Signal Processing*, 64(19):5052–5065, 2016. 7, 23
- 535 [63] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed opti-
536 mization and statistical learning via the alternating direction method of multipliers. *Foundations*
537 *and Trends® in Machine learning*, 3(1):1–122, 2011. 7, 16, 23, 24
- 538 [64] Pascal Sturmfels, Scott Lundberg, and Su-In Lee. Visualizing the impact of feature attribution
539 baselines. *Distill*, 5(1):e22, 2020. 8
- 540 [65] Kaiyu Yang, Jacqueline Yau, Li Fei-Fei, Jia Deng, and Olga Russakovsky. A study of face
541 obfuscation in imagenet. *arXiv preprint arXiv:2103.06191*, 2021. 14
- 542 [66] Yunhao Ge, Yao Xiao, Zhi Xu, Meng Zheng, Srikrishna Karanam, Terrence Chen, Laurent
543 Itti, and Ziyang Wu. A peek into the reasoning of neural networks: Interpreting with structural
544 visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
545 *Recognition*, pages 2195–2204, 2021. 16
- 546 [67] Thomas Fel, Lucas Hervier, David Vigouroux, Antonin Poche, Justin Plakoo, Remi Cadene,
547 Mathieu Chalvidal, Julien Colin, Thibaut Boissin, Louis Béthune, Agustin Picard, Claire
548 Nicodeme, Laurent Gardes, Gregory Flandin, and Thomas Serre. Xplique: A deep learning
549 explainability toolbox. *Workshop, Proceedings of the IEEE Conference on Computer Vision*
550 *and Pattern Recognition (CVPR)*, 2022. 21
- 551 [68] Magnus R Hestenes and Eduard Stiefel. Methods of conjugate gradients for solving. *Journal of*
552 *research of the National Bureau of Standards*, 49(6):409, 1952. 23, 24
- 553 [69] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual
554 networks. In *European conference on computer vision*, pages 630–645. Springer, 2016. 26

555 **Checklist**

- 556 1. For all authors...
- 557 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
558 contributions and scope? [Yes]
- 559 (b) Did you describe the limitations of your work? [Yes]
- 560 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 561 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
562 them? [Yes]
- 563 2. If you are including theoretical results...
- 564 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 565 (b) Did you include complete proofs of all theoretical results? [Yes] In the appendix.
- 566 3. If you ran experiments...
- 567 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
568 mental results (either in the supplemental material or as a URL)? [Yes] As a URL
- 569 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
570 were chosen)? [N/A]
- 571 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
572 ments multiple times)? [Yes] For the experiments comparing TCAV scores to our
573 concept importance score based on Sobol indices
- 574 (d) Did you include the total amount of compute and the type of resources used (e.g., type
575 of GPUs, internal cluster, or cloud provider)? [N/A]
- 576 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 577 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 578 (b) Did you mention the license of the assets? [N/A]
- 579 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 580
- 581 (d) Did you discuss whether and how consent was obtained from people whose data you're
582 using/curating? [N/A]
- 583 (e) Did you discuss whether the data you are using/curating contains personally identifiable
584 information or offensive content? [No] The ILSVRC2012 dataset contain personally
585 identifiable information [65]
- 586 5. If you used crowdsourcing or conducted research with human subjects...
- 587 (a) Did you include the full text of instructions given to participants and screenshots, if
588 applicable? [Yes] Screenshot of the experiments are in the appendix
- 589 (b) Did you describe any potential participant risks, with links to Institutional Review
590 Board (IRB) approvals, if applicable? [N/A]
- 591 (c) Did you include the estimated hourly wage paid to participants and the total amount
592 spent on participant compensation? [Yes]