

593	A	Limitations and broader impact	16
594	A.1	Limitations	16
595	A.2	Broader impact	16
596	B	More results of CRAFT	16
597	B.1	Concept Attribution Maps.	16
598	B.2	Most important concepts.	17
599	B.3	Feature Visualization validation	21
600	C	Backpropagating through the NMF block	22
601	C.1	Alternating Direction Method of Multipliers (ADMM) for NMF	22
602	C.2	Implicit differentiation	23
603	D	Sobol indices for concepts	24
604	E	Human experiments	25
605	F	Fidelity experiments	26
606	G	Additional examples of concepts and sub-concepts	28
607	G.1	Sanity Check	28
608	H	Computational cost	28

609 **A Limitations and broader impact**

610 **A.1 Limitations**

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

617 With CRAFT, we address the question of *what* by showing a cluster of the images that better represent
618 each concept. However, we recognize that it's not perfect: in some cases, concepts are difficult to
619 clearly define – put a label on what it represents –, and might induce some confirmation and selection
620 bias. Feature visualization [18] might help in better illustrating the specific concept (as done in
621 appendix B.3), but we believe there's still space for improvement. For instance, an interesting idea
622 could be to leverage image captioning methods to describe the clusters of image crops, as textual
623 information could help humans in better understanding clusters.

624 Although we believe CRAFT to be a considerable step in the good direction for the field of concept-
625 based XAI, it also have some pitfalls. Namely, we chose the NMF as the activation factorization,
626 which, while drastically improving the quality of extracted concepts, also comes with it's own caveats.
627 For instance, it is known to be NP-hard to compute exactly, and in order to make it scalable, we had
628 to use a tractable approximation by alternating the optimization of U and W through ADMM [63].
629 This approach might indeed yield non-unique solutions. Our experiments (section 4.4), have shown a
630 low variance on between the runs, which comforts us about the stability of our results. However the
631 absence of formal guarantee for uniqueness must be kept in mind: this subject is still an active topic
632 of research and improvement could be expected in the near future. Namely, sparsity constraints and
633 regularization seem to be promising paths. Naturally, we also need enough samples of the class under
634 study to be available for the factorization to construct a relevant concept bank, which might affect the
635 quality of the explanations on frugal applications where data is very scarce.

636 **A.2 Broader impact**

637 We do hope that CRAFT helps in the transition to more human-understandable ways of explain-
638 ing neural network models. It's capacity to find easily understandable concepts inside complex
639 architectures and providing an indication of *where* they are located in the image is – to the best of
640 our knowledge – unmatched. We also think that this method's structure is a step towards reducing
641 confirmation bias: for instance dataset's labels are never used in this method, only the model's
642 predictions. Without claiming to remove confirmation bias, the method focuses on what *the model*
643 *sees* rather than what *we expect the model to see*. We believe this can help end-users build trust on
644 computer vision models, and at the same time, provide ML practitioners with insights into potential
645 sources of bias in the dataset (e.g. the ski pants in the astronaut/shovel example). Other methods
646 in the literature obtaining similar results require very specific architectures [36] or to train another
647 model to generate the explanations [66], so CRAFT provides a considerable advantage in the matter
648 of flexibility in comparison.

649 **B More results of CRAFT**

650 **B.1 Concept Attribution Maps.**

651 We show more examples of Concept Attribution Maps for the classes 'Chain saw' in Figure S2 and
652 'Parachute' in Figure S1.

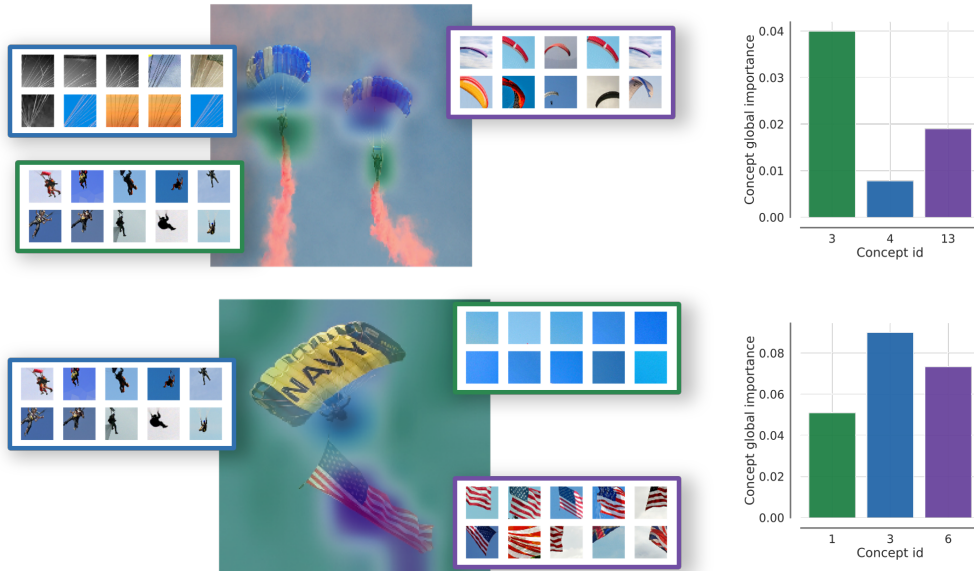


Figure S1: **CRAFT results for the class ‘Parachute’**. The model under study is a ResNet50, and we used the penultimate layer to apply the matrix activation factorization.

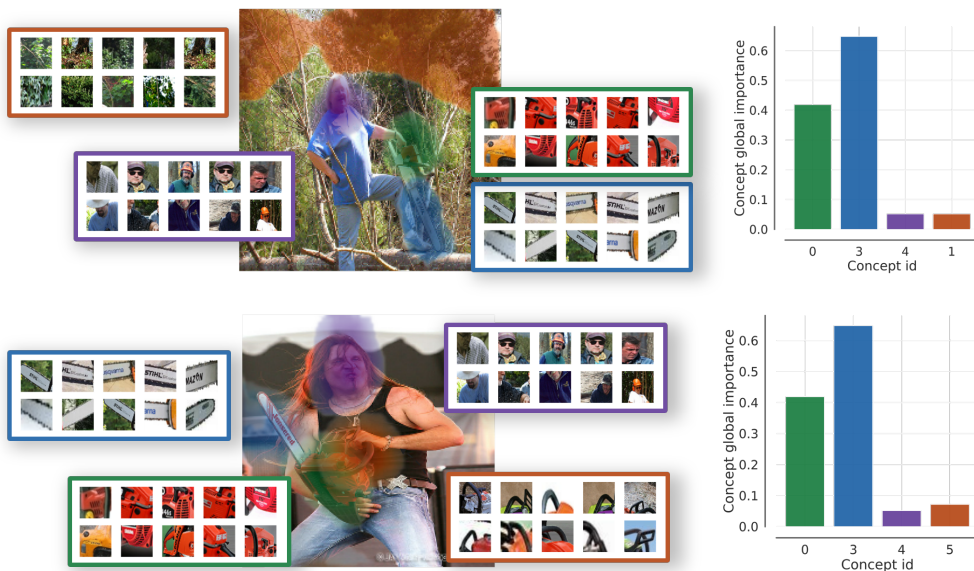


Figure S2: **CRAFT results for the class ‘Chain saw’**. The model under study is a ResNet50 we used the penultimate layer to apply the matrix activation factorization.

653 **B.2 Most important concepts.**

654 We show more example of the 4 most importants concepts for 6 classes: ‘Chain saw’ and ‘English
 655 springer’ (Figure S3), ‘Gas pump’ and ‘Golf ball’ (Figure S4), ‘French horn’ and ‘Garbage Truck’
 656 (Figure S5).

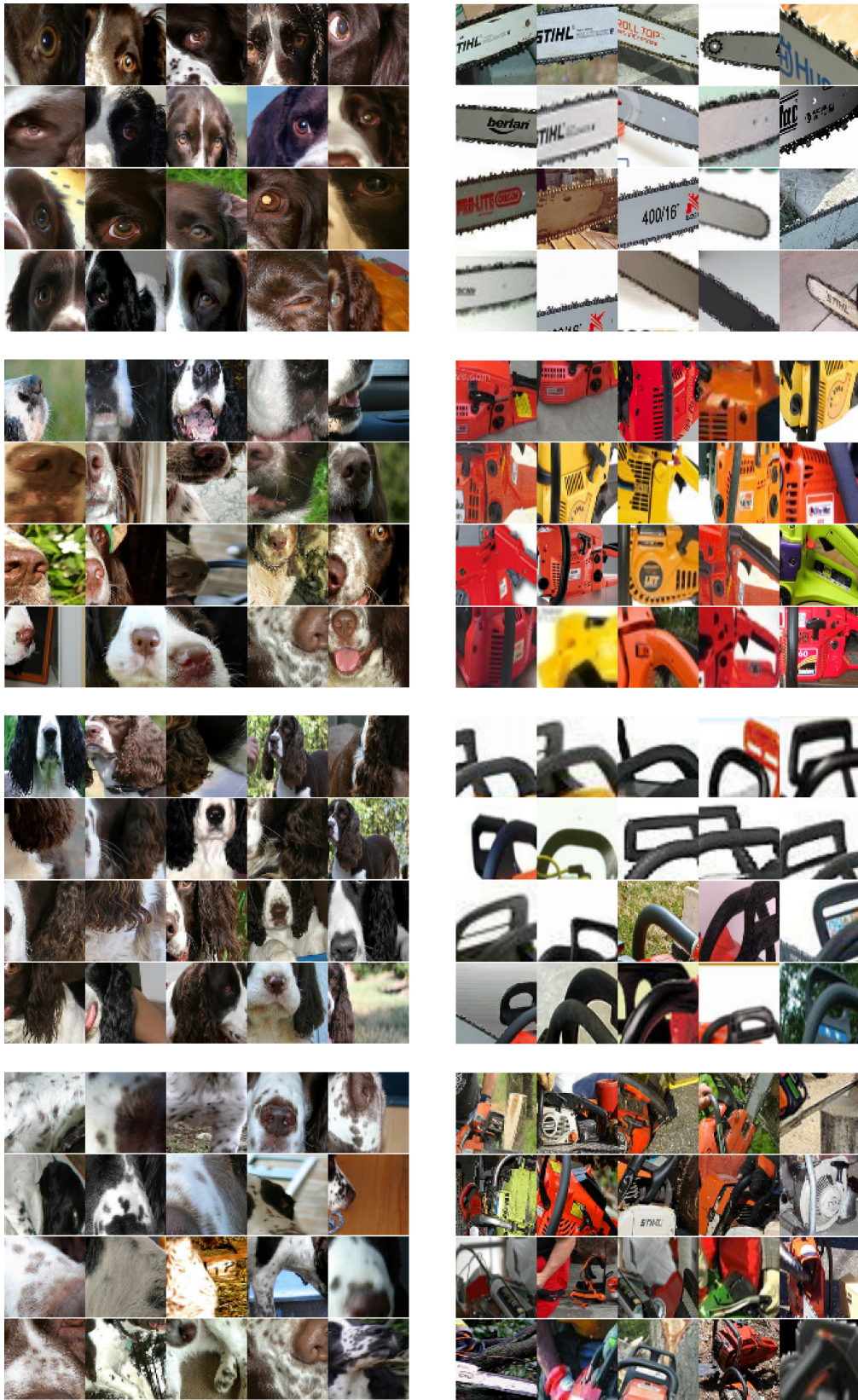


Figure S3: **CRAFT most important concepts.** The 4 most important concepts (higher means more important) for 'English springer' (left) and 'Chain saw' (right).

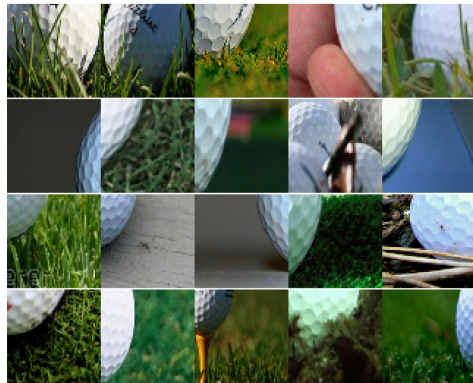


Figure S4: CRAFT most important concepts. The 4 most important concepts (higher means more important) for 'Gas pump' (left) and 'Golf ball' (right).

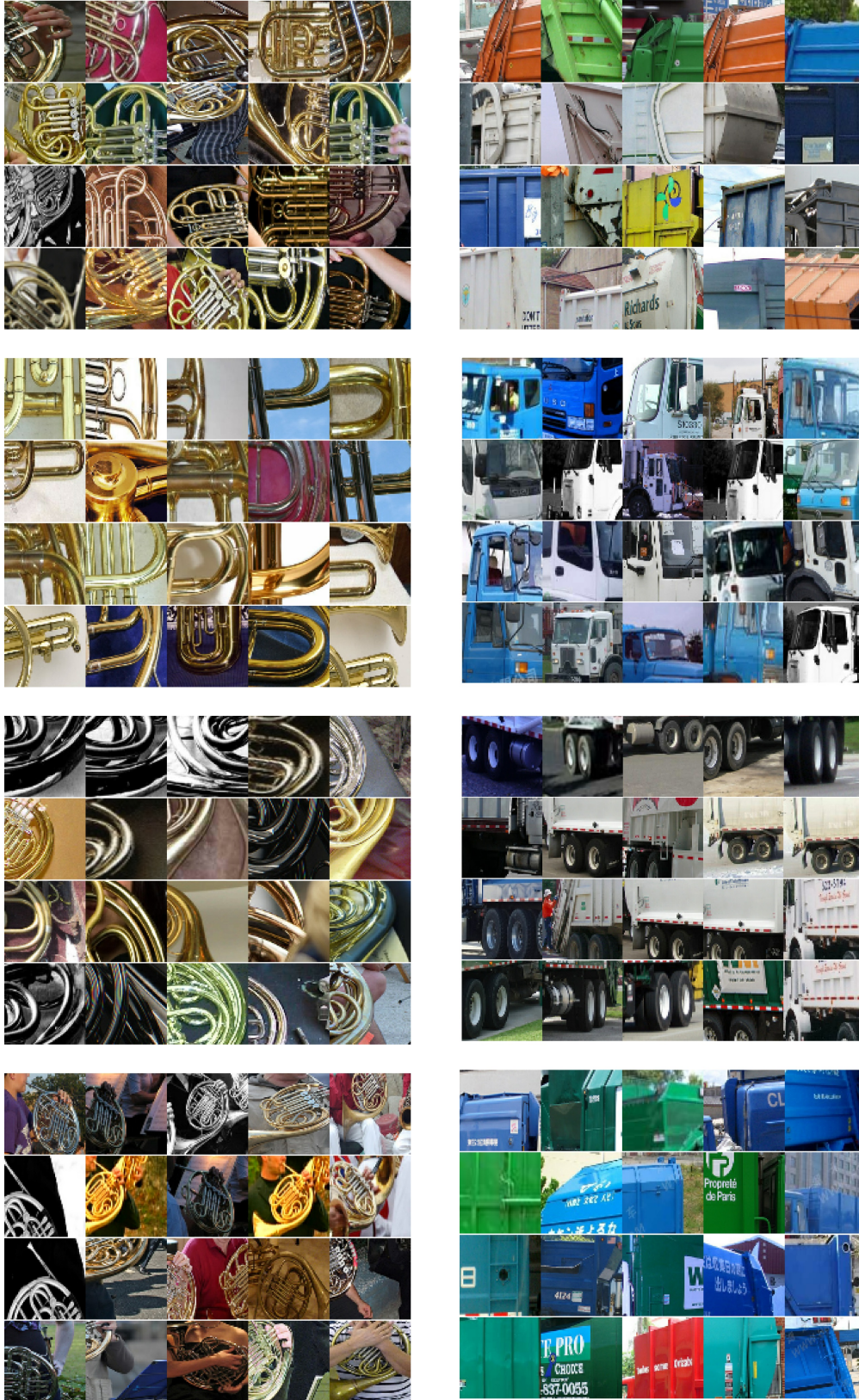


Figure S5: **CRAFT most important concepts.** The 4 most important concepts (higher means more important) for 'French horn' (left) and 'Garbage truck' (right).

657 **B.3 Feature Visualization validation**

658 Another way of interpreting concepts – as per [22] – is to employ feature visualization methods:
659 through optimization, find an image that maximizes an activation pattern. In our case, we used the set
660 of regularization and constraints proposed by [18], which allow us to successfully obtain realistic
661 images. In Figure S6, we showcase these synthetic images obtained through feature visualization,
662 along with the segments that maximize the target concept. We observe that they do reflect the
663 underlying concepts of interest.

664 Concretely, to produce those feature visualization, we are looking for an image \mathbf{x}^* that is optimized to
665 correspond to a concept from the concept bank \mathbf{W}_i . We use the so called ‘dot-cosim’ loss proposed
666 by [18], which give the following objective:

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

667 With $\mathcal{R}(\cdot)$, the regularizations applied to \mathbf{x} – the default regularizations in the **Xplique** library [67].
668 As for the specific parameters, we used Fourier preconditioning on the image with a decay rate of 0.8
669 and an Adam optimizer ($lr = 1e - 1$).

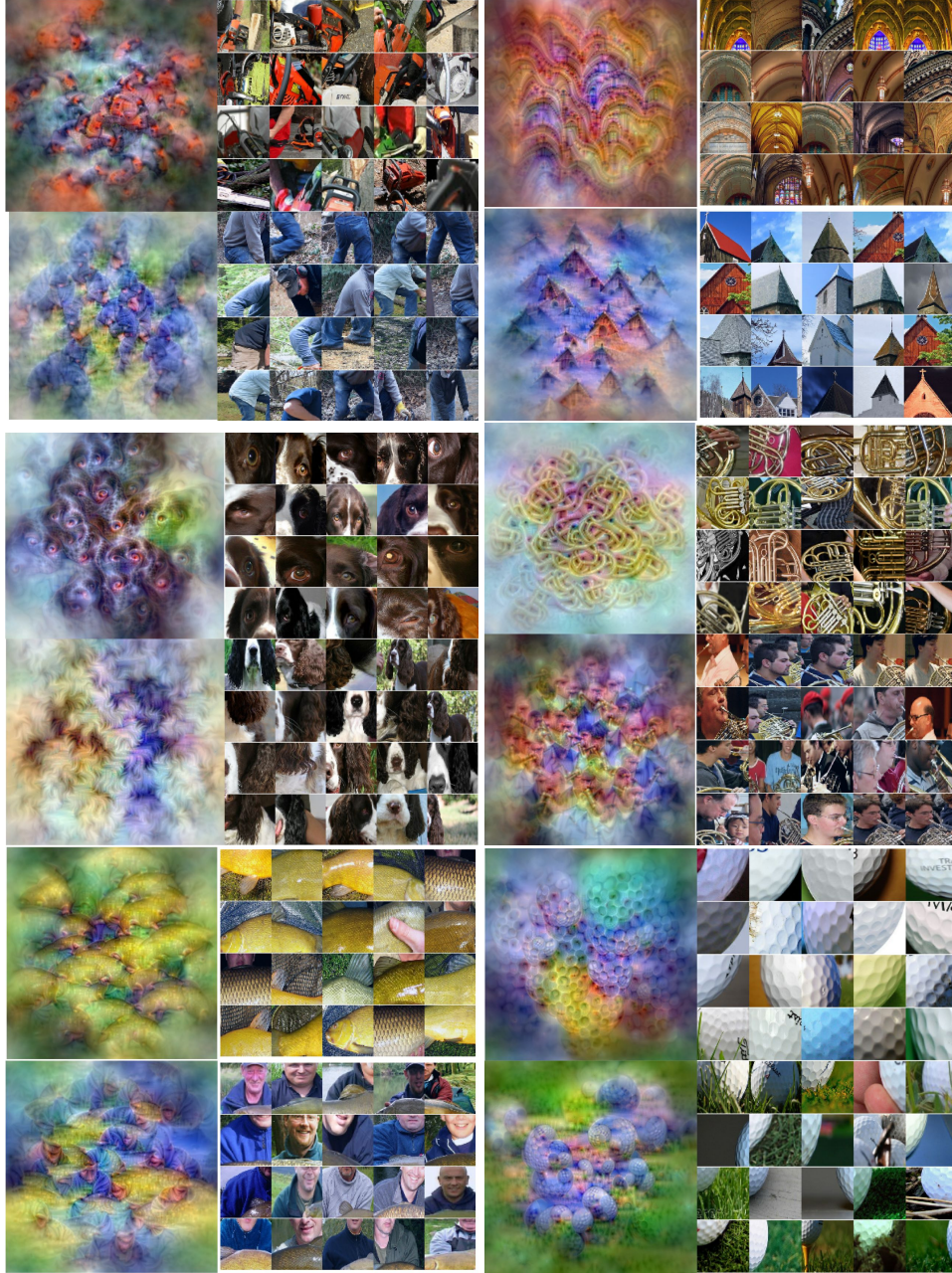


Figure S6: **Feature visualization for CRAFT concepts.** The model under study is a ResNet50 we used the penultimate layer to apply the matrix activation factorization.

670 C Backpropagating through the NMF block

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

672 We recall that NMF decomposes the positive features vector $\mathbf{A} \in \mathbb{R}^{n \times p}$ of n examples lying in
 673 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
 674 $r \ll \min(n, p)$), i.e the solution to the problem:

$$\min_{\mathbf{U} \geq 0, \mathbf{W} \geq 0} \frac{1}{2} \|\mathbf{A} - \mathbf{U}\mathbf{W}^T\|_F^2 \quad (5)$$

675 For simplicity we used a non-regularized version of the NMF objective, following Algorithms 1 and
676 3 in paper [62], based on ADMM [63]. This algorithm transforms the non-linear equality constraints
677 into indicator functions δ . Auxiliary variables \tilde{U} , \tilde{W} are also introduced to separate the optimization
678 of the objective on the one side, and the satisfaction of the constraint on U , W on the other side.
679 The equality constraints $\tilde{U} = U$, $\tilde{W} = W$ are linear and easily handled by the ADMM framework
680 through the associated dual variables \bar{U} , \bar{W} . In our case, the problem in Equation 5 is transformed
681 into:

$$\begin{aligned} \min_{U, \tilde{U}, W, \tilde{W}} \quad & \frac{1}{2} \|A - \tilde{U} \tilde{W}^T\|_F^2 + \delta(U) + \delta(W) \\ \text{s.t.} \quad & \tilde{U} = U, \tilde{W} = W \\ & \text{with } \delta(H) = \begin{cases} 0 & \text{if } H \geq 0 \\ +\infty & \text{otherwise} \end{cases} \end{aligned} \quad (6)$$

682 The constraints are simplified at the cost of a non-smooth (and even a non-finite) objective function
683 $\frac{1}{2} \|A - \tilde{U} \tilde{W}^T\|_F^2 + \delta(U) + \delta(W)$ due to the term $\delta(U) + \delta(W)$. ADMM proceeds to create a
684 so-called *augmented Lagrangian* with l_2 regularization $\rho > 0$:

$$\begin{aligned} \mathcal{L}(A, U, W, \tilde{U}, \tilde{W}, \bar{U}, \bar{W}) = & \frac{1}{2} \|A - \tilde{U} \tilde{W}^T\|_F^2 + \delta(U) + \delta(W) \\ & + \bar{U}^T (\tilde{U} - U) + \bar{W}^T (\tilde{W} - W) \\ & + \frac{\rho}{2} (\|\tilde{U} - U\|_2^2 + \|\tilde{W} - W\|_2^2) \end{aligned} \quad (7)$$

685 The (regularized) problem associated to this Lagrangian is decomposed into a sequence of convex
686 problems that alternate minimization over the U , \tilde{U} , \bar{U} and the W , \tilde{W} , \bar{W} triplets.

$$U_{t+1} = \arg \min_{U=\tilde{U}} \frac{1}{2} \|A - \tilde{U} W_t^T\|_F^2 + \delta(U) + \frac{\rho}{2} \|\tilde{U} - U\|_2^2 \quad (8)$$

$$W_{t+1} = \arg \min_{W=\tilde{W}} \frac{1}{2} \|A - U_t \tilde{W}^T\|_F^2 + \delta(W) + \frac{\rho}{2} \|\tilde{W} - W\|_2^2 \quad (9)$$

687 This guarantees a monotonic decrease of the objective function $\|A - \tilde{U}_t \tilde{W}_t^T\|_F^2$. Each of these
688 sub-problems is thus solved with ADMM separately, by alternating minimization steps of $\frac{1}{2} \|A -$
689 $\tilde{U} W_t^T\|_F^2 + \bar{U}^T (\tilde{U} - U) + \frac{\rho}{2} \|\tilde{U} - U\|_2^2$ over \tilde{U} (i), with minimization steps of $\delta(U) + \frac{\rho}{2} \|\tilde{U} - U\|_2^2$
690 over U (ii), and gradient ascent steps (iii) on the dual variable $\bar{U} \leftarrow \bar{U} + (\tilde{U} - U)$. A similar scheme
691 is used for W updates. Step (i) is a simple convex quadratic program with equality constraints, whose
692 KKT [56, 57] conditions yield a linear system with a Positive Semi-Definite (PSD) matrix. Step (ii)
693 is a simple projection of \tilde{U} onto the convex set $\delta^{-1}(0)$. Finally, step (iii) is inexpensive.

694 Concretely, we solved the quadratic program using Conjugate Gradient [68], from
695 `jax.scipy.sparse.linalg.cg`. This indirect method only involves *matrix-vector* products and can be more
696 GPU-efficient than methods that are based on matrix factorization (such as Cholesky decomposition).
697 Also, we re-implemented the pseudo code in [62] in *Jax* for a fully GPU-compatible program. We
698 used the primal variables U_0 , W_0 returned by `sklearn.decompose.nmf` as a *warm start* for ADMM
699 and observe that the high quality initialization of these primal variables considerably speeds up the
700 convergence of the dual variables.

701 C.2 Implicit differentiation

702 The Lagrangian of the NMF problem reads $\mathcal{L}(U, W, \bar{U}, \bar{W}) = \frac{1}{2} \|A - U W^T\|_F^2 - \bar{U}^T U - \bar{W}^T W$,
703 with dual variables \bar{U} and \bar{W} associated to the constraints $U \geq 0$, $W \geq 0$. It yields a function F
704 based on the KKT conditions [56, 57] whose optimal tuple U , W , \bar{U} , \bar{W} is a root.

705 For single NNLS problem (for example, with optimization over U) the KKT conditions are:

$$\begin{cases} \nabla_U \left(\frac{1}{2} \|\mathbf{A} - \tilde{U}\tilde{W}^T\|_F^2 + \bar{U}^T(-U) \right) = 0, & \text{stationarity} \\ -U \leq 0, & \text{primal feasibility} \\ \bar{U} \odot U = 0, & \text{complementary slackness} \\ \bar{U} \geq 0, & \text{dual feasibility} \end{cases} \quad (10)$$

706 By stacking the KKT conditions of the NNLS problems the we obtain the so-called *optimality*
707 *function F*:

$$F((U, W, \bar{U}, \bar{W}), \mathbf{A}) = \begin{cases} (UW^T - \mathbf{A})W - \bar{U}, \\ (WU^T - \mathbf{A}^T)U - \bar{W}, \\ \bar{U} \odot U, \\ \bar{W} \odot W, \end{cases} \quad (11)$$

708 The implicit function theorem [39] allows us to use implicit differentiation [38, 39, 58] to efficiently
709 compute the Jacobians $\frac{\partial U}{\partial \mathbf{A}}$ and $\frac{\partial W}{\partial \mathbf{A}}$ without requiring to back-propagate through each of the iterations
710 of the NMF solver:

$$\frac{\partial(U, W, \bar{U}, \bar{W})}{\partial \mathbf{A}} = -(\partial_1 F)^{-1} \partial_2 F \quad (12)$$

711 Implicit differentiation requires access to the dual variables of the optimization problem in equation 1,
712 which are not computed by Scikit-learn’s popular implementation. Scikit-learn uses Block coordinate
713 descent algorithm [60, 61], with a randomized SVD initialization. Consequently, we leverage our
714 implementation in Jax based on ADMM [63].

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

718 We use the *Jaxopt* [40] library that allows efficient computation of $\frac{\partial(U, W, \bar{U}, \bar{W})}{\partial \mathbf{A}} = -(\partial_1 F)^{-1} \partial_2 F$.
719 The matrix $(\partial_1 F)^{-1}$ is never explicitly computed – that would be too costly. Instead, the system
720 $\partial_1 F \frac{\partial(U, W, \bar{U}, \bar{W})}{\partial \mathbf{A}} = -\partial_2 F$ is solved with Conjugate Gradient [68] through the use of Jacobian
721 Vector Products (JVP) $v \mapsto (\partial_1 F)v$.

The chain rule yields:

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

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

725 Since r is far smaller ($r = 25$ in all our experiments) than input dimension \mathbf{X} (typically 224×224
726 for ImageNet images), back-propagation is the preferred algorithm in this setting over forward-
727 propagation. We start by computing sequentially the gradients $\nabla_{\mathbf{X}} U_i$ for all concepts $1 \leq i \leq r$.
728 This amounts to compute $v = \nabla_{\mathbf{A}} U_i$ with Implicit Differentiation in Jax, convert the Jax array v
729 into Tensorflow tensor, and then to compute $\nabla_{\mathbf{X}} U_i = \frac{\partial \mathbf{A}}{\partial \mathbf{X}} \nabla_{\mathbf{A}} U_i = \nabla_{\mathbf{X}} (\mathbf{h}_l(\mathbf{X}) \cdot v)$. The latter is
730 easily done in Tensorflow. Finally we stack the gradients $\nabla_{\mathbf{X}} U_i$ to obtain the Jacobian $\frac{\partial U}{\partial \mathbf{X}}$.

731 D Sobol indices for concepts

732 We propose to formally derive the Sobol indices for the estimation of the importance of concepts.
733 Let us define a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ of possible concept perturbations. In order to build
734 these concept perturbations, we start from an original vector of concepts coefficient $\hat{U} \in \mathbb{R}^r$ and
735 use stochastic masks $\mathbf{M} = (M_1, \dots, M_r) \in \mathcal{M} \subseteq [0, 1]^r$, as well as a perturbation operator $\pi :$
736 $\mathcal{A} \times \mathcal{M} \rightarrow \mathcal{A}$ to create stochastic perturbation of \hat{U} that we call concept perturbation $U = \pi(\hat{U}, \mathbf{M})$.

737 Concretely, to create our concept perturbation we consider the inpainting function as our perturbation
738 operator (as in [4, 13, 14]) : $\pi(\tilde{U}, \mathbf{M}) = \tilde{U} \odot \mathbf{M} + (\mathbf{1} - \mathbf{M})\mu$ with \odot the Hadamard product

739 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
 740 function mapping a random concept perturbation \mathbf{U} from the layer l to the output. We denote the
 741 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
 742 perturbation space. Finally, we assume that $\mathbf{f} \in \mathbb{L}^2(\mathcal{A}, \mathbb{P})$ i.e. $|\mathbb{E}(\mathbf{f}(\mathbf{U}))| < +\infty$.

743 The Hoeffding decomposition allows us to express the function \mathbf{f} into summands of increasing
 744 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 $\mathbf{f}(\mathbf{U})$:

$$\begin{aligned} \mathbf{f}(\mathbf{U}) &= \mathbf{f}_{\emptyset} + \sum_i^d \mathbf{f}_i(U_i) + \sum_{1 \leq i < j \leq d} \mathbf{f}_{i,j}(U_i, U_j) + \dots + \mathbf{f}_{1,\dots,r}(U_1, \dots, U_r) \\ &= \sum_{\mathbf{u} \subseteq \mathcal{U}} \mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}}) \end{aligned} \quad (13)$$

745 Eq. 13 consists of 2^r terms and is unique under the following orthogonality constraint:

$$\forall(\mathbf{u}, \mathbf{v}) \subseteq \mathcal{U}^2 \text{ s.t. } \mathbf{u} \neq \mathbf{v}, \quad \mathbb{E}(\mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}})\mathbf{f}_{\mathbf{v}}(\mathbf{U}_{\mathbf{v}})) = 0 \quad (14)$$

746 Furthermore, orthogonality yields the characterization $\mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}}) = \mathbb{E}(\mathbf{f}(\mathbf{U})|\mathbf{U}_{\mathbf{u}}) - \sum_{\mathbf{v} \subset \mathbf{u}} \mathbf{f}_{\mathbf{v}}(\mathbf{U}_{\mathbf{v}})$
 747 and allows us to decompose the model variance as:

$$\begin{aligned} \mathbb{V}(\mathbf{f}(\mathbf{U})) &= \sum_i^d \mathbb{V}(\mathbf{f}_i(U_i)) + \sum_{1 \leq i < j \leq d} \mathbb{V}(\mathbf{f}_{i,j}(U_i, U_j)) + \dots + \mathbb{V}(\mathbf{f}_{1,\dots,r}(U_1, \dots, U_r)) \\ &= \sum_{\mathbf{u} \subseteq \mathcal{U}} \mathbb{V}(\mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}})) \end{aligned} \quad (15)$$

748 Building from Eq. 15, it is natural to characterize the influence of any subset of concepts \mathbf{u} as its own
 749 variance w.r.t. the total variance. This yields, after normalization by $\mathbb{V}(\mathbf{f}(\mathbf{U}))$, the general definition
 750 of Sobol' indices.

751 **Definition D.1** (Sobol indices [46]). The sensitivity index $\mathcal{S}_{\mathbf{u}}$ which measures the contribution of the
 752 concept set $\mathbf{U}_{\mathbf{u}}$ to the model response $\mathbf{f}(\mathbf{U})$ in terms of fluctuation is given by:

$$\mathcal{S}_{\mathbf{u}} = \frac{\mathbb{V}(\mathbf{f}_{\mathbf{u}}(\mathbf{U}_{\mathbf{u}}))}{\mathbb{V}(\mathbf{f}(\mathbf{U}))} = \frac{\mathbb{V}(\mathbb{E}(\mathbf{f}(\mathbf{U})|\mathbf{U}_{\mathbf{u}})) - \sum_{\mathbf{v} \subset \mathbf{u}} \mathbb{V}(\mathbb{E}(\mathbf{f}(\mathbf{U})|\mathbf{U}_{\mathbf{v}}))}{\mathbb{V}(\mathbf{f}(\mathbf{U}))} \quad (16)$$

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

756 Furthermore, the framework of Sobol' indices enables us to easily capture higher-order interactions
 757 between features. Thus, we can view the Total Sobol indices defined in 2 as the sum of of all the Sobol
 758 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
 759 indices using the Jansen estimator [50] and Quasi-Monte carlo Sequence (Sobol LP_{τ} sequence).

760 E Human experiments

761 We first describe how participants were enrolled in the study, then our general experimental design
 762 (See SI for more informations).

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

769 **General study design** For the first experiment – consisting in finding the intruder among elements
 770 of the same concept and an element from a different concept (but of the same class, see Figure S8) –
 771 the choice was randomized in order to avoid any kind of bias due to the order of presentation of the
 772 choices. Moreover, in order to avoid any bias coming from the participants themselves (one group
 773 being more successful than the other) all participants were in the two conditions of finding intruders
 774 in batches of images coming from either concepts or sub-concepts. Concerning experiment 2, the
 775 order was also randomized (see see Figure S9).

776 The participants had to successively find 30 intruders (15 block concepts and 15 block sub-concepts)
 777 for experiment 1 and then make 15 choices (sub-concept vs concept) for experiment 2, see Figure S7.

778 The expert participants are people working in machine learning (researchers, software developers,
 779 engineers) and have participated in the study following an announcement in the authors’ labora-
 780 tory/company. The other participants (Laymen) have no particular competence in machine learning.

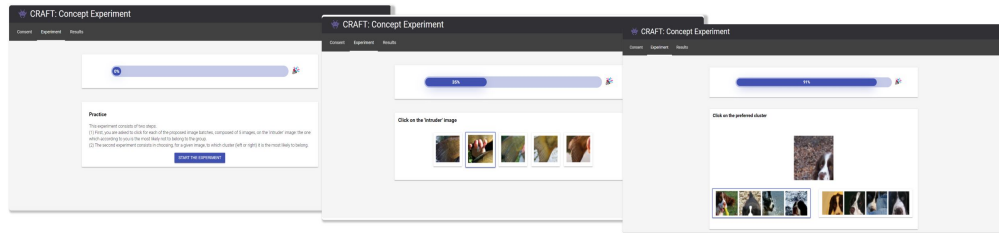


Figure S7: **Human Experiment Website.**

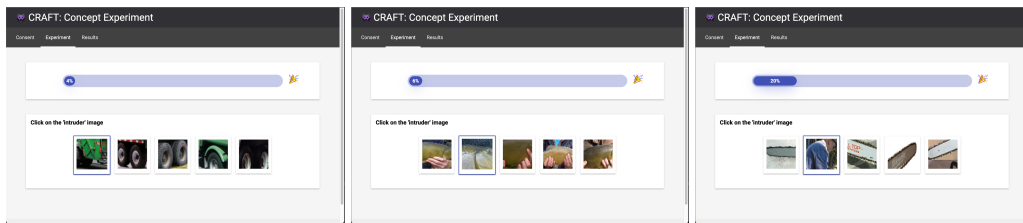


Figure S8: **Binary choice experiment.**



Figure S9: **Intruder experiment.**

781 **F Fidelity experiments**

782 For our experiments on the concept importance measure, we focused on certain classes of IL-
 783 SRVC2012 [27] and used a ResNet50V2 [69] that had already been trained on this dataset. Just like
 784 in [23, 24], we measure the insertion and deletion metrics for our concept extraction technique – as
 785 well as concepts vectors extracted using PCA, ICA and RCA as dimensionality reduction algorithms,
 786 see Figure S10 – and we compare them when we add/remove the concepts as ranked by the TCAV
 787 score [22] and by the Sobol importance score. As originally explained in [13], the objective of
 788 these metrics is to add/remove parts of the input according to how much an explainability method
 789 considers that it is influential and looking at the speed at which the logit for the predicted class
 790 increases/decreases.

791 In particular, for our experimental evaluations, we have randomly chosen 100000 images from
 792 ILSVRC2012 [27] and computed the deletion and insertion metrics for 5 different seeds – for a total
 793 of half a million images. In Figure S10, the shade around the curves represent the standard deviation
 794 over these 5 experiments.

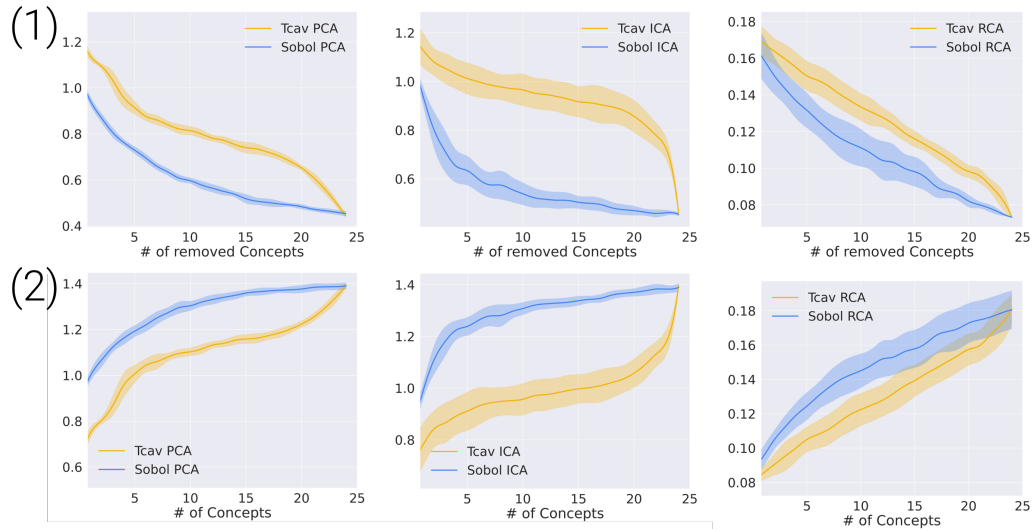


Figure S10: **(1) Deletion curves** for different concept extraction methods, Sobol outperforms TCAV not only for NMF to correctly estimate concept importance (lower is better). **(2) Insertion curves** for different concept extraction methods, Sobol outperforms TCAV to correctly estimate concept importance (higher is better).

795 G Additional examples of concepts and sub-concepts

796 G.1 Sanity Check

797 Following the work from [30], we performed a sanity check on our method, by running the concept
798 extraction pipeline on a randomized model. This procedure was performed on a ResNet-50v2 model
799 with randomized weights. When weights are randomized, concepts are mainly based on color
800 histograms. This might result from skip connections which propagate signal from the inputs.

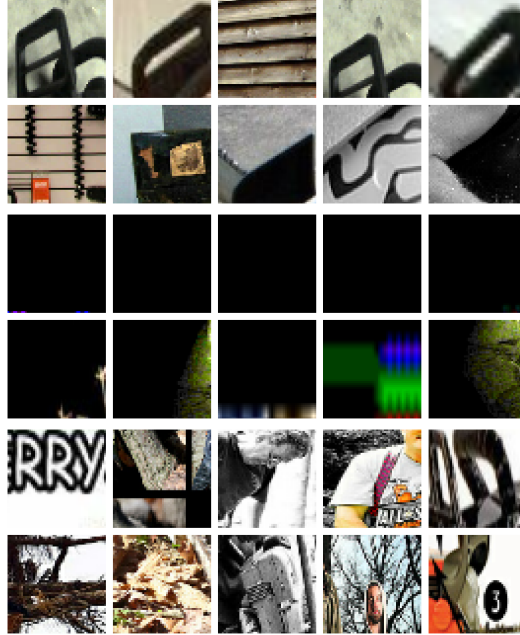


Figure S11: **Sanity check of the method:** we ran the method on a Resnet50 with randomized weights, and extracted the 3 most relevant concepts for the class ‘Chain saw’. When weights are randomized, concepts are mainly based on color histograms. This might result from skip connections which propagate signal from the inputs.

801 H Computational cost

802 Although CRAFT seems like it would require a lot of resources to run, it is actually quite efficient.
803 **Scikit-learn**’s implementation of NMF runs quite fast on the relatively small matrices we work with,
804 and thus, a small amount of steps of ADMM are required; the computation of Sobol indices on
805 only the last layers of the network is not very expensive; and, thanks to the efficiency of **jaxopt**,
806 the concept-wise grad-cAM takes about as much time to calculate as the standard version (for each
807 concept). That being said, the code in its current form doesn’t support batched input images for
808 concept-wise heatmaps, so Smoothgrad [5] and other methods based on the aggregation of gradients
809 will take considerably longer to compute.