

- [56] Bin Ren, Laurent Pueyo, Christine Chen, Élodie Choquet, John H Debes, Gaspard Duchêne, François Ménard, and Marshall D Perrin. Using data imputation for signal separation in high-contrast imaging. *The Astrophysical Journal*, 892(2):74, 2020. 4
- [57] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016. 3, 7, 19
- [58] Andrea Saltelli, Paola Annoni, Ivano Azzini, Francesca Campolongo, Marco Ratto, and Stefano Tarantola. Variance based sensitivity analysis of model output. design and estimator for the total sensitivity index. *Computer physics communications*, 181(2):259–270, 2010. 5
- [59] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 1, 3, 6, 7
- [60] Hua Shen and Ting-Hao Huang. How useful are the machine-generated interpretations to general users? a human evaluation on guessing the incorrectly predicted labels. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 8, pages 168–172, 2020. 2, 3
- [61] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017. 1, 7
- [62] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *In Workshop at International Conference on Learning Representations*. Citeseer, 2014. 1, 3, 6, 7
- [63] Leon Sixt, Maximilian Granz, and Tim Landgraf. When explanations lie: Why many modified bp attributions fail. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020. 2
- [64] Leon Sixt, Martin Schuessler, Oana-Iuliana Popescu, Philipp Weiß, and Tim Landgraf. Do users benefit from interpretable vision? a user study, baseline, and dataset. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2022. 2
- [65] Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 180–186, 2020. 2, 3
- [66] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017. 1, 3, 5, 6, 7
- [67] Ilya M Sobol. Sensitivity analysis for non-linear mathematical models. *Mathematical modelling and computational experiment*, 1:407–414, 1993. 2, 5, 20
- [68] Ilya M Sobol. Global sensitivity indices for nonlinear mathematical models and their monte carlo estimates. *Mathematics and computers in simulation*, 55(1-3):271–280, 2001. 2, 5
- [69] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014. 1, 3, 6
- [70] Pascal Sturmfels, Scott Lundberg, and Su-In Lee. Visualizing the impact of feature attribution baselines. *Distill*, 2020. 3, 13
- [71] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017. 1, 3, 5, 6, 7
- [72] Stefano Tarantola, Debora Gatelli, and Thierry Alex Mara. Random balance designs for the estimation of first order global sensitivity indices. *Reliability Engineering & System Safety*, 91(6):717–727, 2006. 5
- [73] Stephen A Vavasis. On the complexity of nonnegative matrix factorization. *SIAM Journal on Optimization*, 20(3):1364–1377, 2010. 6
- [74] Yu-Xiong Wang and Yu-Jin Zhang. Nonnegative matrix factorization: A comprehensive review. *IEEE Transactions on Knowledge and Data Engineering*, 25(6):1336–1353, 2013. 5
- [75] Peter Wilf, Shengping Zhang, Sharat Chikkerur, Stefan A Little, Scott L Wing, and Thomas Serre. Computer vision cracks the leaf code. *Proceedings of the National Academy of Sciences*, 113(12):3305–3310, 2016. 7
- [76] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014. 1, 3, 7
- [77] Ruihan Zhang, Prashan Madumal, Tim Miller, Krista A Ehinger, and Benjamin IP Rubinstein. Invertible concept-based explanations for cnn models with non-negative concept activation vectors. *arXiv preprint arXiv:2006.15417*, 2020. 2, 3, 5, 8, 23
- [78] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016. 3, 6

A Limitations and broader impact	13	C.2. Implicit differentiation	18
A.1. Limitations	13	D Sobol indices for concepts	19
A.2. Broader impact	13	E Human experiments	21
B More results of CRAFT	13	E.1. Utility evaluation	21
B.1. Qualitative comparison with ACE	13	E.2. Validation of Recursivity	21
B.2. Most important concepts.	13	F. Fidelity experiments	23
B.3. Feature Visualization validation	16	G Sanity Check	24
C Backpropagating through the NMF block	18		
C.1. Alternating Direction Method of Multipliers (ADMM) for NMF	18		

A. Limitations and broader impact

A.1. Limitations

Although we believe concept-based XAI to be a promising research direction, it isn't without pitfalls. It is capable of producing explanations that are ideally easy to understand by humans, but to what extent is a question that remains unanswered. The fact that there is no way to mathematically measure this prevents researchers from easily comparing the different techniques in the literature other than through time consuming and expensive experiments with human subjects. We think that developing a metric should be one of the field's priorities.

With CRAFT, we address the question of *what* by showing a cluster of the images that better represent each concept. However, we recognize that it's not perfect: in some cases, concepts are difficult to clearly define – put a label on what it represents –, and might induce some confirmation and selection bias. Feature visualization [51] might help in better illustrating the specific concept (as done in appendix B.3), but we believe there's still space for improvement. For instance, an interesting idea could be to leverage image captioning methods to describe the clusters of image crops, as textual information could help humans in better understanding clusters.

Although we believe CRAFT to be a considerable step in the good direction for the field of concept-based XAI, it also have some pitfalls. Namely, we chose the NMF as the activation factorization, which, while drastically improving the quality of extracted concepts, also comes with it's own caveats. For instance, it is known to be NP-hard to compute exactly, and in order to make it scalable, we had to use a tractable approximation by alternating the optimization of \mathbf{U} and \mathbf{W} through ADMM [5]. This approach might indeed yield non-unique solutions. Our experiments (section 4.3), have shown a low variance on between the runs, which comforts us about the stability of our results. However the absence of formal guarantee for uniqueness must be kept in mind: this subject is still an active topic of research and improvement could be expected in the near future. Namely, sparsity constraints and regularization seem to be promising paths. Naturally, we also need enough samples of the class under study to be available for the factorization to construct a relevant concept bank, which might affect the quality of the explanations on frugal applications where data is very

scarce.

A.2. Broader impact

We do hope that CRAFT helps in the transition to more human-understandable ways of explaining neural network models. It's capacity to find easily understandable concepts inside complex architectures and providing an indication of *where* they are located in the image is – to the best of our knowledge – unmatched. We also think that this method's structure is a step towards reducing confirmation bias: for instance dataset's labels are never used in this method, only the model's predictions. Without claiming to remove confirmation bias, the method focuses on what *the model sees* rather than what *we expect the model to see*. We believe this can help end-users build trust on computer vision models, and at the same time, provide ML practitioners with insights into potential sources of bias in the dataset (e.g. the ski pants in the astronaut/shovel example). Other methods in the literature obtaining similar results require very specific architectures [6] or to train another model to generate the explanations [21], so CRAFT provides a considerable advantage in the matter of flexibility in comparison.

B. More results of CRAFT

B.1. Qualitative comparison with ACE

Figure S1 compares the examples of concepts found by CRAFT against those found by ACE [24] for 3 classes of Imagenette. For each class the concepts are ordered by importance (the highest being the most important). ACE uses a clustering technique and TCAV to estimate importance, while CRAFT uses the method introduced in 3 and Sobol to estimate importance. These examples illustrate one of the weaknesses of ACE: the segmentation used can introduce biases through the baseline value used [19, 70]. The concepts found by CRAFT seem distinct: (vault, cross, stained glass) for the Church class, (dumpster, truck door, two-wheeler) for the garbage truck, and (eyes, nose, fluffy ears) for the English Springer.

B.2. Most important concepts.

We show more example of the 4 most important concepts for 6 classes: 'Chain saw', 'English springer', 'Gas pump', 'Golf ball', 'French horn' and 'Garbage Truck' (Figure S2).

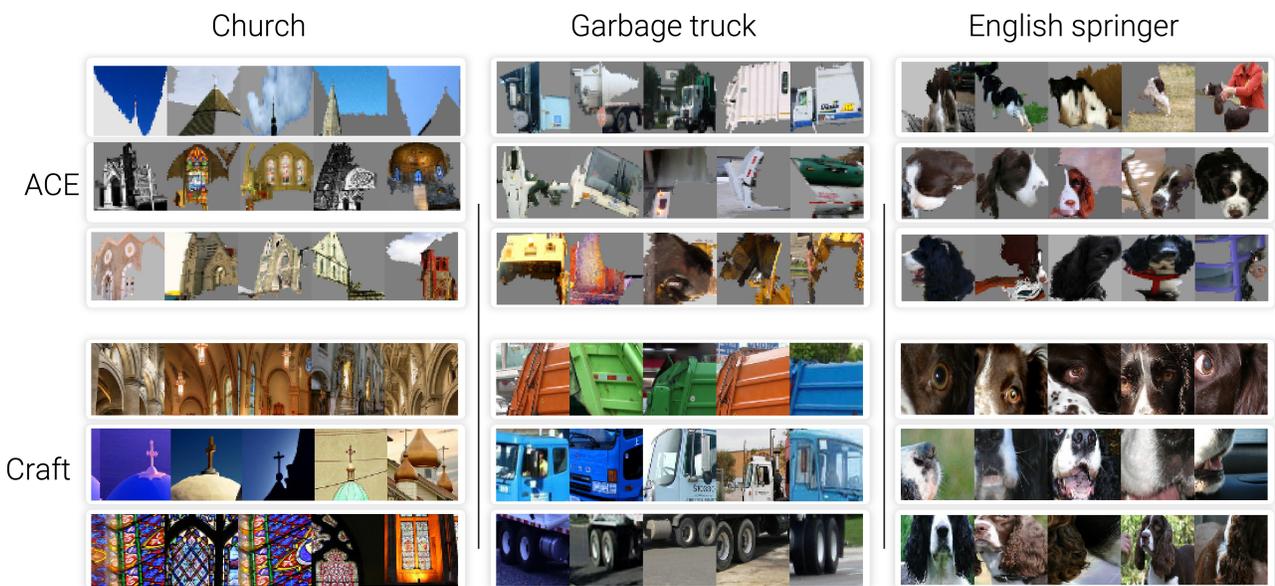


Figure S1. **Qualitative comparison.** We compare concepts found by our method (top) to those extracted with ACE [24] (bottom) for the classes *Church*, *Garbage truck* and *English springer* from ILSVRC2012 [10].

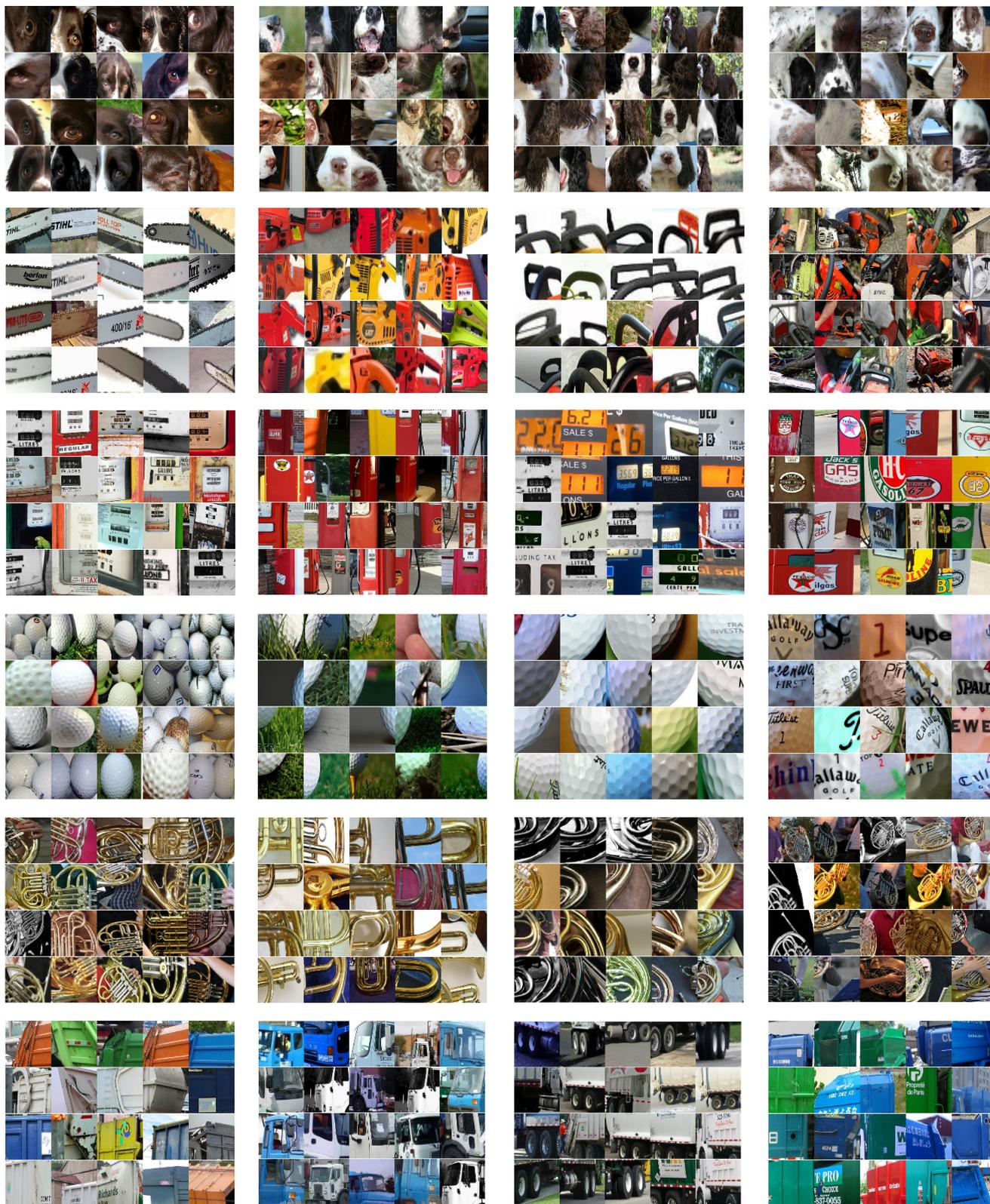


Figure S2. **CRAFT most important concepts.** The 4 most important concepts ranked by importance (left to right) for the following classes: 'English springer', 'Chain saw', 'Gas pump', 'Golf ball', 'French horn', and 'Garbage truck'.

B.3. Feature Visualization validation

Another way of interpreting concepts – as per [40] – is to employ feature visualization methods: through optimization, find an image that maximizes an activation pattern. In our case, we used the set of regularization and constraints proposed by [51], which allow us to successfully obtain realistic images. In Figures [S3-S8], we showcase these synthetic images obtained through feature visualization, along with the segments that maximize the target concept. We observe that they do reflect the underlying concepts of interest.

Concretely, to produce those feature visualization, we are looking for an image x^* that is optimized to correspond to a concept from the concept bank \mathbf{W}_i . We use the so called ‘dot-cosim’ loss proposed by [51], which give the following objective:

$$x^* = \arg \max_{x \in \mathcal{X}} \langle h_l(x), \mathbf{W}_i \rangle \frac{\langle h_l(x), \mathbf{W}_i \rangle^2}{\|h_l(x)\| \|\mathbf{W}_i\|} - \mathcal{R}(x)$$

With $\mathcal{R}(\cdot)$, the regularizations applied to x – the default regularizations in the **Xplique** library [16]. As for the specific parameters, we used Fourier preconditioning on the image with a decay rate of 0.8 and an Adam optimizer ($lr = 1e - 1$).

here



Figure S3. Feature visualization for chainsaw CRAFT concepts.



Figure S4. Feature visualization for Church CRAFT concepts.

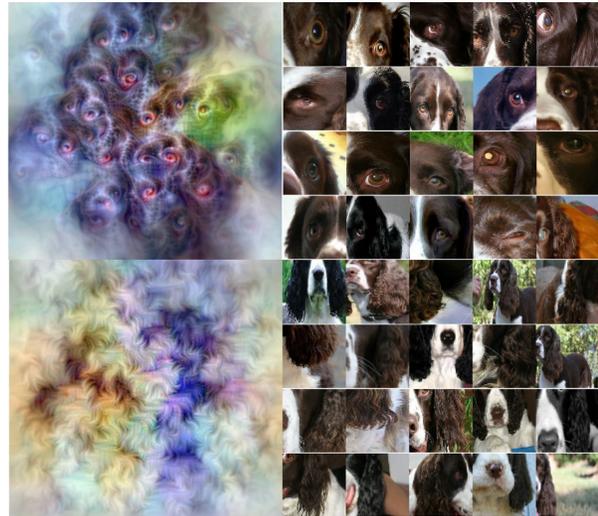


Figure S5. Feature visualization for english springer CRAFT concepts.

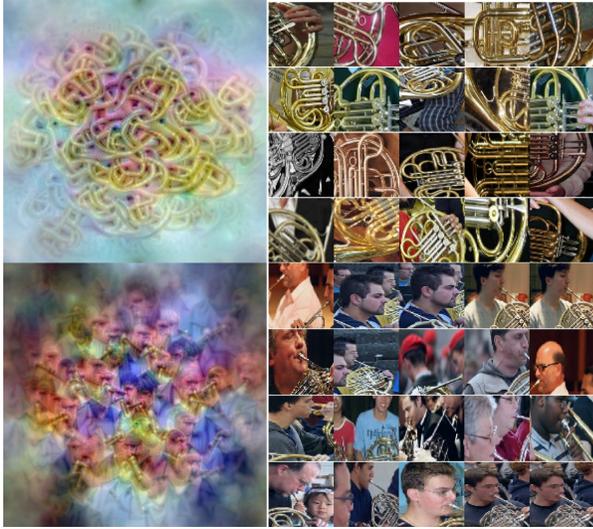


Figure S6. Feature visualization for french horn CRAFT concepts.



Figure S8. Feature visualization for golf CRAFT concepts.

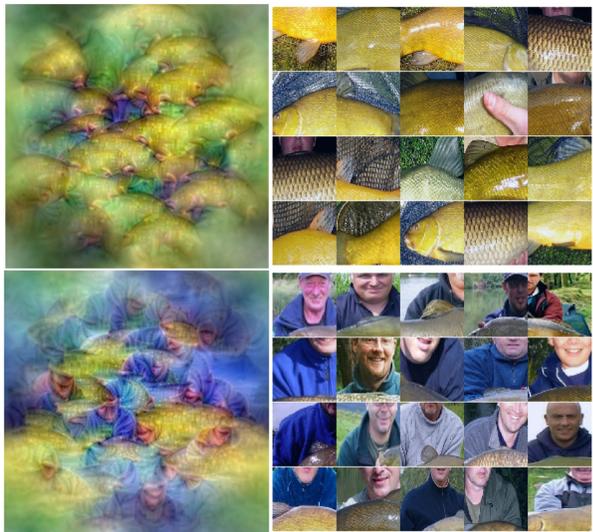


Figure S7. Feature visualization for tench CRAFT concepts.

C. Backpropagating through the NMF block

C.1. Alternating Direction Method of Multipliers (ADMM) for NMF

We recall that NMF decomposes the positive features vector $\mathbf{A} \in \mathbb{R}^{n \times p}$ of n examples lying in dimension p , into a product of positive low rank matrices $\mathbf{U}(\mathbf{A}) \in \mathbb{R}^{n \times r}$ and $\mathbf{W}(\mathbf{A}) \in \mathbb{R}^{p \times r}$ (with $r \ll \min(n, p)$), i.e the solution to the problem:

$$\min_{\mathbf{U} \geq 0, \mathbf{W} \geq 0} \frac{1}{2} \|\mathbf{A} - \mathbf{U}\mathbf{W}^T\|_F^2. \quad (4)$$

For simplicity we used a non-regularized version of the NMF objective, following Algorithms 1 and 3 in paper [32], based on ADMM [5]. This algorithm transforms the non-linear equality constraints into indicator functions δ . Auxiliary variables $\tilde{\mathbf{U}}, \tilde{\mathbf{W}}$ are also introduced to separate the optimization of the objective on the one side, and the satisfaction of the constraint on \mathbf{U}, \mathbf{W} on the other side. The equality constraints $\tilde{\mathbf{U}} = \mathbf{U}, \tilde{\mathbf{W}} = \mathbf{W}$ are linear and easily handled by the ADMM framework through the associated dual variables $\bar{\mathbf{U}}, \bar{\mathbf{W}}$. In our case, the problem in Equation 4 is transformed into:

$$\begin{aligned} \min_{\mathbf{U}, \tilde{\mathbf{U}}, \mathbf{W}, \tilde{\mathbf{W}}} \quad & \frac{1}{2} \|\mathbf{A} - \tilde{\mathbf{U}}\tilde{\mathbf{W}}^T\|_F^2 + \delta(\mathbf{U}) + \delta(\mathbf{W}), \\ \text{s.t.} \quad & \tilde{\mathbf{U}} = \mathbf{U}, \tilde{\mathbf{W}} = \mathbf{W} \\ & \text{with } \delta(\mathbf{H}) = \begin{cases} 0 & \text{if } \mathbf{H} \geq 0, \\ +\infty & \text{otherwise.} \end{cases} \end{aligned} \quad (5)$$

Note that $\tilde{\mathbf{U}}$ and \mathbf{U} (resp. $\tilde{\mathbf{W}}$ and \mathbf{W}) seem redundant: they are meant to be equal thanks to constraints $\tilde{\mathbf{U}} = \mathbf{U}, \tilde{\mathbf{W}} = \mathbf{W}$. This is standard practice within ADMM framework: introducing redundancies allows to disentangle the (unconstrained) optimization of the objective on one side (with $\tilde{\mathbf{U}}$ and $\tilde{\mathbf{W}}$) and constraint satisfaction on the other side with \mathbf{U} and \mathbf{W} . During the optimization process the variables $\tilde{\mathbf{U}}, \mathbf{U}$ (resp. $\tilde{\mathbf{W}}, \mathbf{W}$) are different, and only become equal in the limit at convergence. The dual variables $\bar{\mathbf{U}}, \bar{\mathbf{W}}$ control the balance between optimization of the objective $\frac{1}{2} \|\mathbf{A} - \tilde{\mathbf{U}}\tilde{\mathbf{W}}^T\|_F^2$ and constraint satisfaction $\tilde{\mathbf{U}} = \mathbf{U}, \tilde{\mathbf{W}} = \mathbf{W}$. The constraints are simplified at the cost of a non-smooth (and even a non-finite) objective function $\frac{1}{2} \|\mathbf{A} - \tilde{\mathbf{U}}\tilde{\mathbf{W}}^T\|_F^2 + \delta(\mathbf{U}) + \delta(\mathbf{W})$ due to the term $\delta(\mathbf{U}) + \delta(\mathbf{W})$. ADMM proceeds to create a so-called *augmented Lagrangian* with l_2 regularization $\rho > 0$:

$$\begin{aligned} \mathcal{L}(\mathbf{A}, \mathbf{U}, \mathbf{W}, \tilde{\mathbf{U}}, \tilde{\mathbf{W}}, \bar{\mathbf{U}}, \bar{\mathbf{W}}) = & \\ & \frac{1}{2} \|\mathbf{A} - \tilde{\mathbf{U}}\tilde{\mathbf{W}}^T\|_F^2 + \delta(\mathbf{U}) + \delta(\mathbf{W}) \\ & + \bar{\mathbf{U}}^T (\tilde{\mathbf{U}} - \mathbf{U}) + \bar{\mathbf{W}}^T (\tilde{\mathbf{W}} - \mathbf{W}) \\ & + \frac{\rho}{2} (\|\tilde{\mathbf{U}} - \mathbf{U}\|_2^2 + \|\tilde{\mathbf{W}} - \mathbf{W}\|_2^2). \end{aligned} \quad (6)$$

This regularization ensures that the dual problem is well posed and that it remain convex, even with the non smooth and infinite terms $\delta(\mathbf{U}) + \delta(\mathbf{W})$. Once again, this is standard practice within ADMM framework. The (regularized) problem associated to this Lagrangian is decomposed into a sequence of convex problems that alternate minimization over the $\mathbf{U}, \tilde{\mathbf{U}}, \bar{\mathbf{U}}$ and the $\mathbf{W}, \tilde{\mathbf{W}}, \bar{\mathbf{W}}$ triplets.

$$\mathbf{U}_{t+1} = \arg \min_{\mathbf{U}=\tilde{\mathbf{U}}} \frac{1}{2} \|\mathbf{A} - \tilde{\mathbf{U}}\mathbf{W}_t^T\|_F^2 + \delta(\mathbf{U}) + \frac{\rho}{2} \|\tilde{\mathbf{U}} - \mathbf{U}\|_2^2. \quad (7)$$

$$\mathbf{W}_{t+1} = \arg \min_{\mathbf{W}=\tilde{\mathbf{W}}} \frac{1}{2} \|\mathbf{A} - \mathbf{U}_t\tilde{\mathbf{W}}^T\|_F^2 + \delta(\mathbf{W}) + \frac{\rho}{2} \|\tilde{\mathbf{W}} - \mathbf{W}\|_2^2. \quad (8)$$

This guarantees a monotonic decrease of the objective function $\|\mathbf{A} - \tilde{\mathbf{U}}_t\tilde{\mathbf{W}}_t^T\|_F^2$. Each of these sub-problems is thus solved with ADMM separately, by alternating minimization steps of $\frac{1}{2} \|\mathbf{A} - \tilde{\mathbf{U}}\mathbf{W}_t^T\|_F^2 + \bar{\mathbf{U}}^T (\tilde{\mathbf{U}} - \mathbf{U}) + \frac{\rho}{2} \|\mathbf{U} - \tilde{\mathbf{U}}\|_2^2$ over $\tilde{\mathbf{U}}$ (*i*), with minimization steps of $\delta(\mathbf{U}) + \frac{\rho}{2} \|\mathbf{U} - \tilde{\mathbf{U}}\|_2^2$ over \mathbf{U} (*ii*), and gradient ascent steps (*iii*) on the dual variable $\bar{\mathbf{U}} \leftarrow \bar{\mathbf{U}} + (\tilde{\mathbf{U}} - \mathbf{U})$. A similar scheme is used for \mathbf{W} updates. Step (*i*) is a simple convex quadratic program with equality constraints, whose KKT [38, 45] conditions yield a linear system with a Positive Semi-Definite (PSD) matrix. Step (*ii*) is a simple projection of $\tilde{\mathbf{U}}$ onto the convex set $\delta^{-1}(\mathbf{0})$. Finally, step (*iii*) is inexpensive.

Concretely, we solved the quadratic program using Conjugate Gradient [30], from *jax.scipy.sparse.linalg.cg*. This indirect method only involves *matrix-vector* products and can be more GPU-efficient than methods that are based on matrix factorization (such as Cholesky decomposition). Also, we re-implemented the pseudo code of [32] in *Jax* for a fully GPU-compatible program. We used the primal variables $\mathbf{U}_0, \mathbf{W}_0$ returned by *sklearn.decompose.nmf* as a *warm start* for ADMM and observe that the high quality initialization of these primal variables considerably speeds up the convergence of the dual variables.

C.2. Implicit differentiation

The Lagrangian of the NMF problem reads $\mathcal{L}(\mathbf{U}, \mathbf{W}, \bar{\mathbf{U}}, \bar{\mathbf{W}}) = \frac{1}{2} \|\mathbf{A} - \mathbf{U}\mathbf{W}^T\|_F^2 - \bar{\mathbf{U}}^T \mathbf{U} - \bar{\mathbf{W}}^T \mathbf{W}$, with dual variables $\bar{\mathbf{U}}$ and $\bar{\mathbf{W}}$ associated to the constraints $\mathbf{U} \geq 0, \mathbf{W} \geq 0$. It yields a function \mathcal{F} based on the KKT conditions [38, 45] whose optimal tuple $\mathbf{U}, \mathbf{W}, \bar{\mathbf{U}}, \bar{\mathbf{W}}$ is a root.

For single NNLS problem (for example, with optimization over \mathbf{U}) the KKT conditions are:

$$\begin{cases} \nabla_{\mathbf{U}} \left(\frac{1}{2} \|\mathbf{A} - \tilde{\mathbf{U}}\tilde{\mathbf{W}}^T\|_F^2 + \bar{\mathbf{U}}^T(-\mathbf{U}) \right) = 0, \text{ stationarity,} \\ -\mathbf{U} \leq 0, \text{ primal feasibility,} \\ \bar{\mathbf{U}} \odot \mathbf{U} = 0, \text{ complementary slackness,} \\ \bar{\mathbf{U}} \geq 0, \text{ dual feasibility.} \end{cases} \quad (9)$$

By stacking the KKT conditions of the NNLS problems the we obtain the so-called *optimality function* \mathbf{F} :

$$\mathbf{F}((\mathbf{U}, \mathbf{W}, \bar{\mathbf{U}}, \bar{\mathbf{W}}), \mathbf{A}) = \begin{cases} (\mathbf{U}\mathbf{W}^T - \mathbf{A})\mathbf{W} - \bar{\mathbf{U}}, \\ (\mathbf{W}\mathbf{U}^T - \mathbf{A}^T)\mathbf{U} - \bar{\mathbf{W}}, \\ \bar{\mathbf{U}} \odot \mathbf{U}, \\ \bar{\mathbf{W}} \odot \mathbf{W}. \end{cases} \quad (10)$$

The implicit function theorem [25] allows us to use implicit differentiation [3, 25, 44] to efficiently compute the Jacobians $\frac{\partial \mathbf{U}}{\partial \mathbf{A}}$ and $\frac{\partial \mathbf{W}}{\partial \mathbf{A}}$ without requiring to back-propagate through each of the iterations of the NMF solver:

$$\frac{\partial(\mathbf{U}, \mathbf{W}, \bar{\mathbf{U}}, \bar{\mathbf{W}})}{\partial \mathbf{A}} = -(\partial_1 \mathbf{F})^{-1} \partial_2 \mathbf{F}. \quad (11)$$

Implicit differentiation requires access to the dual variables of the optimization problem in equation 1, which are not computed by Scikit-learn’s popular implementation. Scikit-learn uses Block coordinate descent algorithm [7, 17], with a randomized SVD initialization. Consequently, we leverage our implementation in Jax based on ADMM [5].

Concretely, we perform a two-stage backpropagation *Jax* (2) \rightarrow *Tensorflow* (1) to leverage the advantage of each framework. The lower stage (1) corresponds to feature extraction $\mathbf{A} = \mathbf{h}_l(\mathbf{X})$ from crops of images \mathbf{X} , and upper stage (2) computes NMF $\mathbf{A} \approx \mathbf{U}\mathbf{W}^T$.

We use the *Jaxopt* [4] library that allows efficient computation of $\frac{\partial(\mathbf{U}, \mathbf{W}, \bar{\mathbf{U}}, \bar{\mathbf{W}})}{\partial \mathbf{A}} = -(\partial_1 \mathbf{F})^{-1} \partial_2 \mathbf{F}$. The matrix $(\partial_1 \mathbf{F})^{-1}$ is never explicitly computed – that would be too costly. Instead, the system $\partial_1 \mathbf{F} \frac{\partial(\mathbf{U}, \mathbf{W}, \bar{\mathbf{U}}, \bar{\mathbf{W}})}{\partial \mathbf{A}} = -\partial_2 \mathbf{F}$ is solved with Conjugate Gradient [30] through the use of Jacobian Vector Products (JVP) $\mathbf{v} \mapsto (\partial_1 \mathbf{F})\mathbf{v}$.

The chain rule yields:

$$\frac{\partial \mathbf{U}}{\partial \mathbf{X}} = \frac{\partial \mathbf{A}}{\partial \mathbf{X}} \frac{\partial \mathbf{U}}{\partial \mathbf{A}}.$$

Usually, most Autodiff frameworks (e.g Tensorflow, Pytorch, Jax) handle it automatically. Unfortunately, combining two of those framework raises a new difficulty since they are not compatible. Hence, we re-implement manually the two stages auto-differentiation.

Since r is far smaller ($r = 25$ in all our experiments) than input dimension \mathbf{X} (typically 224×244 for ImageNet

images), back-propagation is the preferred algorithm in this setting over forward-propagation. We start by computing sequentially the gradients $\nabla_{\mathbf{X}} \mathbf{U}_i$ for all concepts $1 \leq i \leq r$. This amounts to compute $\mathbf{v} = \nabla_{\mathbf{A}} \mathbf{U}_i$ with Implicit Differentiation in Jax, convert the Jax array \mathbf{v} into Tensorflow tensor, and then to compute $\nabla_{\mathbf{X}} \mathbf{U}_i = \frac{\partial \mathbf{A}}{\partial \mathbf{X}} \nabla_{\mathbf{A}} \mathbf{U}_i = \nabla_{\mathbf{X}} (\mathbf{h}_l(\mathbf{X}) \cdot \mathbf{v})$. The latter is easily done in Tensorflow. Finally we stack the gradients $\nabla_{\mathbf{X}} \mathbf{U}_i$ to obtain the Jacobian $\frac{\partial \mathbf{U}}{\partial \mathbf{X}}$.

D. Sobol indices for concepts

We propose to formally derive the Sobol indices for the estimation of the importance of concepts. Let us define a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ of possible concept perturbations. In order to build these concept perturbations, we start from an original vector of concepts coefficient $\hat{\mathbf{U}} \in \mathbb{R}^r$ and use i.i.d. stochastic masks $\mathbf{M} = (M_1, \dots, M_r) \in \mathcal{M} \subseteq [0, 1]^r$, as well as a perturbation operator $\pi : \mathcal{A} \times \mathcal{M} \rightarrow \mathcal{A}$ to create stochastic perturbation of $\hat{\mathbf{U}}$ that we call concept perturbation $\mathbf{U} = \pi(\hat{\mathbf{U}}, \mathbf{M})$.

Concretely, to create our concept perturbation we consider the inpainting function as our perturbation operator (as in [14, 55, 57]): $\pi(\hat{\mathbf{U}}, \mathbf{M}) = \hat{\mathbf{U}} \odot \mathbf{M} + (\mathbf{1} - \mathbf{M})\mu$ with \odot the Hadamard product and $\mu \in \mathbb{R}$ a baseline value, here zero. For the sake of notation, we will note $\mathbf{f} : \mathcal{A} \rightarrow \mathbb{R}$ the function mapping a random concept perturbation \mathbf{U} from the layer l to the output. We denote the set $\mathcal{U} = \{1, \dots, r\}$, \mathbf{u} a subset of \mathcal{U} , its complementary $\sim \mathbf{u}$ and $\mathbb{E}(\cdot)$ the expectation over the perturbation space. Finally, we assume that $\mathbf{f} \in \mathbb{L}^2(\mathcal{A}, \mathbb{P})$ i.e. $\mathbb{E}(\mathbf{f}(\mathbf{U})) < +\infty$.

The Hoeffding decomposition allows us to express the function \mathbf{f} into summands of increasing dimension, denoting $\mathbf{f}_{\mathbf{u}}$ the partial contribution of the concepts $\mathbf{U}_{\mathbf{u}} = (U_i)_{i \in \mathbf{u}}$ to the score $\mathbf{f}(\mathbf{U})$:

$$\begin{aligned} \mathbf{f}(\mathbf{U}) &= \mathbf{f}_{\emptyset} \\ &+ \sum_i^r \mathbf{f}_i(U_i) \\ &+ \sum_{1 \leq i < j \leq r} \mathbf{f}_{i,j}(U_i, U_j) + \dots \\ &+ \mathbf{f}_{1, \dots, r}(U_1, \dots, U_r) \\ &= \sum_{\mathbf{u} \subseteq \mathcal{U}} \mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}}). \end{aligned} \quad (12)$$

Eq. 12 consists of 2^r terms and is unique under the following orthogonality constraint:

$$\forall (\mathbf{u}, \mathbf{v}) \subseteq \mathcal{U}^2 \text{ s.t. } \mathbf{u} \neq \mathbf{v}, \quad \mathbb{E}(\mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}}) \mathbf{f}_{\mathbf{v}}(\mathbf{U}_{\mathbf{v}})) = 0. \quad (13)$$

Furthermore, orthogonality yields the characterization $\mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}}) = \mathbb{E}(\mathbf{f}(\mathbf{U}) | \mathbf{U}_{\mathbf{u}}) - \sum_{\mathbf{v} \subsetneq \mathbf{u}} \mathbf{f}_{\mathbf{v}}(\mathbf{U}_{\mathbf{v}})$ and allows us

to decompose the model variance as:

$$\begin{aligned}
\mathbb{V}(\mathbf{f}(\mathbf{U})) &= \sum_i^r \mathbb{V}(\mathbf{f}_i(U_i)) \\
&+ \sum_{1 \leq i < j \leq r} \mathbb{V}(\mathbf{f}_{i,j}(U_i, U_j)) \\
&+ \dots + \mathbb{V}(\mathbf{f}_{1,\dots,r}(U_1, \dots, U_r)) \\
&= \sum_{\mathbf{u} \subseteq \mathcal{U}} \mathbb{V}(\mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}})).
\end{aligned} \tag{14}$$

Building from Eq. 14, it is natural to characterize the influence of any subset of concepts \mathbf{u} as its own variance w.r.t. the total variance. This yields, after normalization by $\mathbb{V}(\mathbf{f}(\mathbf{U}))$, the general definition of Sobol' indices.

Definition D.1 (Sobol indices [67]). The sensitivity index $\mathcal{S}_{\mathbf{u}}$ which measures the contribution of the concept set $\mathbf{U}_{\mathbf{u}}$ to the model response $\mathbf{f}(\mathbf{U})$ in terms of fluctuation is given

by:

$$\begin{aligned}
\mathcal{S}_{\mathbf{u}} &= \frac{\mathbb{V}(\mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}}))}{\mathbb{V}(\mathbf{f}(\mathbf{U}))} \\
&= \frac{\mathbb{V}(\mathbb{E}(\mathbf{f}(\mathbf{U})|\mathbf{U}_{\mathbf{u}})) - \sum_{\mathbf{v} \subset \mathbf{u}} \mathbb{V}(\mathbb{E}(\mathbf{f}(\mathbf{U})|\mathbf{U}_{\mathbf{v}}))}{\mathbb{V}(\mathbf{f}(\mathbf{U}))}.
\end{aligned} \tag{15}$$

Sobol indices give a quantification of the importance of any subset of concepts with respect to the model decision, in the form of a normalized measure of the model output deviation from $\mathbf{f}(\mathbf{U})$. Thus, Sobol indices sum to one : $\sum_{\mathbf{u} \subseteq \mathcal{U}} \mathcal{S}_{\mathbf{u}} = 1$.

Furthermore, the framework of Sobol' indices enables us to easily capture higher-order interactions between features. Thus, we can view the Total Sobol indices defined in 2 as the sum of of all the Sobol indices containing the concept i : $\mathcal{S}_{T_i} = \sum_{\mathbf{u} \subseteq \mathcal{U}, i \in \mathbf{u}} \mathcal{S}_{\mathbf{u}}$. Concretely, we estimate the total Sobol indices using the Jansen estimator [36] and Quasi-Monte carlo Sequence (Sobol LP_{τ} sequence).

E. Human experiments

We first describe how participants were enrolled in our studies, then the general experimental design they went through.

E.1. Utility evaluation

Participants The participants that went through our experiments are users from the online platform Amazon Mechanical Turk (AMT), specifically, we recruit users with high qualifications (number of HIT completed = 5000 and HIT accepted > 98%). All participants provided informed consent electronically in order to perform the experiment (~ 5 – 8 min), for which they received 1.4\$.

For the *Husky vs. Wolf* scenario, $n = 84$ participants passed all our screening and filtering process, respectively $n = 32$ for CRAFT, $n = 22$ for ACE and $n = 22$ for CRAFTCO.

For the *Leaves* scenario, after filtering, we analyzed data from $n = 87$ participants, respectively $n = 32$ for CRAFT, $n = 24$ for ACE and $n = 31$ for CRAFTCO.

For the *"Kit Fox" vs. "Red Fox"* scenario, the results come from $n = 79$ participants who passed all our screening processes, respectively $n = 22$ for CRAFT, $n = 31$ for ACE and $n = 26$ for CRAFTCO.

General study design We followed the experimental design proposed by Colin and Fel et al. [8], in which explanations are evaluated according to their ability to help training participants at getting better at predicting their models' decisions on unseen images.

Each of those participants are only tested on a single condition to avoid possible experimental confounds.

The main experiment is divided into 3 training sessions (with 5 training samples in each) each followed by a brief test. In each individual training trial, an image was presented with the associated prediction of the model, together with an explanation. After a brief training phase (5 samples), participants' ability to predict the classifier's output was evaluated on 7 new samples during a test phase. During the test phase, no explanation was provided. We also

use the reservoir that subjects can refer to during the testing phase to minimize memory load as a confounding factor.

We implement the same 3-stage screening process as in [8]: First we filter participants not successful at the practice session done prior to the main experiment used to teach them the task, then we have them go through a quiz to make sure they understood the instructions. Finally, we add a catch trial in each testing phase –that users paying attention are expected to be correct on– allowing us to catch uncooperative participants.

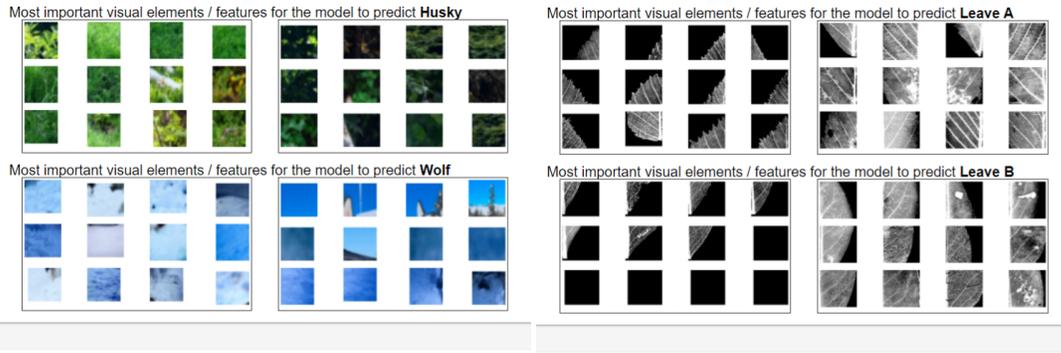
E.2. Validation of Recursivity

Participants Behavioral accuracy data were gathered from $n = 73$ participants. All participants provided informed consent electronically in order to perform the experiment (~ 4 – 6 min). The protocol was approved by the University IRB and was carried out in accordance with the provisions of the World Medical Association Declaration of Helsinki. For each of the 2 experiment tested, we had prepared filtering criteria for uncooperative people (namely based on time), but all participants passed these filters.

General study design For the first experiment – consisting in finding the intruder among elements of the same concept and an element from a different concept (but of the same class, see Figure S10b) – the order of presentation is randomized across participants so that it does not bias the results. Moreover, in order to avoid any bias coming from the participants themselves (one group being more successful than the other) all participants went through both conditions of finding intruders in batches of images coming from either concepts or sub-concepts. Concerning experiment 2, the order was also randomized (see Figure S10c).

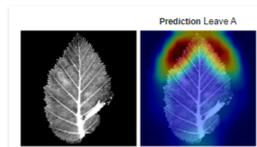
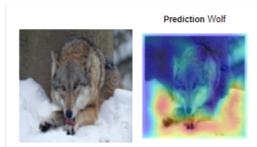
The participants had to successively find 30 intruders (15 block concepts and 15 block sub-concepts) for experiment 1 and then make 15 choices (sub-concept vs concept) for experiment 2, see Figure S10a.

The expert participants are people working in machine learning (researchers, software developers, engineers) and have participated in the study following an announcement in the authors' laboratory/company. The other participants (Laymen) have no expertise in machine learning.

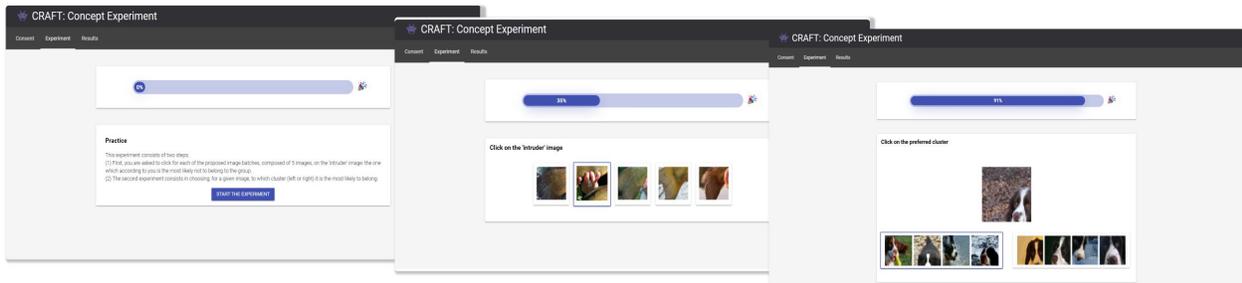


Study examples of classification of the model, try to understand its behavior.

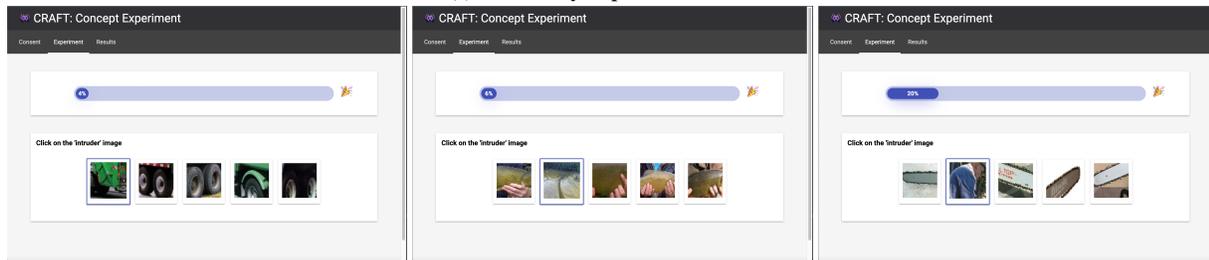
Study examples of classification of the model, try to understand its behavior.



(a) **Utility experiment.** Training trials taken from the *Husky vs. Wolf* scenario (left) and the *Leaves* scenario (right).



(a) **Recursivity Experiment Website.**



(b) **Binary choice experiment.**



(c) **Intruder experiment.**

F. Fidelity experiments

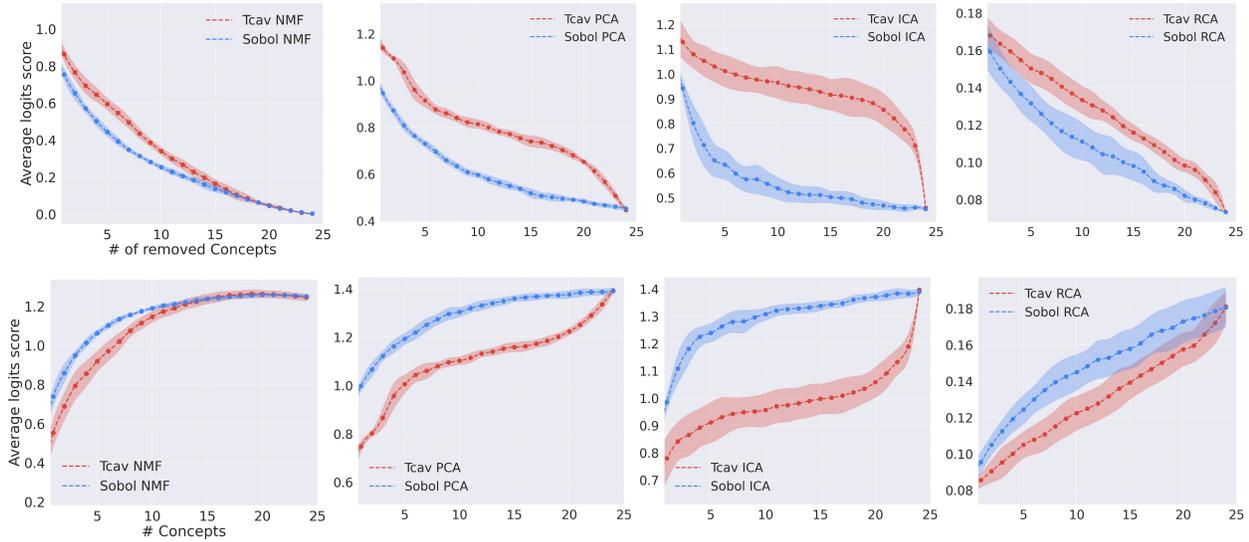


Figure S11. **(1) Deletion curves** for different concept extraction methods, Sobol outperforms TCAV not only for NMF to correctly estimate concept importance (lower is better). **(2) Insertion curves** for different concept extraction methods, Sobol outperforms TCAV to correctly estimate concept importance (higher is better).

For our experiments on the concept importance measure, we focused on certain classes of ILSRVC2012 [10] and used a ResNet50V2 [29] that had already been trained on this dataset. Just like in [24, 77], we measure the insertion and deletion metrics for our concept extraction technique – as well as concepts vectors extracted using PCA, ICA and RCA as dimensionality reduction algorithms, see Figure S11 – and we compare them when we add/remove the concepts as ranked by the TCAV score [40] and by the Sobol importance score. As originally explained in [55], the objective of these metrics is to add/remove parts of the input according to how much an explainability method considers that it is influential and looking at the speed at which the logit for the predicted class increases/decreases.

In particular, for our experimental evaluations, we have randomly chosen 100000 images from ILSVRC2012 [10] and computed the deletion and insertion metrics for 5 different seeds – for a total of half a million images. In Figure S11, the shade around the curves represent the standard deviation over these 5 experiments.

G. Sanity Check

Following the work from [1], we performed a sanity check on our method, by running the concept extraction pipeline on a randomized model. This procedure was performed on a ResNet-50v2 model with randomized weights. As showcased in Figure S12, the concepts drastically differ from trained models, thus proving that CRAFT passes the sanity check.

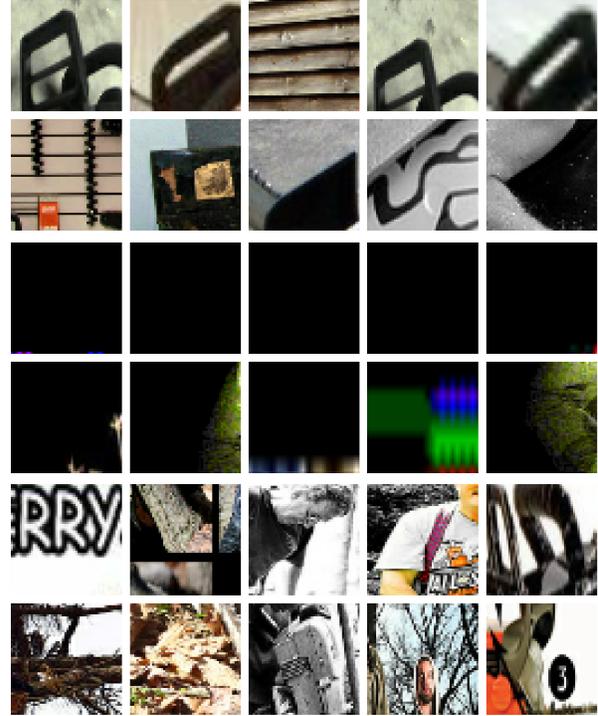


Figure S12. **Sanity check of the method:** we ran the method on a Resnet50 with randomized weights, and extracted the 3 most relevant concepts for the class 'Chain saw'. When weights are randomized, concepts are mainly based on color histograms.