



RELAX: Representation Learning Explainability

Kristoffer K. Wickstrøm¹ · Daniel J. Trosten¹ · Sigurd Løkse¹ · Ahcène Boubekki¹ · Karl øyvind Mikalsen¹ · Michael C. Kampffmeyer¹ · Robert Jenssen¹

Received: 4 February 2022 / Accepted: 16 February 2023 / Published online: 11 March 2023
© The Author(s) 2023

Abstract

Despite the significant improvements that self-supervised representation learning has led to when learning from unlabeled data, no methods have been developed that explain what influences the learned representation. We address this need through our proposed approach, RELAX, which is the first approach for attribution-based explanations of representations. Our approach can also model the uncertainty in its explanations, which is essential to produce trustworthy explanations. RELAX explains representations by measuring similarities in the representation space between an input and masked out versions of itself, providing intuitive explanations that significantly outperform the gradient-based baselines. We provide theoretical interpretations of RELAX and conduct a novel analysis of feature extractors trained using supervised and unsupervised learning, providing insights into different learning strategies. Moreover, we conduct a user study to assess how well the proposed approach aligns with human intuition and show that the proposed method outperforms the baselines in both the quantitative and human evaluation studies. Finally, we illustrate the usability of RELAX in several use cases and highlight that incorporating uncertainty can be essential for providing faithful explanations, taking a crucial step towards explaining representations.

Keywords Representation learning · Explainability · Uncertainty · Self-supervised learning

1 Introduction

Interpretability is of vital importance for designing trustworthy and transparent deep learning-based systems (Pedreschi

et al., 2019; Tonekaboni et al., 2019), and the field of explainable artificial intelligence (XAI) has made great improvements over the last couple of years (Antoran et al., 2021; Schulz et al., 2020). However, there exists no methods for attribution-based explanations of *representations*, despite the tremendous advances in representation learning using e.g self-supervised learning (Chen et al., 2020; Caron et al., 2020; He et al., 2020). Also, modifying existing XAI methods to handle representations is often impractical or not possible at all, as explained in “Appendix A”. This lack of explainability makes representation learning less trustworthy and dependable, and there is therefore a need for representation learning explainability. To be able to explain learned representations would provide crucial information in several use-cases. For instance, a typical clustering approach is applying K-means to the representation produced by a feature extractor trained on unlabeled data (Lin et al., 2021; Wen et al., 2020; Yang et al., 2017), but there is no method for investigating which features are characteristic for the members of a cluster.

Representation learning explainability would also allow for a new approach for evaluating representation learning frameworks. Representation learning frameworks are

Communicated by Gang Hua.

✉ Kristoffer K. Wickstrøm
kristoffer.k.wickstrom@uit.no

Daniel J. Trosten
daniel.j.trosten@uit.no

Sigurd Løkse
sigurd.lokse@uit.no

Ahcène Boubekki
ahcene.boubekki@uit.no

Karl øyvind Mikalsen
karl.o.mikalsen@uit.no

Michael C. Kampffmeyer
michael.c.kampffmeyer@uit.no

Robert Jenssen
robert.jenssen@uit.no

¹ Department of Physics and Technology, UiT The Arctic University of Norway, Hansine Hansens veg 18, Tromsø 9019, Troms, Norway

typically evaluated by training simple classifiers on the representation produced by the feature extractor or through a downstream task (Chen et al., 2020; He et al., 2020; Caron et al., 2020). However, such approaches provide only limited information about the features used by the models, and might ignore important distinctions between them. For instance, a similar accuracy on some downstream task does not necessarily equate to the representations being based on the same features. This highlights the need for an explanatory framework for representations, as many of the current evaluation methods are not sufficient for illuminating differences in the what features are used by different feature extractors.

However, any explanatory framework can make over or under-confident explanations. Hence, uncertainty is a key component for designing trustworthy models, since trusting an explanation without knowing the uncertainty of the explanation might lead to an unjustified trust in the model. A recent survey where clinicians were asked what was necessary for making trustworthy models, found that explainability alone was not enough and that uncertainty was also of high importance (Tonekaboni et al., 2019). Our experiments show that uncertainty can be used to increase the faithfulness of explanations, by removing uncertain parts. Nevertheless, little work has been done on uncertainty in explanations of representations.

In this work, we present the first framework for explaining representations, entitled REpresentation LeArning eXplainability (RELAX), which is also equipped with uncertainty quantification with respect to its own explanations. The framework is illustrated in Fig. 1. RELAX measures the change in the representation of an image when compared with masked versions of itself. The core idea is that when informative parts of the input are masked out, the representation should change significantly. When averaging over numerous

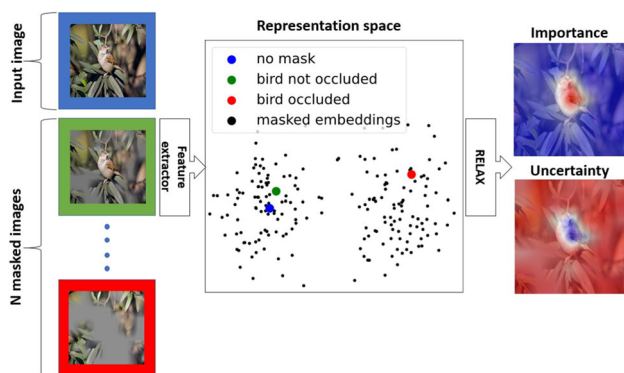


Fig. 1 Conceptual illustration of RELAX. An image is passed through an encoder that produces a new vector representation of the image. Similarly, masked images are embedded in the same latent space. Input feature importance is estimated by measuring the similarity between the representation of the unmasked input with the representations of numerous masked inputs

masked, RELAX reveals the important regions of the input. RELAX is an intuitive and highly versatile framework that can explain any representation, given a suitable similarity function and masking strategy. To provide insight into the geometrical properties of RELAX, we show that the importance of a pixel can be seen as the result of a scoring function based on an inner product between the input and the mean of the masked representations in the representation space. Figure 2 shows an example where RELAX is used to investigate the relevance maps and the corresponding uncertainties for a selection of widely used feature extraction models, which demonstrate that RELAX is a versatile framework for highlighting the emphasis that feature extractors put on pixels and regions in the input (top row).

Our contributions are:

- RELAX, a novel framework for explaining representations that also quantifies its uncertainty.
- A threshold approach called U-RELAX that removes uncertain parts of an explanation and increases the faithfulness of the explanations.
- A theoretical analysis of the framework and an estimation of the number of masks needed to obtain a given level of confidence.
- A comprehensive experimental section that compares widely used supervised and self-supervised feature extraction models and evaluates a number of hyperparameters.
- A user study that examines how well the explanations align with human evaluation.
- Two use cases for RELAX. First, RELAX enables explainability in state-of-the-art incomplete multi-view clustering. This illustrates the usability of RELAX in recent cutting-edge research. Second, RELAX allows for explanation of classic computer vision techniques such as Histogram of Oriented Gradients (HOG). This demonstrates that RELAX is a flexible framework, which is capable of explaining representations produced by any method, not just those produced by deep neural networks.

Code for RELAX is available at <https://github.com/Wickstrom/RELAX>.

2 Related Work

In this section, we present the previous works that are most closely related to our work. The focus will be on attribution-based explanations where each input feature is assigned an importance. Therefore, we will not consider other explainability methods such as example-based explanations (Koh & Liang, 2017; Karimi et al., 2020) or global explanations (Mordvintsev et al., 2015).

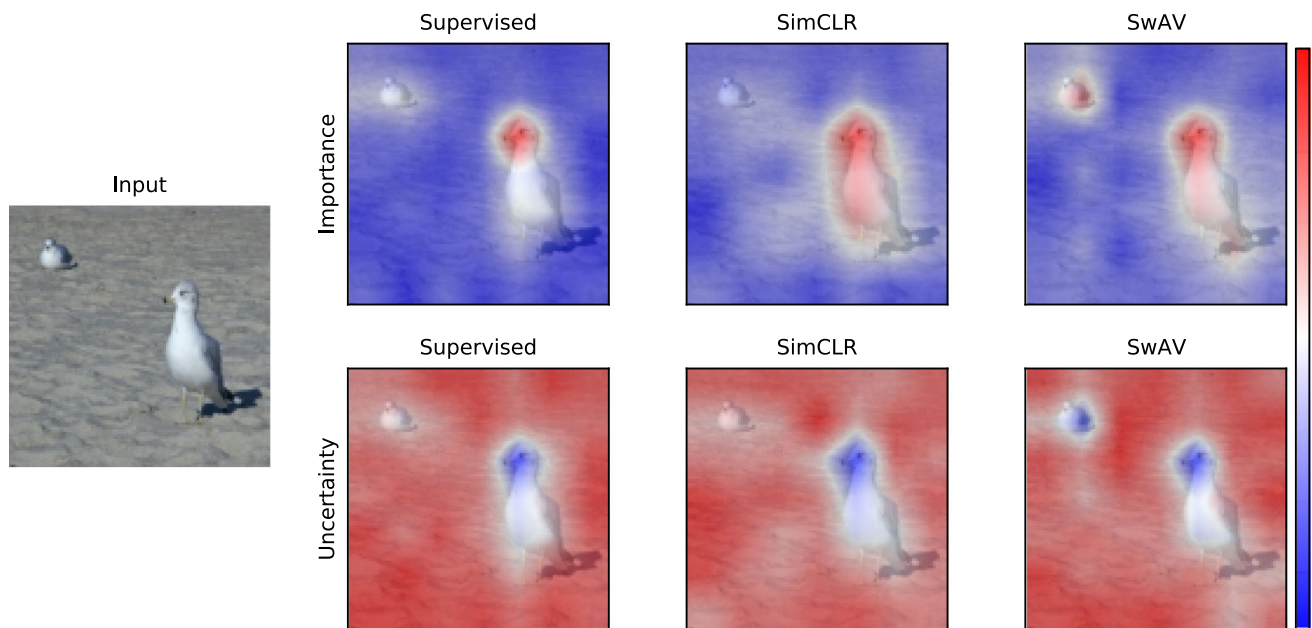


Fig. 2 The figure shows the RELAX importance score and its uncertainty for the representation of the leftmost image for three widely used feature extractors. The first row displays the importance for the representation and the second row shows the uncertainty associated with the different explanations. Red indicates high values and blue indicates low values. In this example, two objects are present in the image, one bird prominently displayed in the foreground, and another more inconspicuous bird in the background. The plots show that all models emphasize

the bird in the foreground with low uncertainty. On the other hand, there is more disagreement on how much emphasis to put on bird in the background, also with a differing degree of uncertainty. The example illustrates that different feature extractors utilize different features in the representation of the image, and with different amounts of uncertainty. The image is taken from VOC (Everingham et al., 2009) (Color figure online)

Occlusion-based explainability There exist a number of occlusion-based explainability methods. Systematically occluding an image with a gray and then measuring the change in activations could be used to provide coarse explanations for CNNs (Zeiler & Fergus, 2014). A more sophisticated occlusion approach can improve explanations, in which smooth masks are generated and accumulated to produce explanations for the prediction of a model (Petsiuk et al., 2018). A slightly different approach is meaningful perturbations, where a spatial perturbation mask that maximally affects the model's output is optimized (Fong & Vedaldi, 2017). A follow up work proposed extremal perturbations, where a perturbation can be considered extremal if it has maximal effect on the network's output among all perturbation of a given, fixed area (Fong et al., 2019). On a different note, an information theoretic approach to XAI has been proposed, where noise is injected in order to measure the information in different regions of the input (Schulz et al., 2020). Similarly, Kolek et al. (2021) introduced a rate-distortion perspective to explainability. Note that none of these methods are capable of providing explanations for representations.

Explaining representations Attribution-based explainability methods are extensively used to explain specific sample predictions (Bach et al., 2015; Petsiuk et al., 2018;

Schulz et al., 2020). However, to the best of our knowledge, no attribution-based explainability method exists for explaining representations. While initial attempts have been made to explain representations such as the Concept Activation Vectors (Kim et al., 2018), which uses directional derivatives to quantify the model prediction's sensitivity, these explanations only relate the representations to high-level concepts and require label information. Similarly, network dissection has been proposed to interpret representations (Bau et al., 2017), but requires predefined concepts and label information without indicating the importance of individual pixels. A different direction is designing models that have the capability to explain their own decisions built into the system (Chen et al., 2019; Alvarez-Melis & Jaakkola, 2018). Two drawbacks of such an approach is that it might lead to models with weaker performance and does not explain representations. Another approach maps semantic concepts to vectorial embedding (Fong & Vedaldi, 2018). However, this requires segmentation masks that are not available in the unsupervised setting. Representations have also been investigated from learnability and describability perspectives (Laina et al., 2020), but this was achieved through human-annotators that are typically not available. Lastly, the inspectability of deep representations have been investigated through an informa-

tion bottleneck approach (Losch et al., 2021), but with a focus on segmentation and predefined concepts.

Uncertainty in explainability Modeling uncertainty in explainability is a rapidly evolving research topic that is receiving an increasing amount of attention. One of the earliest works proposed to use Monte Carlo Dropout (Gal & Ghahramani, 2016) in order to estimate the uncertainty in gradient-based explanations (Wickstrøm et al., 2018, 2020), which was later followed by a similar approach that was based on Layer-wise Relevance Propagation (Bykov et al., 2020). Uncertainties that are inherent in the widely used LIME method (Ribeiro et al., 2016) have been explored (Zhang et al., 2019). Also, ensemble-based approaches, where uncertainty estimates are obtained by taking the standard deviation across the ensemble, have also been proposed (Wickstrøm et al., 2021). Recently, Counterfactual Latent Uncertainty Explanations (CLUE) was presented (Antoran et al., 2021), where uncertainty estimates from probabilistic models can be interpreted. Nevertheless, none of these approaches were designed for quantifying the uncertainty in explanations of representations, as they either require label information or are computationally impractical.

3 Representation Learning Explainability

We present RELAX, our proposed method for explaining representations, equipped with uncertainty quantification. Furthermore, we leverage RELAX’s ability to quantify uncertainty and introduce as a new concept a method for filtering out uncertain parts of the explanations, which we entitle U-RELAX. This is important, as uncertain explanations might give an unwarranted trust in the model. Our framework is inspired by RISE (Petsiuk et al., 2018). However, RISE was designed for explaining predictions and is not transferable for explaining representations or quantifying uncertainty. Note that the proofs of the theorems in this section are given in “Appendix E”.

3.1 RELAX

The central idea of RELAX is that when informative parts are masked out, the representation should change significantly. Let $\mathbf{X} \in \mathbb{R}^{H \times W}$ represent an image¹ consisting of $H \times W$ pixels, and f denote a feature extractor that transforms an image into a representation $\mathbf{h} = f(\mathbf{X}) \in \mathbb{R}^D$. To mask out regions of the input, we apply a stochastic mask $\mathbf{M} \in [0, 1]^{H \times W}$, where each element M_{ij} is drawn from some distribution.

The stochastic variable $\bar{\mathbf{h}} = f(\mathbf{X} \odot \mathbf{M})$, where \odot denotes element-wise multiplication, is a representation of a masked

version of \mathbf{X} . Moreover, we let $s(\mathbf{h}, \bar{\mathbf{h}})$ represent a similarity measure between the unmasked and the masked representation. Intuitively, \mathbf{h} and $\bar{\mathbf{h}}$ should be similar if \mathbf{M} masks *non-informative* parts of \mathbf{X} . Conversely, if *informative* parts are masked out, the similarity between the two representations should be low.

Motivated by this intuition, we define the importance R_{ij} of pixel (i, j) as:

$$R_{ij} = \mathbb{E}_{\mathbf{M}}[s(\mathbf{h}, \bar{\mathbf{h}})M_{ij}]. \quad (1)$$

Equation (1) is core to our framework as it computes the importance of a pixel (i, j) as a weighted similarity score for masked versions of a given image. However, integrating over the entire support of \mathbf{M} is not computationally feasible. Therefore, we approximate the expectation in Eq. (1) by sampling N masks for then to compute the sample mean:

$$\bar{R}_{ij} = \frac{1}{N} \sum_{n=1}^N s(\mathbf{h}, \bar{\mathbf{h}}_n)M_{ij}(n). \quad (2)$$

Here, $\bar{\mathbf{h}}_n$ is the representation of the image masked with mask n , and $M_{ij}(n)$ the value of element (i, j) for mask n . The explanations of RELAX are computed through Eq. (2), and an illustration of RELAX is given in Fig. 1. As a similarity measure we use the cosine similarity

$$s(\mathbf{h}, \bar{\mathbf{h}}) = \frac{\langle \mathbf{h}, \bar{\mathbf{h}} \rangle}{\|\mathbf{h}\| \|\bar{\mathbf{h}}\|}, \quad (3)$$

where $\|\cdot\|$ denotes the Euclidean norm of a vector. There are several motivations for this choice. First, Liu et al. (2021) argued that angular information preserves the essential semantics in neural networks, in contrast to magnitude information. Since the cosine kernel normalizes the representation, essentially discarding magnitude information, such a similarity measure would be suited to capture key information encoded in the representations. We have compared the cosine similarity to the Euclidean distance to examine the effects of including magnitude information, with the results shown in “Appendix B”. Second, the cosine kernel does not rely on hyperparameters that must be selected, which may be beneficial in an unsupervised setting where we cannot do cross validation. Third, a large portion of feature extractors trained using self-supervised learning use the cosine kernel in their loss function (Chen et al., 2020; Chen & He, 2021). Therefore, it is the natural choice for measuring similarities in their latent space. However, based on the two first points, the cosine kernel is still suitable for models trained without the cosine kernel. Lastly, other alternatives for the kernel functions, such as the radial basis function or polynomial kernel, requires careful tuning of hyperparameters. We consider an

¹ To enhance readability, we do not include image channels, but this can be easily included by letting the masks span the channel dimension.

investigation of such alternatives and their hyperparameters as a direction for future research.

Note that we recognize that the masking strategy can introduce a shift in the distribution of pixel intensities. However, in our experiments, we observed that this potential shift did not impact the explanations. An experiment where the distribution is approximately preserved is included in “Appendix C”.

Masking distribution There are several ways to sample the masks in Eq. (2), for instance by letting each $M_{ij}(n)$ be iid. Bernoulli. However, sampling masks with the same size as the input results in a massive sample space, and simultaneously makes it challenging to create smooth masks that cover different portions of the image.²

To avoid these problems, we generate masks as suggested by Petsiuk et al. (2018). Binary masks of smaller size than the input image are generated, where each element of these smaller masks is sampled from a Bernoulli distribution with probability p . These masks are then upsampled using bilinear interpolation to the same size as the image. The distribution for M_{ij} is then a continuous distribution between 0 and 1. Specifically: we sample N binary masks, each with size $h \times w$, where $h < H$ and $w < W$. We upsample these masks to size $(h + 1)C_H \times (w + 1)C_W$, where $C_H \times C_W = \lfloor H/h \rfloor \times \lfloor W/w \rfloor$ is the size of the cell in the upsampled masks. Lastly, we crop the final masks of size $H \times W$ randomly from the $(h + 1)C_H \times (w + 1)C_W$ masks.

Number of masks required In order to minimize the computational cost of RELAX, we derive the following lower bound on the number of masks required for a certain estimation error.

Theorem 3.1 *Suppose $s(\cdot, \cdot)$ is bounded in $(0, 1)$.³ Then, for any $\delta \in (0, 1)$ and $t > 0$, if N in Eq. (2), satisfies:*

$$N \geq -\frac{\ln(\delta/2)}{2t^2}, \quad (4)$$

we have $P(|\bar{R}_{ij} - R_{ij}| \geq t) \leq \delta$.

Theorem 3.1 states that if N satisfies Eq. (4), we are able to estimate R_{ij} to an absolute error of less than t with probability at least $1 - \delta$. See “Appendix E” for proof and verification of bound. In all of our experiments, we generate 3000 masks, which ensures an estimation error below 0.01 with a probability of 0.99.

RELAX from a kernel perspective To provide insights into the geometrical properties of RELAX, we present a kernel viewpoint of Eq. (2).

² See “Appendix D” for evaluation of masking strategies.

³ This holds for the cosine similarity, since the representations considered are assumed to be ReLU outputs (non-negative).

Theorem 3.2 *Suppose the similarity function $s(\cdot, \cdot)$ is a valid Mercer kernel (Mercer, 1909). The importance \bar{R}_{ij} then acts as a linear scoring function between \mathbf{h} , and the weighted mean of $\bar{\mathbf{h}}_1, \dots, \bar{\mathbf{h}}_N$ in the reproducing kernel Hilbert space (RKHS) induced by $s(\cdot, \cdot)$. That is:*

$$\bar{R}_{ij} = \langle \phi(\mathbf{h}), \frac{1}{N} \sum_{n=1}^N \phi(\bar{\mathbf{h}}_n) M_{ij}(n) \rangle_{\mathcal{H}}, \quad (5)$$

where $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$ is the mapping to the RKHS, \mathcal{H} , induced by the kernel $s(\cdot, \cdot)$, and $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is the inner product on \mathcal{H} .

Theorem 3.2 provides interesting insight, as many scoring functions are based on inner-products, e.g. between points of interest and class-conditional means (e.g., Fisher discriminant analysis, Bayes classifier under Gaussian distributions with equal covariance structure). This means that even though RELAX is a novel approach, it is founded in well-known statistical concepts (McCullagh & Nelder, 1989).

Additionally, RELAX has the following interpretation from non-parametric statistics

Theorem 3.3 *Suppose $s(\cdot, \cdot)$ is a valid Parzen window (Theodoridis & Koutroumbas, 2009). Then:*

$$\bar{R}_{ij} \propto p_{ij}(\mathbf{h}), \quad (6)$$

where $p_{ij}(\cdot)$ is a weighted Parzen density estimate (Parzen, 1962) of the density of the masked embeddings:

$$p_{ij}(\cdot) = \frac{1}{\sum_{n'=1}^N M_{ij}(n')} \sum_{n=1}^N s(\cdot, \bar{\mathbf{h}}_n) M_{ij}(n). \quad (7)$$

A high RELAX score is obtained when the unmasked representation \mathbf{h} is close to mean of masked representations, which aligns well with our intuition for RELAX.

3.2 Uncertainty in Explanations

Trusting an explanation without a notion of uncertainty can lead to an unjustified faith in the model. Therefore, we introduce an approach that allows uncertainty quantification to be incorporated into the RELAX framework. Our intuition for this approach stems from what happens when informative and uninformative parts are masked out. If informative parts are masked out, the similarity score will not only drop, but drop with varying degree. If there is a big variation in the similarity scores for a given pixel, it indicates that the explanation for said pixel is uncertain. Based on this intuition, we propose to estimate the uncertainty in input feature importance as:

$$U_{ij} = \text{E}_M[(s(\mathbf{h}, \bar{\mathbf{h}}) - \bar{R}_{ij})^2 M_{ij}]. \quad (8)$$

Again, it is not feasible to integrate over all of \mathbf{M} and U_{ij} is therefore approximated by the sample variance:

$$\bar{U}_{ij} = \frac{1}{N} \sum_{n=1}^N (s(\mathbf{h}, \bar{\mathbf{h}}_n) - \bar{R}_{ij})^2 M_{ij}(n). \tag{9}$$

Equation (9) estimates the uncertainty of the RELAX-score for pixel (i, j) by measuring the spread along M_{ij} between the similarity score and the explanations. In other words, Eq. (9) estimates the uncertainty in the importance scores themselves. To estimate Eq. (9), we must first estimate the importance of a pixel. The uncertainty estimates provided in Eq. (9) can be thought of as measuring the spread of pixel importance values in relation to importance estimated using Eq. (2). There are several benefits of our method. First, it requires no labels, which is sometimes used in other uncertainty estimation methods (Antoran et al., 2021). Secondly, it avoids computationally intense sampling methods, for instance through Monte Carlo sampling (Teye et al., 2018; Gal & Ghahramani, 2016). Lastly, the uncertainty estimation can be combined with the computation of Eq. (2), as explained in Sect. 3.4.

3.3 U-RELAX: Uncertainty Filtered Explanations

All parts of an explanation do not have the same level of uncertainty associated with it. In such cases, it could be beneficial to remove input features that are indicated as important but also have high uncertainty, while only keeping important input features with low uncertainty. This could increase the faithfulness of an explanation and provide clearer explanations. Therefore, we propose a thresholding approach where explanations with high uncertainty are removed from the explanation. We define our U-RELAX importance score as:

$$\bar{R}'_{ij} = \begin{cases} \bar{R}_{ij}, & \text{if } \bar{U}_{ij} < \epsilon \\ 0, & \text{otherwise} \end{cases}, \tag{10}$$

where ϵ is a threshold chosen by the user. Essentially, Eq. (10) provides the possibility to only consider explanations of a particular certainty level, depending on ϵ . We propose two ways of choosing epsilon. First as:

$$\epsilon = \frac{\gamma}{HW} \sum_i^H \sum_j^W \bar{U}_{ij}, \tag{11}$$

that is, the average uncertainty for a particular image, weighted by hyperparameter γ . This provides a simple and intuitive way of selecting the threshold, which is motivated by only wanting to consider pixels that have high importance and low uncertainty. Alternatively, ϵ can be computed by taking the median uncertainty for a particular image. Using

mean or median statistics to select hyperparameters is a common approach in machine learning. For instance, in kernel methods the kernel width is often chosen by taking the mean or median distance between all samples in the training data Nordhaug Myhre et al. (2018); Shi et al. (2009). Determining which approach will give the best performance is dependent on the distribution of the data, in this case the distribution of the uncertainty estimates for a given image. The median is more robust to outliers in the data Leys et al. (2013), and could therefore be a better choice for noisy or challenging samples. If the distribution is symmetric the mean is usually preferred Leys et al. (2013). In Sect. 5.4, we conduct a thorough examination of the mean versus median thresholding approach for U-RELAX.

We refer to this uncertainty-filtered version of RELAX as U-RELAX. Figure 3 shows an example of the U-RELAX explanation compared with the RELAX explanation. In this case, the emphasis on the bird in the background is removed as the uncertainty was too high for this part of the explanation.

3.4 One-Pass Version of RELAX

Computing Eq. (9) requires first computing Eq. (2), since the uncertainty estimation requires an estimate of the importance in order to be computed. This introduces additional computational overhead. We refer to computing Eq. (2) followed by Eq. (9) as the *two-pass* version of RELAX. To improve computational efficiency, we propose an online version of RELAX where importance and uncertainty is computed simultaneously, which we refer to as the *one-pass* version of RELAX. One-pass RELAX is based on well-known estimators of running mean and variance (West, 1979). Importance is computed as:

$$\begin{aligned} \bar{R}_{ij}^{(n)} &= \bar{R}_{ij}^{(n-1)} \\ &+ M_{ij}(n) \frac{s(\mathbf{h}, \bar{\mathbf{h}}_n)(n) - \bar{R}_{ij}^{(n-1)}}{W_{ij}(n)}, \end{aligned} \tag{12}$$

where $\bar{R}_{ij}^{(n)}$ is the importance of pixel (i, j) at mask n , and $W_{ij}(n) = \sum_{n'=0}^n M_{ij}(n')$ is the sum of the mask elements (i, j) for the first n masks. Uncertainty is computed as:

$$\begin{aligned} \bar{U}_{ij}^{(n)} &= \bar{U}_{ij}^{(n-1)} + M_{ij}(n)(s(\mathbf{h}, \bar{\mathbf{h}}_n) \\ &- \bar{R}_{ij}^{(n)})(s(\mathbf{h}, \bar{\mathbf{h}}_n) - \bar{R}_{ij}^{(n-1)}), \end{aligned} \tag{13}$$

where $\bar{U}_{ij}^{(n)}$ is the uncertainty in the importance of pixel (i, j) after the n th mask. Pseudo-code is shown in Algorithm 1. All experiments are carried out using the one-pass version of RELAX. See ‘‘Appendix F’’ for a comparison of the one-pass versus two-pass version.

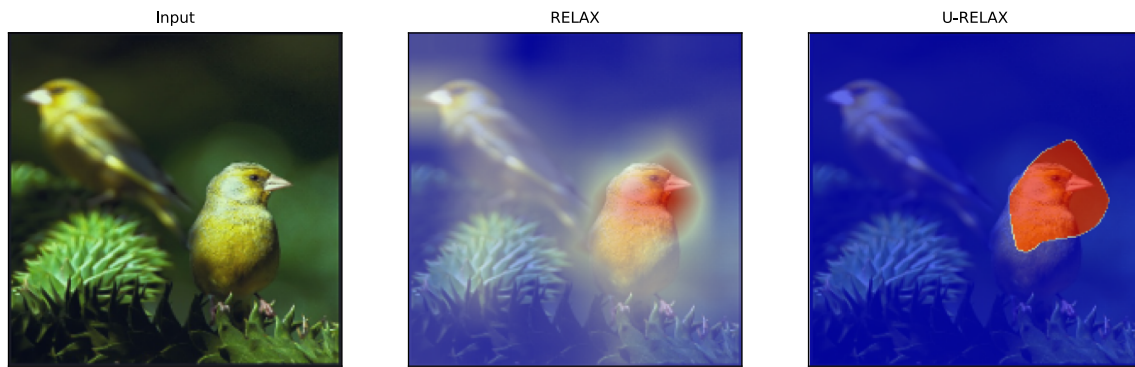


Fig. 3 Comparison of RELAX and U-RELAX on an image taken from PASCAL VOC, where red indicates high importance and blue indicates low importance. In this case, the emphasis on the bird in the background is removed as the uncertainty was too high for this part of the explanation (Color figure online)

Algorithm 1 Pytorch-like pseudocode for RELAX.

```

# f           - feature extractor
# X[1,C,H,W] - input image
# R[H,W]     - importance (init as
                zeros)
# U[H,W]     - uncertainty (init as
                zeros)
# W[H,W]     - sum of masks (init
                with
                small positive number
                ) <
for mask in mask_generator: # [1,1,
                            H,W]
    W += mask
    h, h_mask = f(x), f(x*mask)
    s = cosine_similarity(h, h_mask)

    R_prev = R
    R += m*(s-R)/W
    U += (s-R)*(s-R_prev)*m
return R, U/(W-1)

```

4 Evaluation and Baseline

4.1 Evaluation of Explanations

Evaluation is a developing subfield of XAI, and a unifying score is not agreed upon Doshi-Velez and Kim (2017), even more so for explanations of representations. To evaluate the explanations, we use two of the most widely used explainability evaluation scores, namely localisation and faithfulness (Samek et al., 2017; Petsiuk et al., 2018; Fong et al., 2019; Schulz et al., 2020). All scores are computed using the Quantus toolbox.⁴ Evaluating these metrics is not just important for comparison, but also to ensure the correctness and rigour of RELAX, similarly as done in other works Selvaraju et al.

⁴ <https://github.com/understandable-machine-intelligence-lab/Quantus>.

(2017). By measuring the localisation and faithfulness scores of the explanations created from RELAX we empirically investigate the correctness and reliability of RELAX.

Localisation The explanations should put emphasis on input regions corresponding to the objects present in an image. Localisation measures to which degree the explanation agrees with the ground truth location of an object. High performance in localisation indicates that the explanations often align with the bounding boxes or segmentation masks provided by human annotators. We consider three localisation scores, the *pointing game* (Zhang et al., 2017), *top-k intersection*, and *relevance rank accuracy* (Arras et al., 2022). The pointing game measures whether the pixel with the highest importance is located within the object location. Top-k intersection considers the binarized version of the top-k most important pixels and measures the intersection with the ground truth mask. Relevance rank accuracy is measured by taking the ratio of high intensity relevances within the ground truth mask. Since RELAX operates in the unsupervised setting we do not have explanations for individual classes. Therefore, the bounding boxes/segmentation masks are collected into one unified bounding box/segmentation mask. This results in an unsupervised version of localisation that is suitable for explaining representations.

Faithfulness Pixels assigned with high importance should be indicative of “true” importance. Faithfulness is typically measured by monitoring the classification accuracy of a classifier as input features are iteratively removed. High faithfulness indicates that the explanation is capable of identifying features that are important for classifying an image correctly. We measure faithfulness using the *monotonicity* score. Nguyen and Martinez (2020) proposed to measure monotonicity by computing the correlation of the absolute values of the attributions and the uncertainty in the probability estimation. This will indicate if an explanation is correctly highlighting important features in the input.



Fig. 4 Comparison of RELAX and saliency explanation for an image from PASCAL VOC. The example shows how both explanations focus on the dog, but the saliency explanation is much more erratic and unfocused than the RELAX explanations

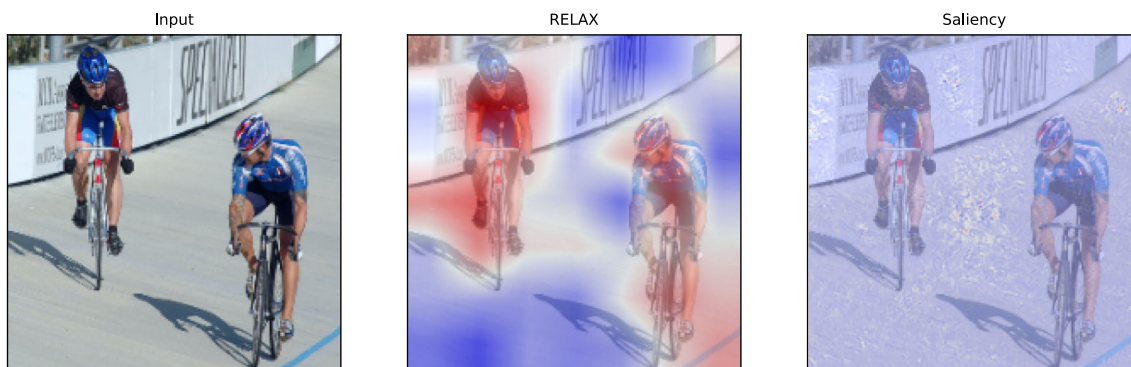


Fig. 5 Comparison of RELAX and Saliency explanation for an image from PASCAL VOC. The example shows how RELAX captures information about both objects, while the saliency explanation is focused on the gap in between the two objects

4.2 Representation Explainability Baseline

While there are no existing methods that provide attribution-based explanations for representations, it is possible to adopt certain methods to provide such explanations. One of the most common baselines in the field of explainability is saliency explanations (Springenberg et al., 2015; Adebayo et al., 2018), which utilize gradient information to attribute importance. An explanation is obtained by computing the gradient for a prediction with respect to the input. However, it is not trivial to extend such methods for explaining representations. We propose the following for a saliency approach:

$$S = \frac{1}{D} \sum_{d=1}^D \nabla f(\mathbf{X})_d, \tag{14}$$

where D is the dimensionality of the representation and S_{ij} is the importance of pixel (i, j) for the given representation. The gradient for each dimension of the representation will give an explanation, and Eq. (14) takes the mean across all explanations. This is the most straight-forward and intuitive approach for explaining representations with gradients.

It also illustrates the challenges that arise when adopting gradient-based explanations for representation, as some form of agglomeration of the explanations is required. Figures 4 and 5 shows a qualitative comparison between the RELAX and saliency explanation for a representation of an image. Both figures illustrate how RELAX provides more intuitive and clear explanations that are able to capture information related to the objects in the image, when compared with the saliency explanation.

Once the saliency approach from Eq. (14) have been established, it is also possible to adopt improvements of the standard saliency explanations. For instance, Guided Backpropagation is a widely used explainability technique that uses gradient information (Springenberg et al., 2015). Guided Backpropagation differs from Eq. (14) by zeroing out negative gradients in the backward pass of the backpropagation scheme. We define the Guided Backpropagation procedure for representations as:

$$S_{GB} = \frac{1}{D} \sum_{d=1}^D \nabla_{GB} f(\mathbf{X})_d. \tag{15}$$

Second, SmoothGrad is another gradient-based explainability method that can be adopted from Eq. 14 (Smilkov et al.,

2017). SmoothGrad injects noise into the input and produces an explanation by averaging over multiple explanations. We define SmoothGrad for representation as:

$$\mathbf{S}_{\text{SG}} = \frac{1}{M} \sum_{m=1}^M \frac{1}{D} \sum_{d=1}^D \nabla f(\mathbf{X}_m)_d, \quad (16)$$

where M is the number of explanations computed based on the noisy input.

Adaptation of state-of-the-art methods As explained in “Appendix A”, many of the existing explanation methods are not trivially extended to the representation learning explainability setting. Nevertheless, using the baselines introduced above we can construct adaptation of the state-of-the-art algorithms integrated gradients Sundararajan et al. (2017) and GRAD-CAM Selvaraju et al. (2017). For the integrated gradients explanations, we follow their proposed procedure but compute gradients using Eq. 14. For the GRAD-CAM explanations, the upsampled output of the global average pooling layer is typically weighted by the class weights. However, these class weights are not available in our unsupervised representation learning setting. Therefore, we weight all parts equally, and gradients are computed using Eq. 14.

5 Experiments

To evaluate RELAX, we conduct numerous experiments and report both quantitative and qualitative results. We evaluate several features extraction models, both deep and non-deep, and trained with and without supervision. Our experiments show the advantages of RELAX compared to the baselines, and illustrate how RELAX enables new approaches for analysing and understanding representation learning.

Implementation details. For the supervised model, we use the pretrained model from Pytorch (Paszke et al., 2019). For the models trained without labels but with self-supervision, we use the SimCLR (Chen et al., 2020) and SwAV (Caron et al., 2020) frameworks, both of which have seen recent widespread use. These methods are chosen to represent two major types of self-supervised learning frameworks, namely contrastive instance learning (SimCLR) and clustering-based learning (SwAV). For SimCLR and SwAV, we use the pretrained models from Pytorch Lightning Bolts (Falcon & Cho, 2020). We use a ResNet50 (He et al., 2016) as the backbone for the feature extractors, and all models are trained on ImageNet (Deng et al., 2009). Additionally, we also perform experiments with recent Vision Transformer architectures (Dosovitskiy et al., 2021). The results of these experiments are shown in “Appendix G”.

Similarly as in previous works (Fong et al., 2019; Schulz et al., 2020), we use the test split of the PASCAL VOC07 (VOC) (Everingham et al., 2009) and the validation split of MSCOCO2014 (COCO) (Lin et al., 2014) for evaluating the localisation scores, since they contain information about the location of the objects in the images. For the faithfulness score, we use the validation set of ImageNet (Deng et al., 2009). For all datasets, we randomly sample 1000 images for evaluation and repeat all experiments 3 times. Since we are interested in investigating how RELAX and U-RELAX vary due to the stochastic masking process, we use the same 1000 images across the repeated experiments. We generate 3000 masks to ensure a low estimator error. We set $h = w = 7$ and resize all images to $H = W = 224$, as suggested by Zhang et al. (2017). For the monotonicity score, we use Alexnet (Krizhevsky et al., 2012) as the classifier, as suggested by Samek et al. (2017). We also experiment with the VGG13 (Simonyan & Zisserman, 2015) as the classifier for monotonicity score. These results are reported in “Appendix H”. The threshold for U-RELAX is determined with median aggregation and $\gamma = 1.0$, based on the empirical evaluation conducted in Sect. 5.4.

5.1 Qualitative Results

Figures 2 and 6 displays the explanation and the uncertainty in the explanations provided by RELAX for an image from the PASCAL VOC and MS COCO dataset, respectively. See “Appendix J” for additional qualitative results. The input to the feature extractors is shown on the left, the first row shows the explanations, and the second row shows the uncertainties.

Are all instances of the same object equally important?

Figure 2 shows an example with two objects, one bird prominently displayed in the foreground, and another more inconspicuous bird in the background. An interesting question that RELAX allows us to answer is: are both of these birds important for the representation of this image? And, are both of them equally important? First, all models indicate that the bird in the foreground is important, and that the importance scores for this bird have low uncertainty. Second, SimCLR puts little emphasis on the bird in the background. In contrast, both the supervised feature extractor and SwAV are highlighting the second bird as having an influence on the representation. However, the uncertainty estimates for the second bird is slightly higher than those of the first bird, but still low compared to the remaining parts of the image.

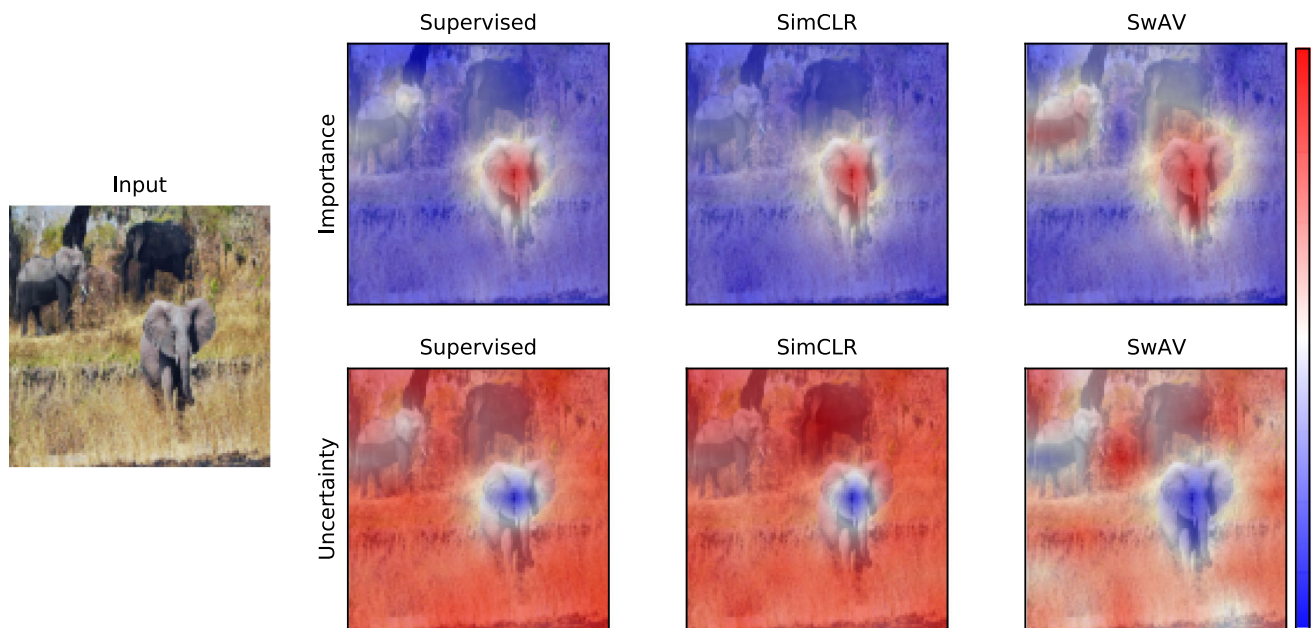


Fig. 6 The figure shows the RELAX explanation and its uncertainty for the representation of the leftmost image for a number of widely used feature extractors. The first row displays the explanations for the representation and the second row shows the uncertainty associated with the different explanations. Red indicates high values and blue indicates low values. In this example, three elephants are visible in the image. The

results show that all models highlight the elephant in the foreground as important for the representation, but there is more disagreement about the elephants in the background. Moreover, the uncertainty of the explanation for the elephant in the foreground is very low compared to the remaining regions of the image. Image is taken from MS COCO (Color figure online)

What features are important in complex images with numerous objects?

Figure 6 shows an image with 3 elephants, one in the foreground and two in the background. Additionally, the background is more diverse and the objects have different lighting and perspective. Again, RELAX enables investigation of interesting aspects of the representations, such as: are the models capable of recognizing all elephants and to utilize the information? Does the models focus on background information instead of the objects? All models highlight the elephant in the foreground as important with high certainty. However, there is little emphasis on the shaded elephant, and the associated region of the image also has a high degree of uncertainty. Both the supervised model and SwAV put some importance on the third elephant with some degree of certainty, while SimCLR uses little or no information about the third elephant.

In both Figs. 2 and 6, the SwAV feature extractor is focusing on several regions in the input, but with some regions of high uncertainty. While it is difficult to say exactly why, we hypothesize that it can be related to its self-supervised training procedure. SwAV relies on matching image views to a set of prototypes. Therefore, different parts of the input can be related to different prototypes, which we conjecture can lead to SwAV considering several regions of the input.

5.2 Quantitative Results

Tables 1 and 2 displays the quantitative evaluation of our proposed methodology compared with the gradient-based baselines described in Sect. 4.2. The results show how the proposed method outperforms the baselines across all scores. The low standard deviation for RELAX show that the proposed methodology is robust to the stochasticity in the masks. Furthermore, the feature extractor trained using supervised learning achieves the highest performance compared to the feature extractors trained using self-supervised learning, which illustrates that label information does provide additional useful information for these scores.

For the localisation scores, RELAX provides the highest performance. The segmentation masks or bounding boxes can in many cases be large, and U-RELAX might remove uncertain points close to the boundaries of the segmentation masks. This might be desirable from a human perspective, as it provides clearer explanations with less uncertainty, but it will decrease the localisation scores. For the faithfulness score, U-RELAX provides a significant boost in performance for two encoders. The removal of uncertain explanations allows the classifier to focus on a smaller subset of highly relevant features. This can lead to the classifier having a more stable decrease in accuracy and a higher faithfulness score.

Table 1 Pointing game, top k, and relevance rank scores in percentages and averaged over 3 runs

Scores	Methods	Supervised		SimCLR		SwAV	
		COCO	VOC	COCO	VOC	COCO	VOC
Pointing game	Saliency	67.1 ± 0.0	82.8 ± 0.0	59.9 ± 0.0	75.9 ± 0.0	60.0 ± 0.0	76.3 ± 0.0
	Smooth Saliency	62.8 ± 0.0	79.5 ± 0.0	60.1 ± 0.0	75.9 ± 0.0	59.8 ± 0.0	76.4 ± 0.0
	Guided Saliency	66.6 ± 0.0	82.9 ± 0.0	58.4 ± 0.0	73.3 ± 0.0	59.5 ± 0.0	75.8 ± 0.0
	Integrated Gradients	47.8 ± 0.0	59.1 ± 0.0	32.9 ± 0.0	48.2 ± 0.0	36.5 ± 0.0	51.5 ± 0.0
	Grad CAM	66.8 ± 0.4	78.7 ± 0.5	47.7 ± 0.7	57.0 ± 0.6	48.7 ± 1.0	58.6 ± 0.8
	RELAX	72.6 ± 0.1	86.6 ± 0.2	68.7 ± 0.3	85.2 ± 0.3	67.8 ± 0.2	84.7 ± 0.2
	U-RELAX	72.1 ± 0.3	86.4 ± 0.4	68.6 ± 0.2	85.0 ± 0.5	66.7 ± 0.7	84.1 ± 0.4
Top k	Saliency	62.2 ± 0.0	80.1 ± 0.0	56.5 ± 0.0	71.3 ± 0.0	56.5 ± 0.0	71.4 ± 0.0
	Smooth Saliency	59.2 ± 0.0	74.1 ± 0.0	56.4 ± 0.0	71.1 ± 0.0	56.4 ± 0.0	71.3 ± 0.0
	Guided Saliency	62.2 ± 0.0	80.2 ± 0.0	55.1 ± 0.0	69.0 ± 0.0	56.3 ± 0.0	71.1 ± 0.0
	Integrated Gradients	47.7 ± 0.0	61.0 ± 0.0	35.4 ± 0.0	52.8 ± 0.0	33.2 ± 0.0	49.0 ± 0.0
	Grad CAM	64.0 ± 0.0	78.3 ± 0.0	43.6 ± 0.0	55.3 ± 0.0	43.1 ± 0.1	54.8 ± 0.0
	RELAX	72.8 ± 0.4	86.9 ± 0.1	69.0 ± 0.3	85.6 ± 0.2	68.1 ± 0.4	85.1 ± 0.2
	U-RELAX	72.2 ± 0.4	86.5 ± 0.2	68.8 ± 0.4	85.3 ± 0.1	66.6 ± 0.4	84.2 ± 0.3
Relevance rank	Saliency	46.8 ± 0.0	59.5 ± 0.0	41.2 ± 0.0	53.6 ± 0.0	40.9 ± 0.0	53.4 ± 0.0
	Smooth Saliency	42.6 ± 0.0	54.6 ± 0.0	41.1 ± 0.0	53.4 ± 0.0	40.9 ± 0.0	53.3 ± 0.0
	Guided Saliency	46.8 ± 0.0	59.8 ± 0.0	40.6 ± 0.0	53.0 ± 0.0	40.9 ± 0.0	53.3 ± 0.0
	Integrated Gradients	38.4 ± 0.0	51.9 ± 0.0	31.9 ± 0.0	47.2 ± 0.0	32.3 ± 0.0	48.3 ± 0.0
	Grad CAM	46.0 ± 0.0	60.2 ± 0.0	37.5 ± 0.0	50.7 ± 0.0	37.8 ± 0.0	50.9 ± 0.0
	RELAX	56.4 ± 0.0	70.2 ± 0.1	54.2 ± 0.2	69.8 ± 0.1	52.4 ± 0.1	69.1 ± 0.0
	U-RELAX	52.4 ± 0.0	64.7 ± 0.1	50.7 ± 0.1	63.3 ± 0.1	46.2 ± 0.1	59.5 ± 0.0

Higher is better and bold numbers highlight the top performance. Results show that our method improves on the baseline across all scores

Table 2 Monotonicity scores averaged over 3 runs

Scores	Methods	Supervised	SimCLR	SwAV
Monotonicity	Saliency	12.8 ± 0.2	14.8 ± 0.5	14.6 ± 0.3
	Smooth Saliency	15.4 ± 0.1	14.3 ± 0.3	14.0 ± 0.3
	Guided Saliency	15.3 ± 0.3	15.3 ± 0.2	14.2 ± 0.6
	Integrated Gradients	12.4 ± 0.3	11.9 ± 0.3	13.5 ± 0.5
	Grad CAM	19.6 ± 0.2	18.6 ± 0.3	18.5 ± 0.3
	RELAX	18.3 ± 0.5	20.2 ± 0.4	21.3 ± 0.4
	U-RELAX	23.6 ± 0.4	22.9 ± 0.1	18.3 ± 0.6

Higher is better and bold numbers highlight the top performance. Results show that our method improves on the baseline

Table 3 Human evaluation of representation explainability methods across 10 images from the PASCAL VOC dataset

	RELAX	U-RELAX	Saliency	Smooth Saliency	Guided Saliency	Random
Counts	79	29	9	4	8	1

Results show that the majority of the votes were cast for RELAX and U-RELAX

5.3 Human Evaluation

The localisation and faithfulness scores are both proxies for human evaluation that allow for quantitative analysis. However, the ultimate goal of XAI is to provide explanations that are understandable for people and align well with human intuition. Therefore, we conduct a user study with human

evaluation of explanations. Our study is inspired by the localisation scores but rely on evaluation of individual humans instead of segmentation masks or bounding boxes. In this user study, 13 people were asked to select their preferred explanation from a selection of explanations across 10 different images. See “Appendix I” for a detailed description of the user study.

Table 4 Evaluation of U-RELAX hyperparameters in terms of pointing game, top k, and relevance rank scores in percentages and averaged over 3 runs

Scores	(aggregation, γ)	Supervised		SimCLR		SwAV	
		COCO	VOC	COCO	VOC	COCO	VOC
Pointing game	(mean, 0.95)	71.1 ± 0.4	86.5 ± 0.2	67.6 ± 0.1	83.9 ± 0.3	63.3 ± 0.7	81.1 ± 0.5
	(mean, 0.99)	71.8 ± 0.4	86.4 ± 0.5	68.6 ± 0.4	85.0 ± 0.4	66.4 ± 0.6	84.2 ± 0.4
	(mean, 1.0)	71.7 ± 0.1	86.5 ± 0.2	68.6 ± 0.1	85.0 ± 0.3	66.7 ± 0.7	84.1 ± 0.2
	(median, 0.95)	71.2 ± 0.2	86.6 ± 0.1	67.6 ± 0.4	84.2 ± 0.2	63.6 ± 0.2	80.9 ± 0.1
	(median, 0.99)	71.8 ± 0.3	86.5 ± 0.4	68.8 ± 0.3	85.0 ± 0.2	66.3 ± 0.6	84.0 ± 0.3
	(median, 1.0)	72.1 ± 0.3	86.4 ± 0.4	68.6 ± 0.2	85.0 ± 0.5	66.7 ± 0.7	84.1 ± 0.4
Top k	(mean, 0.95)	71.3 ± 0.4	86.2 ± 0.2	67.1 ± 0.1	83.2 ± 0.3	62.8 ± 0.2	79.5 ± 0.4
	(mean, 0.99)	72.2 ± 0.4	86.6 ± 0.2	68.8 ± 0.3	85.2 ± 0.2	66.4 ± 0.2	84.0 ± 0.3
	(mean, 1.0)	72.2 ± 0.4	86.5 ± 0.2	68.8 ± 0.4	85.3 ± 0.1	66.7 ± 0.4	84.3 ± 0.2
	(median, 0.95)	71.2 ± 0.4	86.1 ± 0.2	67.1 ± 0.2	83.2 ± 0.4	62.7 ± 0.2	79.1 ± 0.4
	(median, 0.99)	72.2 ± 0.4	86.5 ± 0.2	68.7 ± 0.3	85.2 ± 0.2	66.4 ± 0.2	83.9 ± 0.3
	(median, 1.0)	72.2 ± 0.4	86.5 ± 0.2	68.8 ± 0.4	85.3 ± 0.1	66.6 ± 0.4	84.2 ± 0.3
Relevance rank	(mean, 0.95)	45.9 ± 0.0	55.7 ± 0.0	41.6 ± 0.1	52.3 ± 0.1	39.6 ± 0.1	51.0 ± 0.0
	(mean, 0.99)	50.3 ± 0.0	61.2 ± 0.1	48.6 ± 0.1	59.8 ± 0.1	44.0 ± 0.1	56.0 ± 0.1
	(mean, 1.0)	51.4 ± 0.1	63.0 ± 0.1	50.3 ± 0.1	62.2 ± 0.1	45.6 ± 0.1	58.2 ± 0.1
	(median, 0.95)	46.8 ± 0.0	57.2 ± 0.1	42.4 ± 0.1	53.3 ± 0.1	40.4 ± 0.1	52.1 ± 0.0
	(median, 0.99)	51.2 ± 0.0	63.0 ± 0.1	49.1 ± 0.1	60.8 ± 0.1	44.6 ± 0.1	57.3 ± 0.1
	(median, 1.0)	52.4 ± 0.0	64.7 ± 0.1	50.7 ± 0.1	63.3 ± 0.1	46.2 ± 0.1	59.5 ± 0.0

Higher is better and bold numbers highlight the top performance

Table 5 Evaluation of U-RELAX hyperparameters in terms of monotonicity score in percentages and averaged over 3 runs

Scores	(aggregation, γ)	Supervised	SimCLR	SwAV
Monotonicity	(mean, 0.95)	16.3 ± 0.5	11.8 ± 0.3	12.4 ± 0.3
	(mean, 0.99)	22.2 ± 0.2	20.4 ± 0.5	16.2 ± 0.3
	(mean, 1.0)	23.2 ± 0.1	21.8 ± 0.3	18.0 ± 0.0
	(median, 0.95)	17.9 ± 0.7	12.8 ± 0.2	13.5 ± 0.2
	(median, 0.99)	23.0 ± 0.7	21.1 ± 0.1	17.1 ± 0.4
	(median, 1.0)	23.6 ± 0.4	22.9 ± 0.1	18.3 ± 0.6

Higher is better and bold numbers highlight the top performance

Table 3 reports the results of the human evaluation. The results clearly indicate that RELAX and U-RELAX were the methods that aligned most closely with human intuition. Some participants highlighted that when both RELAX and the gradient-based methods indicated an object as important, they often preferred the more object focused explanation of RELAX, as opposed to the more edge focused explanations of the baselines. It was also noted that for some images the participants disagreed with most explanations, and would have provided a different explanation if possible. We believe that these are valuable insights that will be useful for improving explainability methods and also for designing future user studies.

5.4 U-RELAX Hyperparameter Evaluation

Tables 4 and 5 reports localisation and faithfulness scores for different values of the hyperparameters in U-RELAX. Mean versus median aggregation is considered, and a selection of values for γ . The results indicate that setting γ to less than 1, typical degrades performance. This can be understood by the thresholding being to strict and removing to many pixel indicated as important. Also, the differences between mean and median aggregation of the uncertainties is mostly low, but median aggregation gives a slight improvement, particularly for the relevance rank score and the monotonicity score.

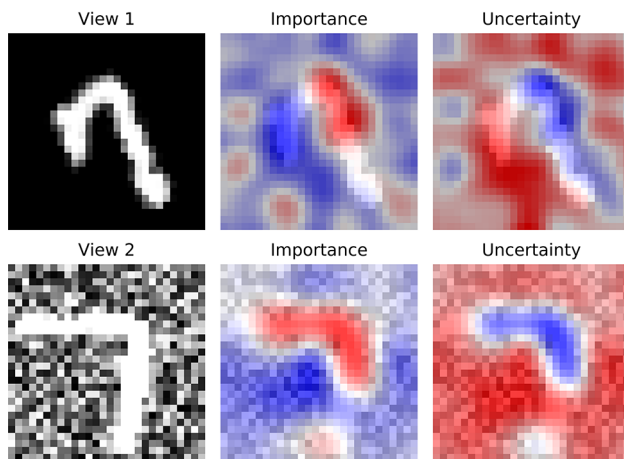


Fig. 7 RELAX explanation and uncertainty for the representation of an example from Noisy MNIST image for a number of widely used feature extractors. The first row displays input, explanation, and uncertainty for view 1, and the second row for view 2. Red indicates high values and blue indicates low values. The figure shows that Completer is extracting complementary information from the two views for creating its unified representation (Color figure online)

5.5 Use Case I: Multi-View Clustering

To further illustrate the ability of RELAX to obtain insights into new tasks, we conduct an experiment on multi-view clustering. We learn a feature extractor using the Completer framework (Lin et al., 2021), which uses an information theoretic approach to fuse several views into a new representation. Completer uses individual encoders for each view, and concatenates the representation from each encoder to produce a unified representation. Clustering is performed by applying K-means to the learned representations. To adopt RELAX for such a setting, we generate individual masks for each view and monitor the change in the representation in the unified representation space. While there is no way to investigate which parts of the different views that influence the unified representation in the Completer framework, using RELAX allows us to answer this question. Figure 7 shows an example

on Noisy MNIST (Wang et al., 2015), where one view is a digit and the other view is a noisy version of the same digit. The result shows that the Completer framework is exploiting information from both views to produce a new representation, even if one view contains more noise. Such insights would not be obtainable without RELAX.

5.6 Use Case II: Explaining HOG Features

RELAX is not limited to representations produced by deep neural networks. It can be used to explain the representation produced by any function that transform an image into a vector representation. To illustrate the versatility of RELAX, we explain representation produced by the Histogram of Oriented Gradients (HOG) feature extraction method (Dalal & Triggs, 2005), which have been used extensively in the computer vision literature. Figures 8 and 9 shows two examples where the relevance map for the HOG representation is compared with the SimCLR and SwAV representations. We consider the representations from these two methods since they are also unsupervised like the HOG features.

Features produced by deep neural networks typically allow for higher performance than those from algorithms such as HOG and other handcrafted feature extraction methods. RELAX provides insights into why this is. In Fig. 8, both the SimCLR and the SwAV feature extractors focus on the cat in the center of the images. The HOG algorithm has a more widespread focus. Also, much of the emphasis is put on the cord going along the staircase. Since the HOG algorithm is utilizing gradient information, these sharp lines will have a big influence on the representation, and it is therefore not surprising that the cat receives less attention. In Fig. 8, both SimCLR and SwAV focus on the bird, while the HOG features are more focused on other regions in the image. For instance, the iron rod and a tree in the background and are indicated as being important for the representation of this image. Both examples provide insights into why HOG features lead to inferior performance, when compared with

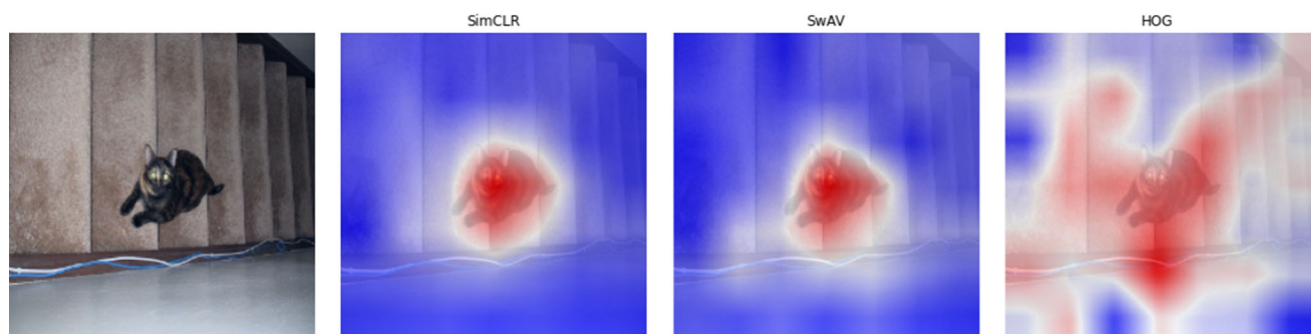


Fig. 8 The figure shows the RELAX explanation for two deep learning-based feature extractors compared with the traditional HOG algorithm. Figure shows how HOG features focus on more indistinct regions in the input, while deep learning methods focus mainly on the cat. Image is taken from PASCAL VOC

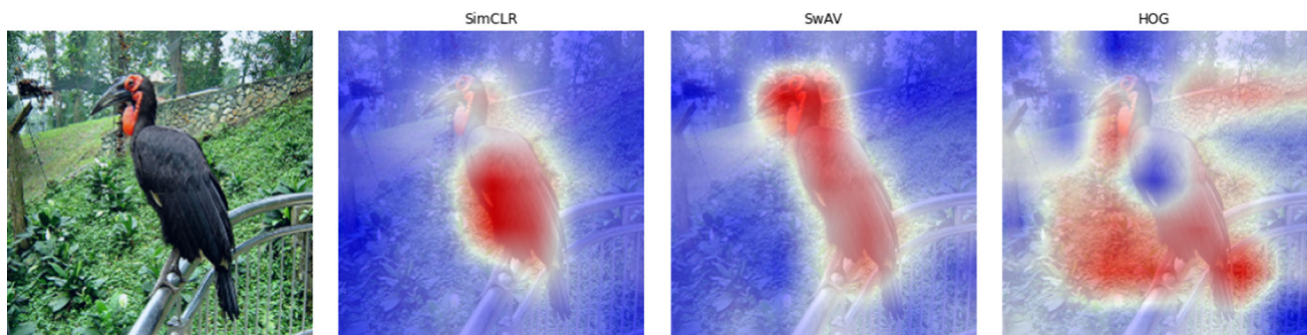


Fig. 9 The figure shows the RELAX explanation for two deep learning-based feature extractors compared with the traditional HOG algorithm. Figure shows how HOG features puts little attention on the bird and mostly focus on the background. Image is taken from PASCAL VOC

features produced by deep neural networks. This information would not be available without the proposed RELAX framework.

6 Conclusion

In this work, we presented RELAX, a framework for explaining representations produced by any feature extractor. RELAX is based on masking out parts of an image and for then to measure the similarity with an unmasked version in the representation space. We introduced a principled approach to quantifying uncertainty in explanations. RELAX was evaluated by comparing several widely used feature extractors. Results indicate that there can be a big difference in the quality of the explanations. It was shown that filtering out parts of an explanation based on its uncertainty can improve the faithfulness, and that RELAX can have a facilitating role, providing explainability for several downstream applications such as multi-view clustering. We consider the evaluation of RELAX to other use-cases, such as for the investigation of a models failure cases, as an interesting direction for future research. We believe that RELAX can be an important addition in the intersection between XAI and representation learning.

Acknowledgements The authors would like to thank Nils Midtbø, Caroline Granås, Theodor Ross, Suaiba Salahuddin, Sigrid Vold Jensen, Thomas Johansen, Julianne Nyvold, Kristoffer Furøy, Jostein Henriksen, Erland Grimstad, Jonas Mørch-Lampe, Inger Solheim, and Andreas Kvammen for participating in the user study through human evaluation.

Funding Open access funding provided by UiT The Arctic University of Norway (incl University Hospital of North Norway) Open access funding provided by UiT The Arctic University of Norway (incl University Hospital of North Norway) This work was financially supported by the Research Council of Norway (RCN), through its Centre for Research-based Innovation funding scheme (Visual Intelligence, Grant No. 309439), and Consortium Partners. The work was further partially funded by RCN FRIPRO Grant No. 315029, RCN IKTPLUSS Grant No. 303514, and the UiT Thematic Initiative “Data-Driven Health Technology”.

Data Availability The data that supports these findings are all publicly available in online repositories. PASCAL VOC (Everingham et al., 2009): <http://host.robots.ox.ac.uk/pascal/VOC/index.html> MS COCO (Lin et al., 2014): <https://cocodataset.org/#home>. ImageNet (Deng et al., 2009): <https://www.image-net.org/>. MNIST: <http://yann.lecun.com/exdb/mnist/>.

Declarations

Conflict of interest The authors have no competing interests to declare that are relevant to the content of this article.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

Appendix A

The field of XAI have grown significantly in recent years, with a recent review article listing 46 different XAI methods (Samek et al., 2021). Nevertheless, the majority are designed with classification tasks in mind, and it is not straight forward to adapt these methods for the task of explaining representations. The difficulty stems from handling outputs in the form of vectors instead of scalars, since many method rely on propagating the scalar score from the output to the input or in other ways exploit the classification score for generating an explanation. For instance, surrogate models are a popular family of explainability methods, with LIME being the most well known approach. Molnar (2022) list the following steps for constructing relevance maps using surrogate models:

1. Select a dataset X .
2. For the selected dataset X , get the predictions of the black box model.
3. Select an interpretable model type.
4. Train the interpretable model on the dataset X and its predictions.
5. Measure how well the surrogate model replicates the predictions of the black box model.
6. Interpret the surrogate model.

Notice that the prediction of the model and labels are required in this procedure, neither of which are available in our representation learning scenario. Therefore, it is not obvious how such a procedure could be adapted to the unsupervised representation learning setting.

Another popular family of explanation methods is attribution-based explanation methods, where Layer-wise relevance propagation (LRP) is one of the most popular techniques. There exists several variants of LRP, but the core computation is to distribute the relevance of a neuron from the output layer back to the input. For instance, to explain the prediction for an input for a given class, the output of the class-neuron can be propagated back to the input. However, in the representation setting, the output is not a scalar, but typically a vector representation of the image. Therefore, each element of the vector must be propagated back, which would give an explanation for each of the elements. These explanations could be aggregated together as described in Sect. 4.2 of the main manuscript. However, this introduced a lot of noise in the explanation, since many of the elements that are being explained will not contain relevant information about the input. Both of these examples illustrate how it is challenging or not possible to trivially extend current explanation methods the context of representation learning explainability.

Appendix B

We compare the cosine similarity with the Euclidean distance for computing the similarity between the masked and unmasked representation. The results based on 100 samples are shown in Table 6 and show how the cosine similarity gives better performance. This suggests that angular information is more important in the context of representation learning explainability, which could be a result of angular information encoding semantically relevant information in neural networks Liu et al. (2021).

Appendix C

An alternative approach for creating the random variable $\bar{\mathbf{h}}$ is the following:

$$\bar{\mathbf{h}} = f(\mathbf{X} \odot \mathbf{M} - \mathbf{D}(1 - \mathbf{M})), \quad (\text{C1})$$

where each element of \mathbf{D} follows $N(\mu_{x_{ij}}, \sigma_{x_{ij}})$. The mean $\mu_{x_{ij}}$ and standard deviation $\sigma_{x_{ij}}$ is estimated by averaging across all samples in the data. Such a strategy could avoid potential distribution shifts that might occur when zeroing out large parts of the image, but also required determining the mean and variance of the data distribution.

Table 7 displays localisation scores the two masking strategies outlines in Sect. 3.1, namely zero masking or insertion of normally distributed noise. While there is some variation in the results, masking out with zeros provide the highest performance overall.

Appendix D

Figure 10 shows alternative strategies for masking out part of the input. One alternative is to apply Bernoulli noise to the input, which is equivalent to using Dropout (Srivastava et al., 2014) on the input. However, this does not introduce noise with spatial awareness, and therefore results in failing to explain the representation of the image. Another option is to drop regions of the input, such that objects could be fully or partially removed from the input. This could be achieved using the DropBlock algorithm (Ghiasi et al., 2018). However, this requires tuning the size of the mask on the input, which will be highly dependent on the objects present in the image. Such a per-image tuning would be impractical in most scenarios. Table 8 displays localisation scores for the Block drop strategy with different patch sizes. It shows how the patch-based strategy can both improve and decrease the performance of RELAX.

Appendix E

In this section we present the proofs for all theorems in the main paper.

Proof of Theorem 3.1

Proof First, let the Bounded difference assumption be defined as follows:

Definition 11.1 (Bounded difference assumption) Let A be some set and $f : A^N \rightarrow \mathbb{R}$. The function f satisfies the bounded differences assumption if there exists real numbers

$c_1, \dots, c_N \geq 0$ so that for all $i = 1, \dots, N$,

$$\sup_{x_1, \dots, x_N, x'_i \in A} |f(x_1, \dots, x_N, x'_i) - f(x_1, \dots, x_N, x_i)| \quad (E1)$$

Lemma 11.2 (McDiarmid’s inequality) *Let X_1, \dots, X_N be arbitrary independent random variables on set A and $f : A^N \rightarrow \mathbb{R}$ satisfies the bounded difference assumption. Then,*

We then have the following lemmas:

Table 6 Pointing game, top k, and relevance rank scores in percentages and averaged over 3 runs

Scores	Methods	Supervised		SimCLR		SwAV	
		COCO	VOC	COCO	VOC	COCO	VOC
Pointing game	RELAX	72.7 ± 1.5	85.0 ± 0.0	71.3 ± 3.1	88.7 ± 0.6	64.0 ± 3.0	87.0 ± 0.0
	RELAX (E)	20.0 ± 3.0	31.7 ± 1.5	24.3 ± 7.0	34.7 ± 8.6	22.3 ± 9.7	34.0 ± 1.0
	U-RELAX	72.3 ± 2.3	85.0 ± 0.0	71.0 ± 3.6	88.3 ± 1.2	64.0 ± 3.6	86.7 ± 0.6
	U-RELAX (E)	18.3 ± 2.9	28.3 ± 1.5	33.0 ± 4.4	41.0 ± 6.0	27.3 ± 5.0	41.7 ± 3.8
Top k	RELAX	73.0 ± 0.6	84.6 ± 0.1	71.5 ± 0.8	88.2 ± 0.2	65.8 ± 0.6	86.8 ± 0.2
	RELAX (E)	19.3 ± 0.5	30.3 ± 2.4	21.2 ± 4.6	30.3 ± 3.9	19.2 ± 8.4	29.6 ± 1.1
	U-RELAX	72.2 ± 0.7	84.6 ± 0.1	70.0 ± 1.5	87.3 ± 0.2	65.3 ± 0.8	86.3 ± 0.3
	U-RELAX (E)	18.5 ± 1.9	30.2 ± 3.5	27.4 ± 4.3	41.1 ± 5.1	24.2 ± 6.2	41.0 ± 3.0
Relevance rank	RELAX	58.2 ± 0.1	72.1 ± 0.4	53.8 ± 0.2	70.5 ± 0.3	53.1 ± 0.1	72.4 ± 0.3
	RELAX (E)	22.6 ± 0.6	34.5 ± 1.1	29.1 ± 2.6	43.5 ± 1.9	26.8 ± 3.2	43.6 ± 2.1
	U-RELAX	52.9 ± 0.2	66.8 ± 0.1	46.9 ± 0.4	59.9 ± 0.3	49.2 ± 0.3	65.1 ± 0.5
	U-RELAX (E)	26.9 ± 0.4	43.4 ± 0.9	31.5 ± 1.9	45.2 ± 0.8	31.5 ± 1.7	45.9 ± 0.4

Scores are computed using the cosine similarity and the Euclidean distance (indicated by (E))

Table 7 Evaluation of zero versus noise masking strategy in terms of pointing game, top k, and relevance rank scores in percentages and averaged over 3 runs

Scores	Methods	Supervised		SimCLR		SwAV	
		COCO	VOC	COCO	VOC	COCO	VOC
Pointing game	RELAX (zeros)	72.6 ± 0.1	86.6 ± 0.2	68.7 ± 0.3	85.2 ± 0.3	67.8 ± 0.2	84.7 ± 0.2
	RELAX (noise)	72.0 ± 0.5	86.0 ± 0.3	66.6 ± 0.1	84.3 ± 0.7	67.7 ± 0.5	85.1 ± 0.3
Top k	RELAX (zeros)	72.8 ± 0.4	86.9 ± 0.1	69.0 ± 0.3	85.6 ± 0.2	68.1 ± 0.4	85.1 ± 0.2
	RELAX (noise)	72.4 ± 0.4	86.5 ± 0.1	66.0 ± 0.3	84.2 ± 0.2	68.2 ± 0.3	85.3 ± 0.2
Relevance rank	RELAX (zeros)	56.4 ± 0.0	70.2 ± 0.1	54.2 ± 0.2	69.8 ± 0.1	52.4 ± 0.1	69.1 ± 0.0
	RELAX (noise)	56.7 ± 0.0	70.1 ± 0.1	53.5 ± 0.1	68.5 ± 0.0	52.8 ± 0.1	69.2 ± 0.0

Higher is better and bold numbers highlight the top performance. Results indicate that zero masking provides the best performance

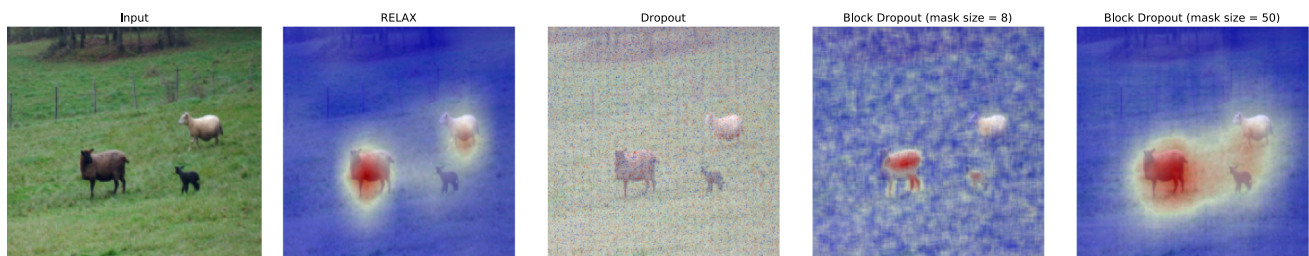


Fig. 10 Comparison of different masking strategies. Leftmost image shows input, and second to left is the RELAX explanations with the masking presented in the main paper. The center image is with Bernoulli-noise (Dropout) directly on the input, and the remaining two

images are with Block Dropout with different block size. The example illustrates that other masking strategies either fail completely, or require per-image parameter tuning, which is impractical in most scenarios

for all $t > 0$

$$P(|f(X_1, \dots, X_N) - E[f(X_1, \dots, X_N)]| \geq t) \leq 2e^{-\frac{2t^2}{\sum_{n=1}^N c_n^2}} \tag{E2}$$

Proof See McDiarmid (1989).

Lemma 11.3 Let X_1, \dots, X_N and f be defined as in Lemma 11.2, then if each X_n satisfies $X_n \in (a_n, b_n)$ and $f(X_1, \dots, X_N) = \sum_{n=1}^N X_n$, then $c_n = b_n - a_n$.

Proof See McDiarmid (1989).

We are now ready to prove the theorem. First, let

$$X_n = \frac{s(\mathbf{h}, \bar{\mathbf{h}}_n)M_{ij}(n)}{N}, \tag{E3}$$

and

$$f(X_1, \dots, X_n) = \sum_{n=1}^N X_n. \tag{E4}$$

Since $s(\cdot, \cdot)$ is bounded in $(0, 1)$ (we use the cosine similarity between vectors with non-negative elements (ReLU outputs)), we have $a_n = 0$ and $b_n = 1/N$, which gives $c_n = 1/N$ by Lemma 11.3.

Now, observe that

$$f(X_1, \dots, X_n) = \frac{1}{N} \sum_{n=1}^N s(\mathbf{h}, \bar{\mathbf{h}}_n)M_{ij}(n) = \bar{R}_{ij}. \tag{E5}$$

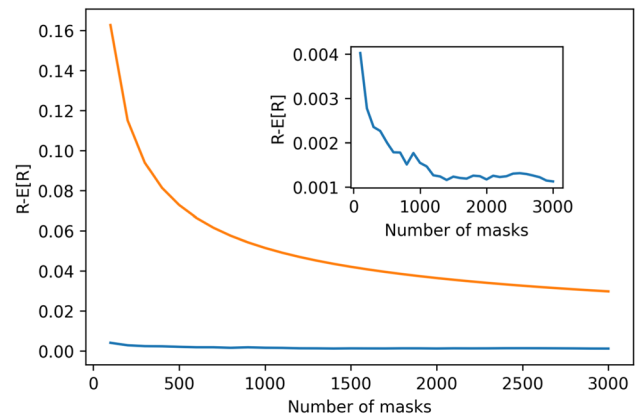


Fig. 11 Empirical evaluation of the derived bound for the number of masks necessary for low estimation error. We calculate the absolute error as the number of masks increase, average over 10 randomly samples images from the PASCAL VOC dataset. To obtain a value for R_{ij} , we use 10,000 masks and average over 10 runs for a single sample. Results indicate that the estimation error is much lower than the predicted bound

Combining Lemmas 11.2 and 11.3 then gives

$$P(|\bar{R}_{ij} - R_{ij}| \geq t) \leq 2e^{-\frac{2t^2}{\sum_{n=1}^N (1/N)^2}} \tag{E6}$$

for all $t > 0$. Inserting $N = -\ln(\delta/2)/2t^2$ gives

$$P(|\bar{R}_{ij} - R_{ij}| \geq t) \leq 2e^{-\frac{2t^2}{\sum_{n=1}^N (1/N)^2}} \tag{E7}$$

$$= 2e^{-2t^2 \left(-\frac{\ln(\delta/2)}{2t^2} \right)} \tag{E8}$$

$$= 2e^{\ln(\delta/2)} \tag{E9}$$

$$= \delta, \tag{E10}$$

Table 8 Pointing game, top k, and relevance rank scores in percentages and averaged over 3 runs

Scores	Methods	Supervised		SimCLR		SwAV	
		COCO	VOC	COCO	VOC	COCO	VOC
Pointing game	RELAX (16)	78.0±1.0	84.0±0.6	60.8±0.8	77.3±0.7	66.9±0.5	78.8±0.4
	RELAX (20)	79.1±0.5	87.7±0.3	70.2±0.6	83.2±0.4	70.7±0.5	84.6±0.4
	RELAX (50)	78.9±0.8	88.1±0.2	69.9±0.1	84.2±0.4	67.0±1.2	85.3±1.0
	RELAX	72.6±0.1	86.6±0.2	68.7±0.3	85.2±0.3	67.8±0.2	84.7±0.2
Top k	RELAX (16)	75.9±0.1	83.9±0.0	60.6±0.6	76.3±0.3	65.6±0.1	78.6±0.2
	RELAX (20)	78.7±0.1	87.3±0.2	69.2±0.3	82.7±0.4	69.3±0.4	84.0±0.1
	RELAX (50)	77.8±0.4	87.7±0.1	69.0±0.6	84.2±0.2	66.6±0.8	85.0±0.4
	RELAX	72.8±0.4	86.9±0.1	69.0±0.3	85.6±0.2	68.1±0.4	85.1±0.2
Relevance rank	RELAX (16)	52.0±0.0	60.0±0.0	43.9±0.0	56.8±0.0	46.0±0.1	56.7±0.0
	RELAX (20)	57.9±0.0	67.3±0.0	49.5±0.1	64.2±0.1	49.8±0.1	64.0±0.1
	RELAX (50)	58.4±0.1	72.5±0.0	51.8±0.1	69.2±0.0	52.5±0.2	69.8±0.0
	RELAX	56.4±0.0	70.2±0.1	54.2±0.2	69.8±0.1	52.4±0.1	69.1±0.0

Scores computed with the masking strategy outline in the main manuscript and using the DropBlock algorithm with different patch size (indicated in parenthesis)

Table 9 Pointing game, top k, and relevance rank scores in percentages and averaged over 3 runs for vision transformers

Scores	Methods	ViT		MAE		
		COCO	VOC	COCO	VOC	
Pointing game	Saliency	37.1 ± 0.0	46.2 ± 0.0	20.8 ± 0.0	28.3 ± 0.0	
	Smooth Saliency	37.1 ± 0.0	46.2 ± 0.0	35.1 ± 1.1	51.6 ± 1.2	
	Guided Saliency	36.4 ± 0.0	45.3 ± 0.0	19.0 ± 0.0	29.6 ± 0.0	
	Integrated Gradients	52.2 ± 0.5	65.2 ± 2.5	14.9 ± 0.0	26.0 ± 0.0	
	Grad CAM	37.7 ± 0.0	48.2 ± 0.0	42.9 ± 0.0	63.3 ± 0.0	
	RELAX	69.1 ± 0.3	85.2 ± 0.5	60.0 ± 1.4	74.9 ± 0.9	
	U-RELAX	68.7 ± 0.3	84.9 ± 0.5	60.0 ± 0.8	75.3 ± 0.8	
	Top k	Saliency	36.8 ± 0.0	48.5 ± 0.0	24.0 ± 0.0	39.8 ± 0.0
Top k	Smooth Saliency	36.8 ± 0.0	48.5 ± 0.0	37.6 ± 0.0	56.5 ± 0.1	
	Guided Saliency	36.3 ± 0.0	46.2 ± 0.0	24.6 ± 0.0	40.3 ± 0.0	
	Integrated Gradients	47.8 ± 0.1	61.5 ± 0.3	23.4 ± 0.0	41.0 ± 0.0	
	Grad CAM	35.7 ± 0.0	47.8 ± 0.0	39.2 ± 0.0	55.8 ± 0.0	
	RELAX	70.0 ± 0.1	84.7 ± 0.4	60.2 ± 0.5	74.6 ± 0.6	
	U-RELAX	69.6 ± 0.0	84.4 ± 0.4	60.1 ± 0.4	74.7 ± 0.6	
	Relevance rank	Saliency	34.7 ± 0.0	48.9 ± 0.0	29.6 ± 0.0	44.5 ± 0.0
	Relevance rank	Smooth Saliency	34.7 ± 0.0	48.9 ± 0.0	34.4 ± 0.0	49.1 ± 0.0
Guided Saliency		33.9 ± 0.0	47.9 ± 0.0	29.7 ± 0.0	44.6 ± 0.0	
Integrated Gradients		36.7 ± 0.0	50.2 ± 0.0	30.5 ± 0.0	45.6 ± 0.0	
Grad CAM		33.1 ± 0.0	47.4 ± 0.0	32.8 ± 0.0	47.3 ± 0.0	
RELAX		55.6 ± 0.0	70.0 ± 0.1	51.6 ± 0.1	66.6 ± 0.1	
U-RELAX		52.7 ± 0.1	65.3 ± 0.1	51.0 ± 0.2	63.9 ± 0.2	

Higher is better and bold numbers highlight the top performance. Results show that our method improves on the baseline across all scores

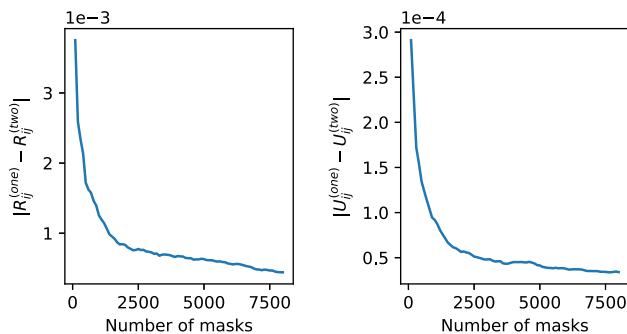


Fig. 12 Absolute error of one-pass versus two-pass version of RELAX for importance (leftmost figure) and uncertainty (rightmost figure), averaged over 50 images from the VOC dataset. The figure shows how the difference between the versions is small for both the importance and uncertainty estimates

which concludes our proof.

In Fig. 11 we show an empirical validation the bound. We calculate the absolute error as the number of masks increase, averaged over 10 randomly sampled images from the PAS-CAL VOC dataset. To obtain a value for R_{ij} , we use 10,000 masks and average over 10 runs for a single sample. The results indicate that the true error is much lower than the proposed bound, which we attribute to setting $a_n = 0$. While

Table 10 Monotonicity scores averaged over 3 runs for vision transformers

Scores	Methods	ViT	MAE
Monotonicity	Saliency	14.2 ± 0.0	3.6 ± 0.3
	Smooth Saliency	14.6 ± 0.6	3.4 ± 0.2
	Guided Saliency	4.7 ± 0.0	3.9 ± 0.0
	Integrated Gradients	12.8 ± 0.3	1.9 ± 0.4
	Grad CAM	4.5 ± 0.0	4.6 ± 0.0
	RELAX	17.7 ± 0.4	15.4 ± 0.3
	U-RELAX	24.9 ± 0.3	23.7 ± 0.5

Higher is better and bold numbers highlight the top performance. Results show that our method improves on the baseline

it is possible to obtain a similarity of 0, it is highly unlikely since our masking strategy never removes all information in an image.

Proof of Theorem 3.2

Proof Since $s(\cdot, \cdot)$ is a valid Mercer kernel, we can write $s(\mathbf{h}, \bar{\mathbf{h}}_n) = \langle \phi(\mathbf{h}), \phi(\bar{\mathbf{h}}_n) \rangle_{\mathcal{H}}$. This gives

$$\bar{R}_{ij} = \frac{1}{N} \sum_{n=1}^N \langle \phi(\mathbf{h}), \phi(\bar{\mathbf{h}}_n) \rangle_{\mathcal{H}} M_{ij}(n) \tag{E11}$$

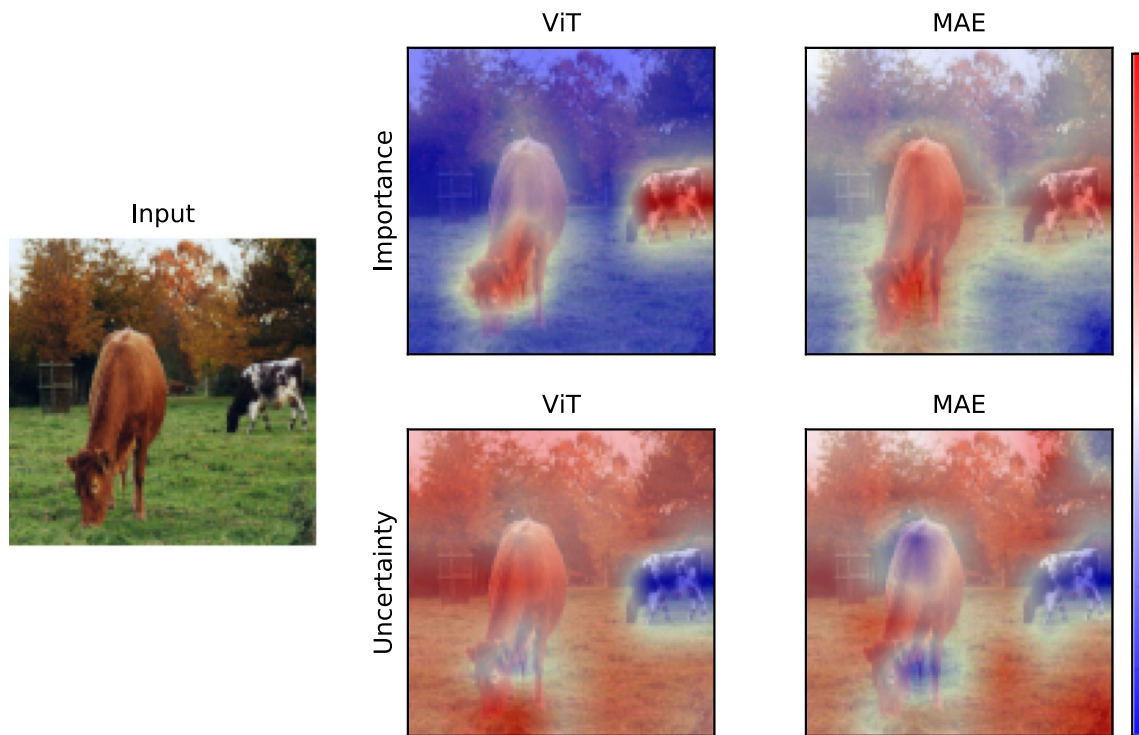


Fig. 13 The figure shows the RELAX explanation and its uncertainty for the representation of the leftmost image with two recent vision transformer architectures used as feature extractors. The first row displays the explanations for the representation and the second row shows the uncertainty associated with the different explanations. Red indicates high values and blue indicates low values. In this example, the ViT

that is trained with label information produces a more focused explanation, while MAE focus on larger parts of the image. Interestingly, both models are more certain about importance of the rightmost cow in the background. The image is taken from VOC Everingham et al. (2009) (Color figure online)

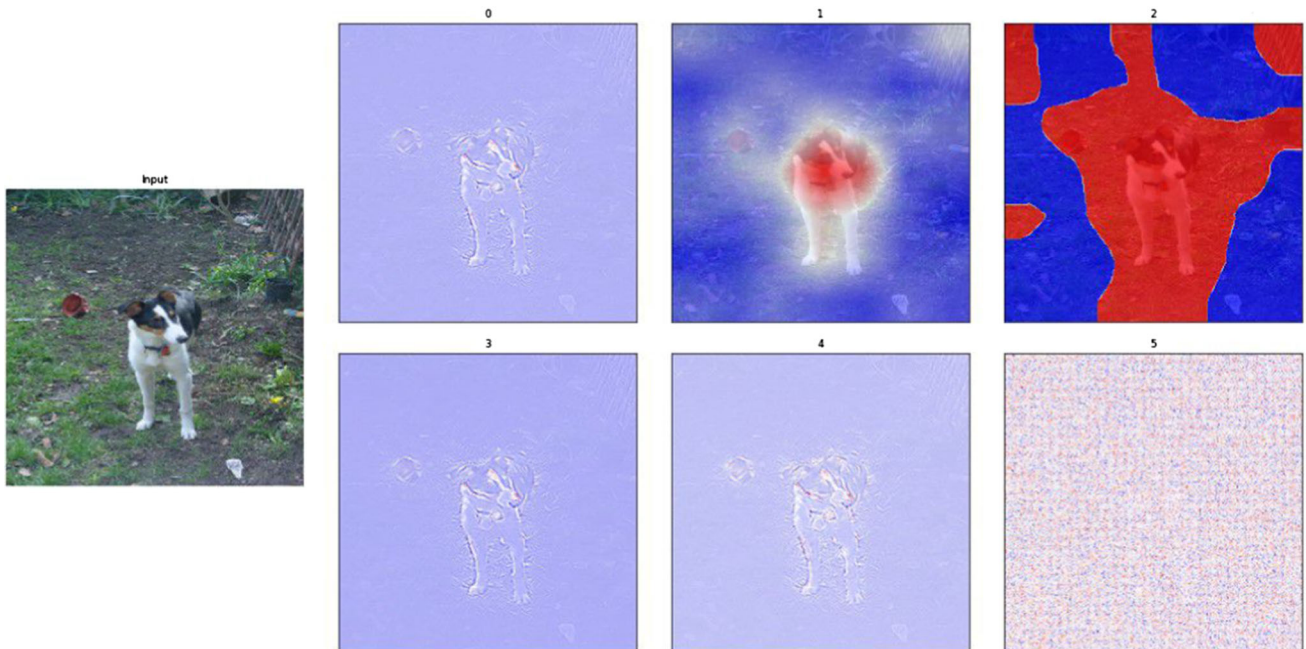


Fig. 14 Example from the human evaluation experiment. Participants were asked to select which explanation they preferred out of the 6 alternatives. For each of the images, the explanations were shuffled in a

random order. One of the explanations for each image was randomly sampled from random noise, in order to assess if any participants would select a nonsensical explanation

Table 11 Monotonicity scores averaged over 3 runs with a VGG13 classifier

Scores	Methods	Supervised	SimCLR	SwAV
monotonicity	Saliency	11.2 ± 0.0	12.7 ± 0.5	12.8 ± 0.3
	Smooth Saliency	14.4 ± 0.7	12.6 ± 0.6	12.3 ± 0.6
	Guided Saliency	13.2 ± 0.4	12.5 ± 0.1	11.6 ± 0.9
	Integrated Gradients	10.9 ± 0.2	6.3 ± 0.5	8.3 ± 0.4
	Grad CAM	19.2 ± 0.4	16.2 ± 0.1	15.8 ± 0.3
	RELAX	24.0 ± 0.4	28.3 ± 0.7	24.0 ± 0.0
	U-RELAX	29.3 ± 0.4	26.1 ± 0.7	18.7 ± 0.8

Higher is better and bold numbers highlight the top performance. Results show that our method improves on the baseline

$$= \langle \phi(\mathbf{h}), \frac{1}{N} \sum_{n=1}^N \phi(\bar{\mathbf{h}}_n) M_{ij}(n) \rangle_{\mathcal{H}} \tag{E12}$$

by the bilinearity of the inner product on \mathcal{H} .

Proof of Theorem 3.3

Proof Observe that

$$\bar{R}_{ij} \cdot \frac{N}{\sum_{n'=1}^N M_{ij}(n')} \tag{E13}$$

$$= \frac{N}{\sum_{n'=1}^N M_{ij}(n')} \cdot \frac{1}{N} \sum_{n=1}^N s(\mathbf{h}, \bar{\mathbf{h}}_n) M_{ij}(n) \tag{E14}$$

$$= \frac{1}{\sum_{n'=1}^N M_{ij}(n')} \sum_{n=1}^N s(\mathbf{h}, \bar{\mathbf{h}}_n) M_{ij}(n) \tag{E15}$$

$$= p_{ij}(\mathbf{h}) \tag{E16}$$

\bar{R}_{ij} is therefore proportional to $p_{ij}(\mathbf{h})$.

Appendix F

We investigate the potential differences between the one-pass and two-pass version of RELAX. For a given image, we calculate the absolute error between the one-pass and two-pass estimates for different number of masks. The results are shown in Fig. 12 and illustrate that the difference between the two methods is very small, particularly as the number of masks increases. However, since the one-pass version computes both the importance and uncertainty in one pass through the data, it requires only half the number of masks compared to the two pass version, thus increasing the computational efficiency of RELAX.

Appendix G

We use RELAX to investigate the representation of recent vision transformer architectures. Specifically, we investigate the Vision Transformer (ViT) (Dosovitskiy et al., 2021) and the Masked Autoencoder (MAE) (He et al., 2022), with the same setting as in the main manuscript. Tables 9 and 10 shows the localisation and faithfulness scores, and Fig. 13 shows a qualitative comparison of the RELAX explanation. Again, the scores show how RELAX outperforms all other methods. Also, the representation created by the ViT model give slightly superior performance to the MAE model, which again illustrates that label information is useful to create high quality representations.

Appendix H

Table 11 displays the monotonicity score with a VGG13 classifier. Similarly as with an Alexnet classifier, the proposed methods improves upon the baseline.

Appendix I

The user study in the main manuscript was conducted by having a group of participants select among competing explanations for a random selection of images from the PASCAL VOC dataset. The group of participants consisted of men and women, where some had knowledge of machine learning and other were uneducated. None of the participants have been involved in the development of this work. Figure 14 displays an example from the study. The participants were shown an image with 6 competing explanations, and asked to chose which one they preferred. To determine which explanation each participant judged to be the “best”, they were told to ask themselves the following questions:

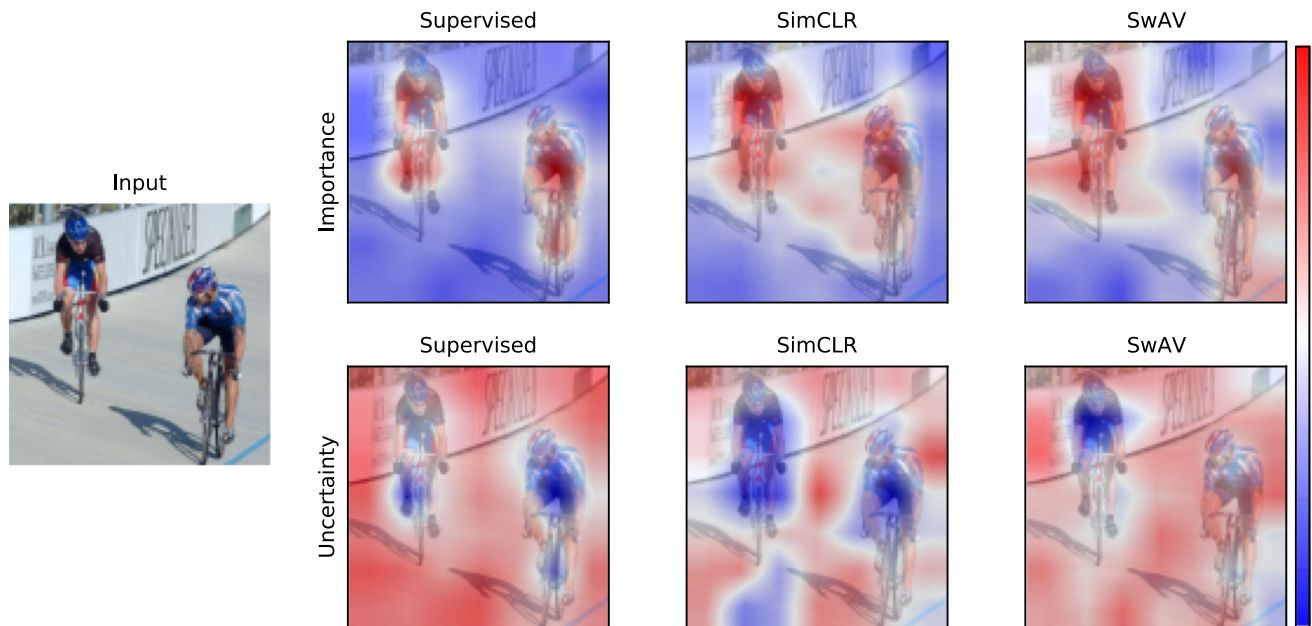


Fig. 15 Example from the VOC dataset

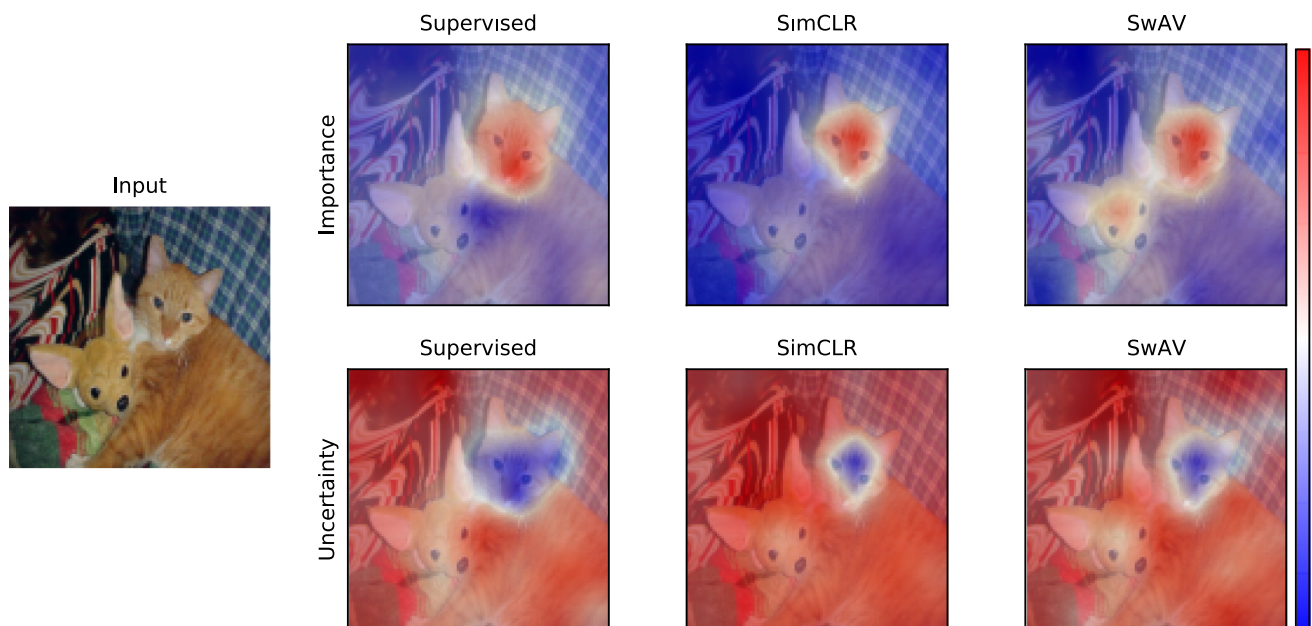


Fig. 16 Example from the COCO dataset

“Which of these explanations agree the most with how you would explain the important content in the given image?”

For each image, the explanations were shuffled randomly. The participants were shown 10 images, and asked to only pick on explanation. Overall, 13 people participated in the study.

There are several limitations. Both the number of images and the number of participants could have been greater. The

participants had to chose one explanation, when in some cases they might have wanted to select none or more explanations. Also, the images could have been selected from other datasets. There are also potential biases with the study. Most participants are from one country and from a limited age segment. Lastly, we did not control the type of screen that participants performed their evaluation on, which could also have an undesirable affect.

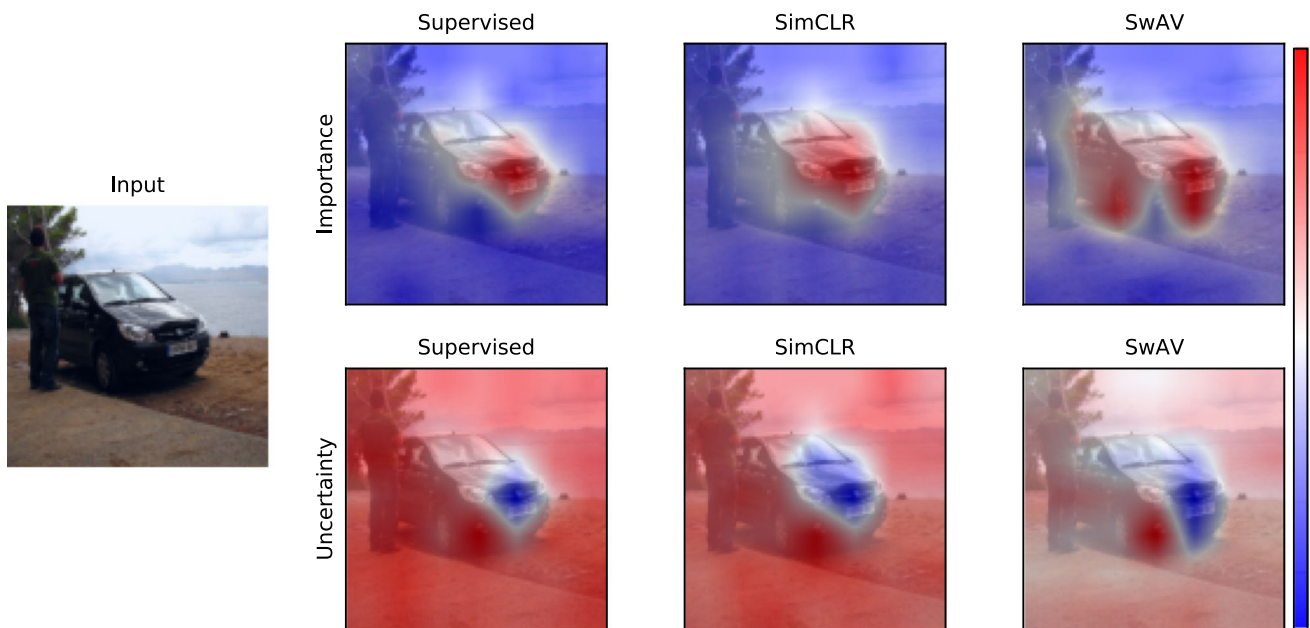


Fig. 17 Example from the VOC dataset

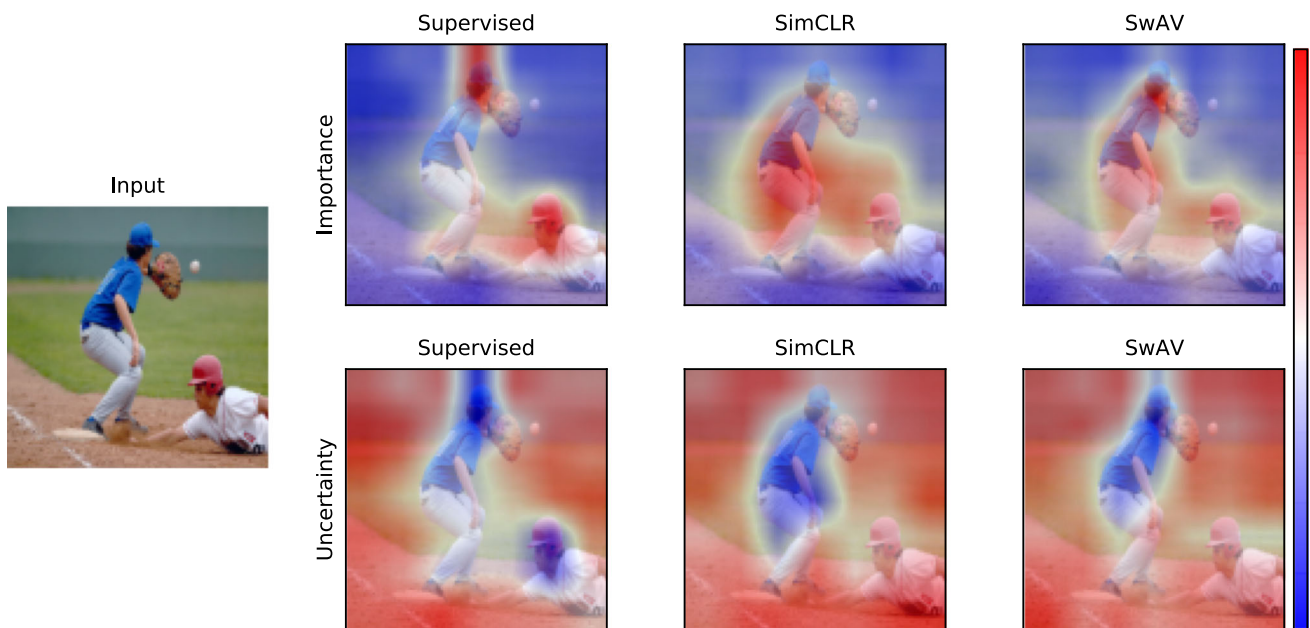


Fig. 18 Example from the COCO dataset

Appendix J

This section presents additional qualitative results. Figures 15, 16, 17, 18, 19, 20, 21, 22, 23 and 24 displays examples of explanations and their associated uncertainty, provided by RELAX, for images from the VOC and COCO dataset. Figure 15 displays an example where all feature extractors agree in terms of importance, but the degree of

uncertainty varies. Figure 16 shows an example where only SwAV highlight both objects as important for the representation. Similarly, Fig. 17 displays an example where only SwAV is considering both the person and the car as important for the representation. Figures 17, 18, 19, 20, 21, 22, 23 and 24 shows similar examples where RELAX provides insights into the different feature extractors.

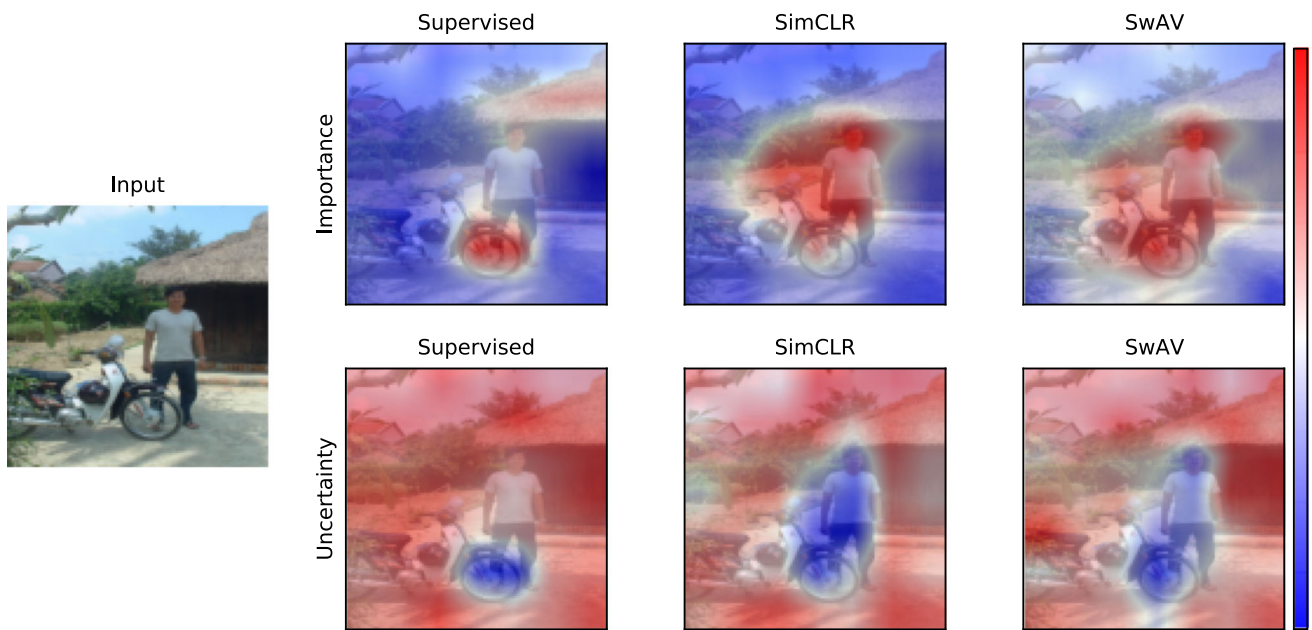


Fig. 19 Example from the VOC dataset

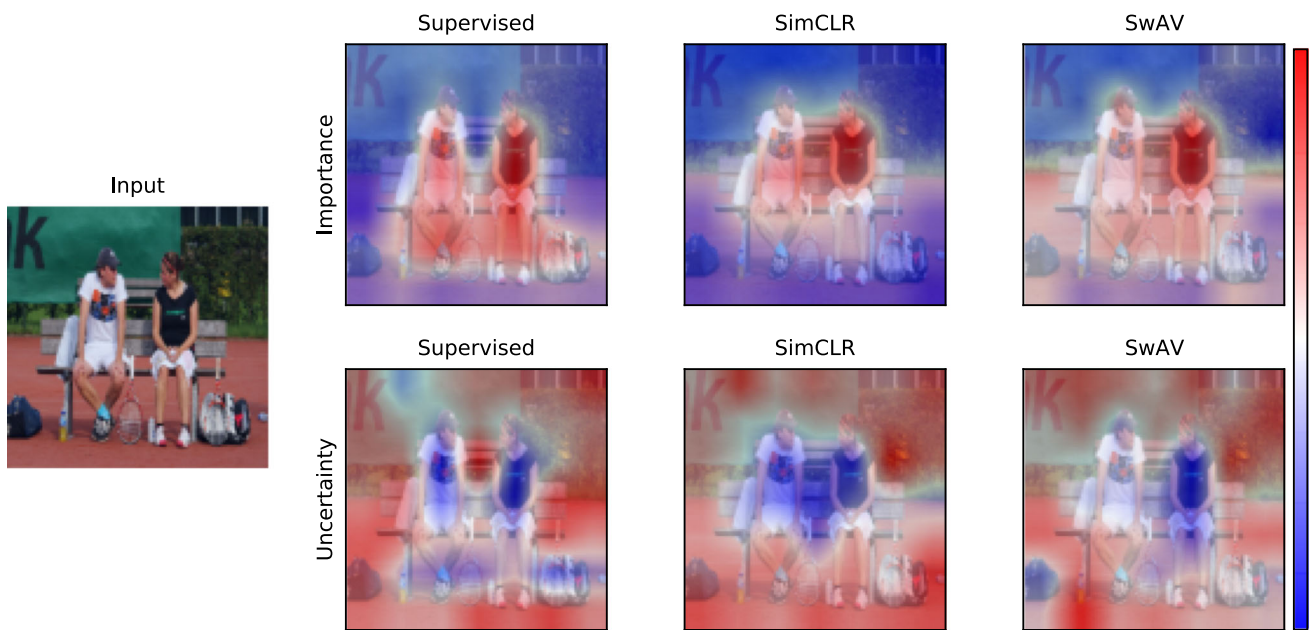


Fig. 20 Example from the COCO dataset

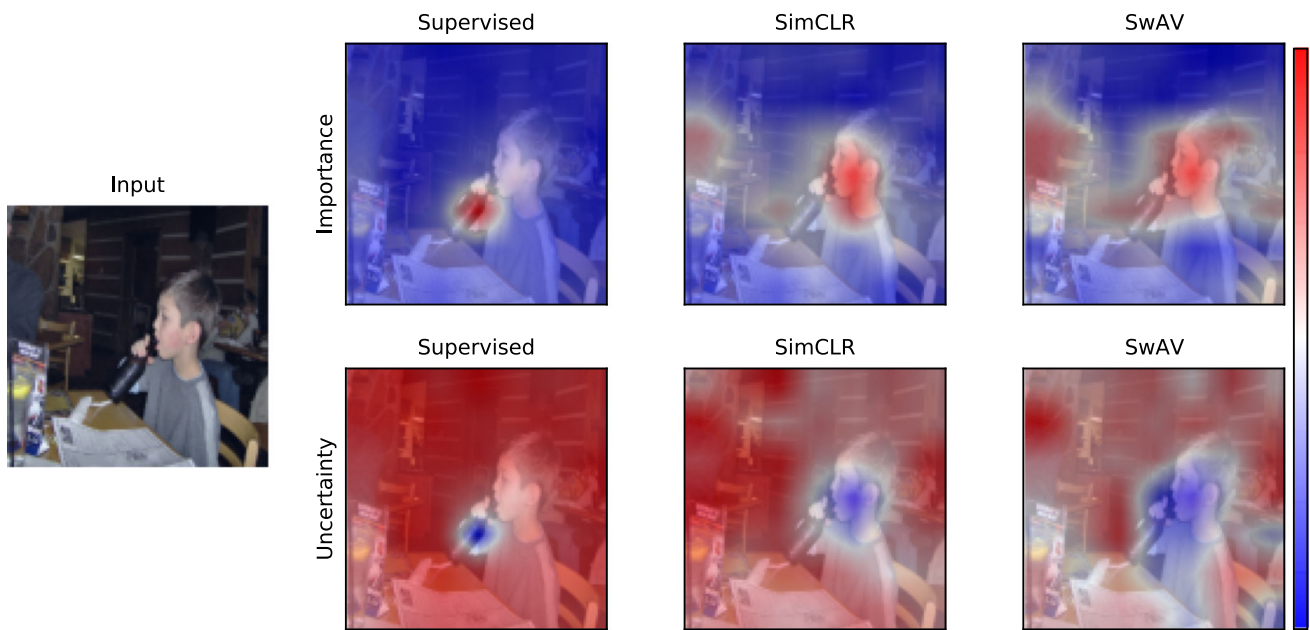


Fig. 21 Example from the VOC dataset

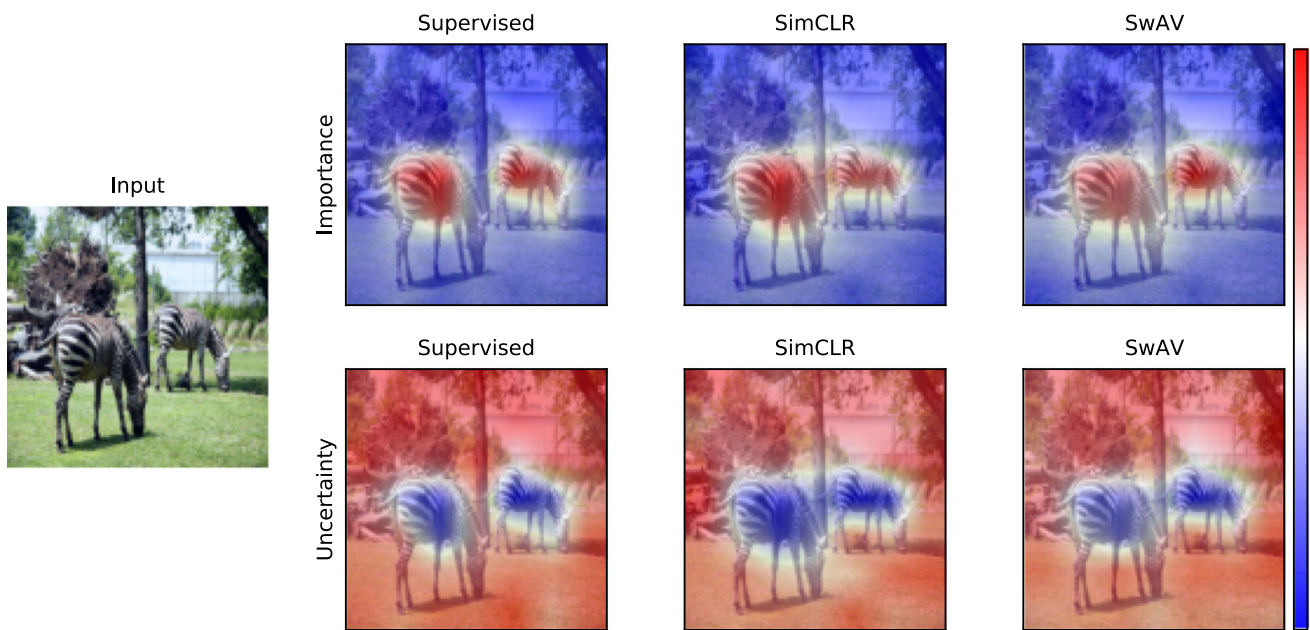


Fig. 22 Example from the COCO dataset

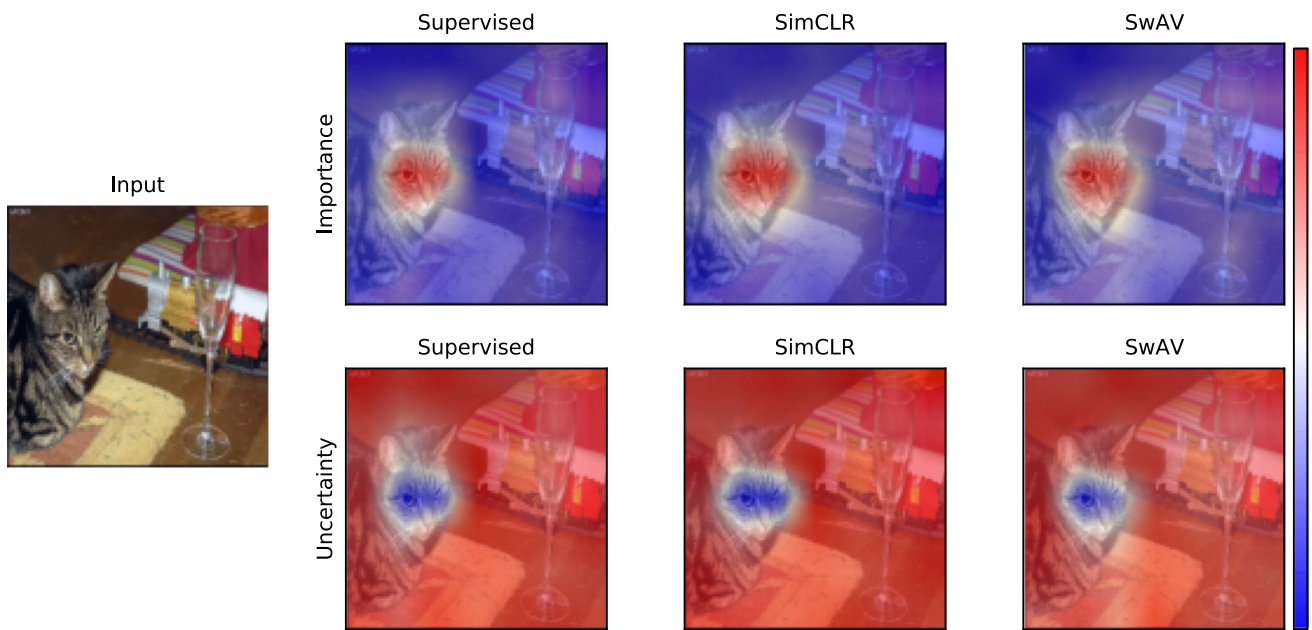


Fig. 23 Example from the VOC dataset

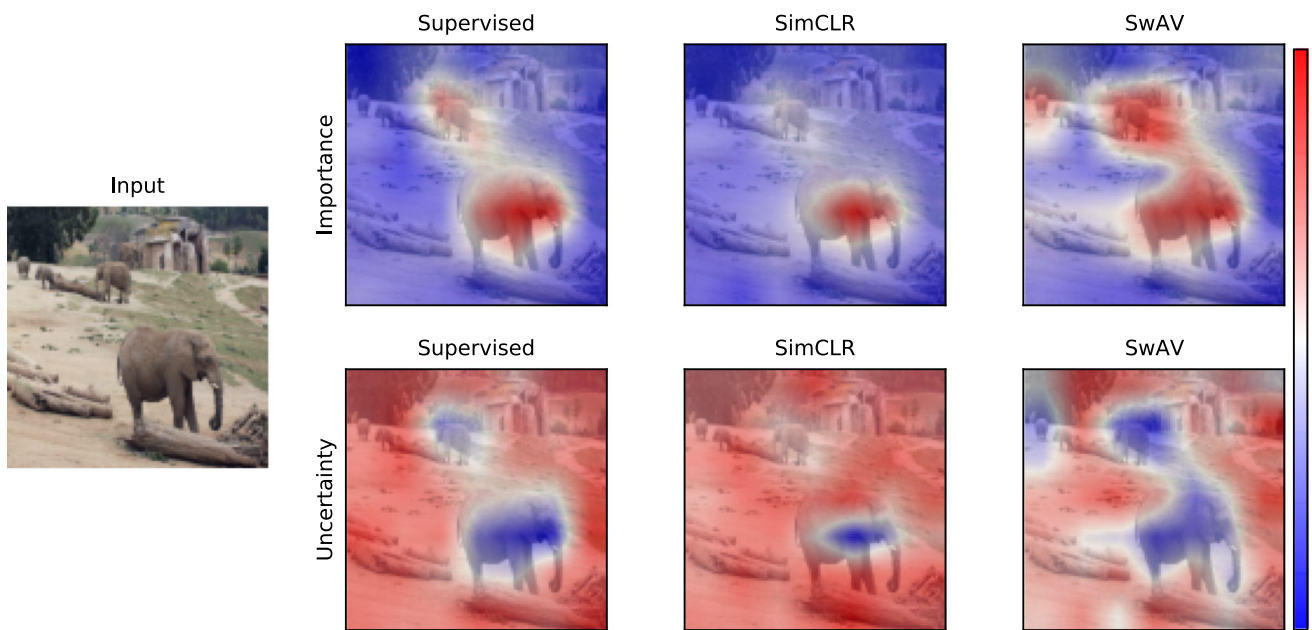


Fig. 24 Example from the COCO dataset

References

- Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., & Kim, B. (2018). Sanity checks for saliency maps. In *Advances in neural information processing systems*. Curran Associates, Inc.
- Alvarez-Melis, D., & Jaakkola, T. S. (2018). Towards robust interpretability with self-explaining neural networks. In *Proceedings of the 32nd international conference on neural information processing systems* (pp. 7786–7795). Curran Associates Inc., Red Hook, NY, USA, NIPS'18.
- Antoran, J., Bhatt, U., Adel, T., Weller, A., & Hernandez-Lobato, J. M. (2020). Getting a clue: A method for explaining uncertainty e (2021). Getting a clue: A method for explaining uncertainty estimates. In *International conference on learning representations*.
- Arras, L., Osman, A., & Samek, W. (2022). Clevr-xai: A benchmark dataset for the ground truth evaluation of neural network explanations. *Information Fusion*, 81, 14–40. <https://doi.org/10.1016/j.inffus.2021.11.008>
- Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K. R., & Samek, W. (2015). On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE*, 10(7), e0130140. <https://doi.org/10.1371/journal.pone.0130140>
- Bau, D., Zhou, B., Khosla, A., Oliva, A., & Torralba, A. (2017). Network dissection: Quantifying interpretability of deep visual representations. In *IEEE computer vision and pattern recognition*.
- Bykov, K., Höhne, M. M. C., Müller, K. R., Nakajima, S., & Kloft, M. (2020) How much can I trust you?—Quantifying uncertainties in explaining neural networks. CoRR [arXiv:2006.09000](https://arxiv.org/abs/2006.09000)
- Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., & Joulin, A. (2020). Unsupervised learning of visual features by contrasting cluster assignments. In *Advances in neural information processing systems* (pp. 9912–9924).
- Chen, C., Li, O., Tao, D., Barnett, A., Rudin, C., & Su, J. K. (2019). This looks like that: Deep learning for interpretable image recognition. In *International conference on neural information processing systems*.
- Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). A simple framework for contrastive learning of visual representations. In *International conference on machine learning* (pp. 1597–1607).
- Chen, X., & He, K. (2021). Exploring simple siamese representation learning. In *IEEE computer vision and pattern recognition* (pp. 15750–15758).
- Dalal, N., & Triggs, B. (2005). Histograms of oriented gradients for human detection. *IEEE Computer Vision and Pattern Recognition*, 1, 886–893. <https://doi.org/10.1109/CVPR.2005.177>
- Deng, J., Dong, W., Socher, R., Li, L. J., Li, K., & Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *IEEE Computer Vision and Pattern Recognition* (pp. 248–255). IEEE.
- Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. [arXiv:1702.08608](https://arxiv.org/abs/1702.08608)
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., & Houshy, N. (2021) An image is worth 16x16 words: Transformers for image recognition at scale. In *International conference on learning representations*. <https://openreview.net/forum?id=YicbFdNTTy>
- Everingham, M., Van Gool, L., Williams, C. K., Winn, J., & Zisserman, A. (2009). The pascal visual object classes (VOC) challenge. *International Journal of Computer Vision*, 88, 303–338. <https://doi.org/10.1007/s11263-009-0275-4>
- Falcon, W., & Cho, K. (2020). A framework for contrastive self-supervised learning and designing a new approach. [arXiv preprint arXiv:2009.00104](https://arxiv.org/abs/2009.00104).
- Fong, R., & Vedaldi, A. (2018). Net2vec: Quantifying and explaining how concepts are encoded by filters in deep neural networks. In *IEEE computer vision and pattern recognition* (pp. 8730–8738). <https://doi.org/10.1109/CVPR.2018.00910>
- Fong, R., Patrick, M., & Vedaldi, A. (2019). Understanding deep networks via extremal perturbations and smooth masks. In *IEEE International Conference on Computer Vision*.
- Fong, R. C., & Vedaldi, A. (2017). Interpretable explanations of black boxes by meaningful perturbation. In *IEEE international conference on computer vision* (pp. 3449–3457). <https://doi.org/10.1109/ICCV.2017.371>
- Gal, Y., & Ghahramani, Z. (2016). Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In *International conference on machine learning* (pp. 1050–1059).
- Ghiasi, G., Lin, T. Y., & Le, Q. V. (2018). Dropblock: A regularization method for convolutional networks. In *International conference on neural information processing systems* (pp. 10750–10760).
- He, K., Zhang, X., Ren, S., & Sun, J. (2016) Deep residual learning for image recognition. In *2016 CVPR* (pp. 770–778). <https://doi.org/10.1109/CVPR.2016.90>
- He, K., Fan, H., Wu, Y., Xie, S., & Girshick, R. (2020). Momentum contrast for unsupervised visual representation learning. In *IEEE computer vision and pattern recognition*.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., & Girshick, R. (2022). Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)* (pp. 16000–16009).
- Karimi, A. H., Barthe, G., Balle, B., & Valera, I. (2020). Model-agnostic counterfactual explanations for consequential decisions. In *International conference on artificial intelligence and statistics* (pp. 895–905).
- Kim, B., Wattenberg, M., Gilmer, J., Cai, C., Wexler, J., & Viegas, F. (2018). Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). In *International conference on machine learning* (pp. 2673–2682).
- Koh, P. W., & Liang, P. (2017). Understanding black-box predictions via influence functions. In *International conference on machine learning* (pp. 1885–1894).
- Kolek, S., Nguyen, D. A., Levie, R., Bruna, J., & Kutyniok, G. (2021). A rate-distortion framework for explaining black-box model decisions. [arXiv:2110.08252](https://arxiv.org/abs/2110.08252)
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In: F. Pereira, C. Burges, L. Bottou, et al. (Eds.), *Advances in Neural Information Processing Systems* (Vol 25). Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>
- Laina, I., Fong, R., & Vedaldi, A. (2020). Quantifying learnability and describability of visual concepts emerging in representation learning. In *Advances in neural information processing systems*.
- Leys, C., Ley, C., Klein, O., Bernard, P., & Licata, L. (2013). Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median. *Journal of Experimental Social Psychology*, 49(4), 764–766. <https://doi.org/10.1016/j.jesp.2013.03.013>
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P. & Lawrence Zitnick, C. (2014). Microsoft COCO: Common objects in context. In: *Computer Vision—ECCV 2014* (pp. 740–755). Springer International Publishing. https://doi.org/10.1007/978-3-319-10602-1_48
- Lin, Y., Gou, Y., Liu, Z., Li, B., Lv, J., & Peng, X. (2021). Completer: Incomplete multi-view clustering via contrastive prediction. In *IEEE computer vision and pattern recognition* (pp. 11174–11183).
- Liu, W., Lin, R., Liu, Z., Xiong, L., Scholkopf, B., & Weller, A. (2021). Learning with hyperspherical uniformity. In *Proceedings of the 24th international conference on artificial intelligence and statistics, Proceedings of machine learning research* (Vol.

- 130, pp. 1180–1188). PMLR. <http://proceedings.mlr.press/v130/liu21d.html>
- Losch, M., Fritz, M., & Schiele, B. (2021). Semantic bottlenecks: Quantifying and improving inspectability of deep representations. *International Journal of Computer Vision*, 129(11), 3136–3153. <https://doi.org/10.1007/s11263-021-01498-0>
- McCullagh, P., & Nelder, J. (1989). *Generalized linear models* (2nd ed.). Chapman & Hall.
- McDiarmid, C. (1989). *On the method of bounded difference* (pp. 148–188). Cambridge University Press. <https://doi.org/10.1017/CBO9781107359949.008>
- Mercer, J. (1909). Functions of positive and negative type, and their connection with the theory of integral equations. *Philosophical Transactions of the Royal Society, London*, 209, 415–446.
- Molnar, C. (2022). *Interpretable machine learning*. 2nd edn. <https://christophm.github.io/interpretable-ml-book>
- Mordvintsev, A., Olah, C., & Tyka, M. (2015). Inceptionism: Going deeper into neural networks. <https://research.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>
- Nguyen, A., & Martinez, M. R. (2020). On quantitative aspects of model interpretability. [arXiv:2007.07584](https://arxiv.org/abs/2007.07584)
- Nordhaug Myhre, J., Øvind Mikalsen, K., & Løkse, S. (2018). Robust clustering using a kNN mode seeking ensemble. *Pattern Recognition*, 76, 491–505. <https://doi.org/10.1016/j.patcog.2017.11.023>
- Parzen, E. (1962). On estimation of a probability density function and mode. *The Annals of Mathematical Statistics*, 33(3), 1065–1076. <https://doi.org/10.1214/aoms/1177704472>
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Benoit Steiner, L., Fang, J. B., & Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems* (pp. 8024–8035).
- Pedreschi, D., Giannotti, F., Guidotti, R., Monreale, A., & Ruggieri, S. (2019). Meaningful explanations of black box AI decision systems. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33, 9780–9784. <https://doi.org/10.1609/aaai.v33i01.33019780>
- Petsiuk, V., Das, A., & Saenko, K. (2018). Rise: Randomized input sampling for explanation of black-box models. In *Proceedings of the British machine vision conference*.
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “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* (pp. 1135–1144), San Francisco, CA, USA, August 13–17, 2016.
- Samek, W., Binder, A., Montavon, G., Lapuschkin, S., & Müller, K. R. (2017). Evaluating the visualization of what a deep neural network has learned. In *IEEE TNNLS* (pp. 2660–2673). <https://doi.org/10.1109/TNNLS.2016.2599820>
- Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J., & Müller, K. R. (2021). Explaining deep neural networks and beyond: A review of methods and applications. In *Proceedings of the IEEE* (pp. 247–278).
- Schulz, K., Sixt, L., Tombari, F., & Landgraf, T. (2020). Restricting the flow: Information bottlenecks for attribution. In *International conference on learning representations*.
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. In *2017 IEEE international conference on computer vision (ICCV)* (pp. 618–626). <https://doi.org/10.1109/ICCV.2017.74>
- Shi, T., Belkin, M., & Yu, B. (2009). Data spectroscopy: Eigenspaces of convolution operators and clustering. *The Annals of Statistics*, 37(6B), 3960–3984. <https://doi.org/10.1214/09-AOS700>
- Simonyan, K., & Zisserman, A. (2015). Very deep convolutional networks for large-scale image recognition. In *International conference on learning representations*.
- Smilkov, D., Thorat, N., Kim, B., Viegas, F., & Wattenberg, M. (2017). Smoothgrad: Removing noise by adding noise. In *International conference on machine learning visualization workshop*.
- Springenberg, J. T., Dosovitskiy, A., Brox, T., & Riedmiller, M. (2015). Striving for simplicity: The all convolutional net. In *ICLR Workshop*.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15, 1929–1958.
- Sundararajan, M., Taly, A., & Yan, Q. (2017). Axiomatic attribution for deep networks. In D. Precup, & Y. W. Teh (Eds.), *Proceedings of the 34th international conference on machine learning, ICML 2017, Sydney, NSW, Australia, 6–11 August 2017, Proceedings of Machine Learning Research* (Vol. 70, pp. 3319–3328). PMLR. <http://proceedings.mlr.press/v70/sundararajan17a.html>
- Teye, M., Azizpour, H., & Smith, K. (2018). Bayesian uncertainty estimation for batch normalized deep networks. In *International conference on machine learning* (pp. 4907–4916).
- Theodoridis, S., & Koutroumbas, K. (2009). *Pattern recognition* (4th ed.). Academic Press.
- Tonekaboni, S., Joshi, S., McCradden, M. D., & Goldenberg, A. (2019). What clinicians want: Contextualizing explainable machine learning for clinical end use. In *Machine learning for healthcare conference* (pp. 359–380).
- Wang, W., Arora, R., Livescu, K., & Bilmes, J. (2015) On deep multi-view representation learning. In *International conference on machine learning* (pp. 1083–1092).
- Wen, J., Wu, Z., Zhang, Z., Fei, L., Zhang, B., & Xu, Y. (2020). Cdimnet: Cognitive deep incomplete multi-view clustering network. In *International joint conference on artificial intelligence*.
- West, D. H. D. (1979). Updating mean and variance estimates: An improved method. *Communications of the ACM*, 22(9), 532–535. <https://doi.org/10.1145/359146.359153>
- Wickstrøm, K., Kampffmeyer, M., & Jenssen, R. (2018). Uncertainty modeling and interpretability in convolutional neural networks for polyp segmentation. In *IEEE International workshop on machine learning for signal processing* (pp. 1–6).
- Wickstrøm, K., Kampffmeyer, M., & Jenssen, R. (2020). Uncertainty and interpretability in convolutional neural networks for semantic segmentation of colorectal polyps. *Medical Image Analysis*, 60, 101619.
- Wickstrøm, K., Mikalsen, K., Kampffmeyer, M., Revhaug, A., & Jenssen, R. (2021). Uncertainty-aware deep ensembles for reliable and explainable predictions of clinical time series. *IEEE Journal of Biomedical and Health Informatics*, 25(7), 2435–2444. <https://doi.org/10.1109/JBHI.2020.3042637>
- Yang, B., Fu, X., Sidiropoulos, N. D., & Hong, M. (2017.) Towards k-means-friendly spaces: Simultaneous deep learning and clustering. In *International conference on machine learning* (pp. 3861–3870).
- Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In D. Fleet, T. Pajdla, B. Schiele, et al. (Eds.), *European conference on computer vision* (pp. 818–833).
- Zhang, J., Bargal, S. A., Lin, Z., Brandt, J., Shen, X., & Sclaroff, S. (2017). Top-down neural attention by excitation backprop. *International Journal of Computer Vision*, 126(10), 1084–1102. <https://doi.org/10.1007/s11263-017-1059-x>
- Zhang, Y., Song, K., Sun, Y., Tan, S., & Udel, M. (2019). “Why should you trust my explanation?” Understanding uncertainty in LIME explanations. In *Workshop on AI for social good*.