

RELAX: Representation Learning Explainability

Kristoffer K. Wickstrøm, Daniel J. Trosten, Sigurd Løkse, Karl Øyvind Mikalsen,
Michael C. Kampffmeyer, Robert Jenssen
Department of Physics and Technology, UiT The Arctic University of Norway

<https://machine-learning.uit.no/>

Abstract

Despite the significant improvements that representation learning via self-supervision has led to when learning from unlabeled data, no methods exist 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 and significantly outperforming the gradient-based baseline. 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. Finally, we illustrate the usability of RELAX in multi-view clustering and highlight that incorporating uncertainty can be essential for providing low-complexity explanations, taking a crucial step towards explaining representations.

1. Introduction

Interpretability is of vital importance for designing trustworthy and transparent deep learning-based systems [37, 48], and the field of explainable artificial intelligence (XAI) has made great improvements over the last couple of years [30, 42]. However, there exists no methods for attribution-based explanations of *representations*, despite the tremendous developments in representation learning using e.g. self-supervised learning [7, 8, 19]. Therefore, there is 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 [28, 50, 55], but there is no method for investigating which features are characteristic for the members of a cluster.

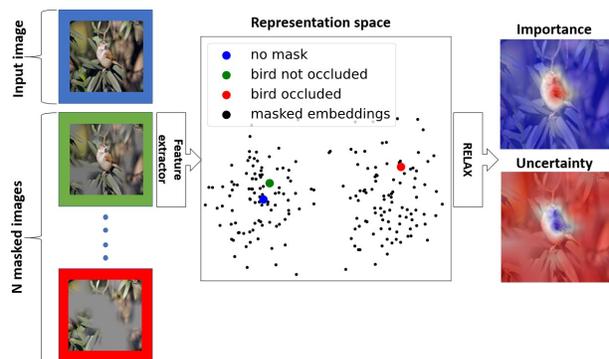


Figure 1. Conceptual illustration of RELAX. An image is passed through an encoder that produces a new vector representation of the image. Similarly, masked images are embedded in the same latent space. Input feature importance is estimated by measuring the similarity between the representation of the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

Representation learning explainability would also allow for a new approach for evaluating representation learning frameworks. Representation learning frameworks are typically evaluated by training simple classifiers on the representation produced by the feature extractor or through a downstream task [7, 8, 19]. 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 nec-

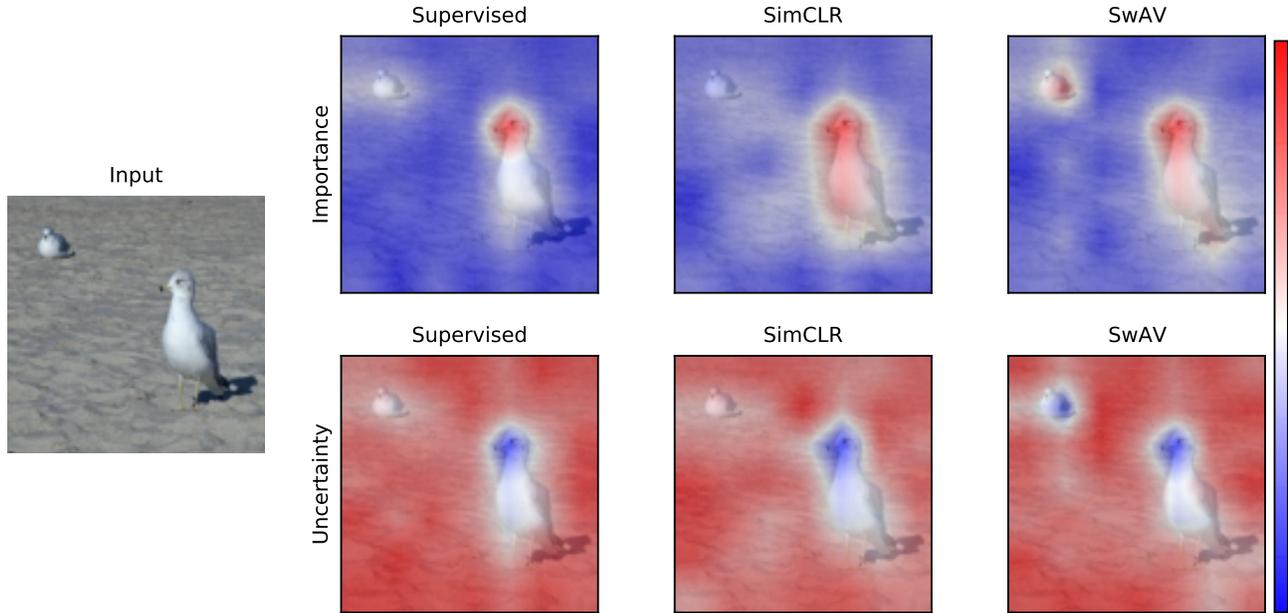


Figure 2. The figure shows the RELAX explanation and its uncertainty for the representation of the leftmost image for a number of widely used feature extractors. The first row displays the explanations for the representation and the second row shows the uncertainty associated with the different explanations. Red indicates high values and blue indicates low values. In this example, 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 [12].

essary for making trustworthy models, found that explainability alone was not enough and that uncertainty was also of high importance [48]. Uncertainty can also be used to reduce the complexity of explanations, as it allows for removal of uncertain parts of an explanations. Nevertheless, little work has been done on uncertainty in explanations of representations.

In this work, we present a new framework for explaining representations, entitled REpresentation LeArning eXplainability (RELAX), which is the first representation learning XAI method 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 masks, 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. Fig. 2 shows an example where RELAX is used to investigate the explanations 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 reduced the complexity of explanations.
- A theoretical analysis of the framework and derivation of an expression for the number of masks needed to obtain estimates with low error.
- A comprehensive experimental section that compares several widely used feature extraction models and a use case where RELAX enables explainability in incomplete multi-view clustering.

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 [21, 23] or global explanations [34].

Occlusion-based explainability. There exist a number of occlusion-based explainability methods. Systematically occluding an image with a gray rectangle and then measuring the change in activations could be used to provide coarse explanations for CNNs [56]. 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 [38]. A slightly different approach is meaningful perturbations, where a spatial perturbation mask that maximally affects the model’s output is optimized [16]. 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 [14]. 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 [42]. Similarly, [24] introduced a rate-distortion perspective to explainability. However, 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 [3, 38, 42]. 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 [22], 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 [4], but requires predefined concepts and label information without indicating the importance of individual pixels. Another approach mapped semantic concepts to vectorial embedding [15], but requires segmentation masks that are not available in the unsupervised setting. Lastly, representations have also been investigated from learnability and descriptibility perspectives [26], but this was achieved through human-annotators that are typically not available.

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 [17] in order to estimate the uncertainty in gradient-based explanations [52, 53], which was later followed by a similar approach that was based on Layer-wise Relevance Propaga-

tion [6]. Uncertainties inherent in the widely used LIME method [39] have been explored [58], and ensemble-based approaches, where uncertainty estimates are obtained by taking the standard deviation across the ensemble, have also been proposed [54]. 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 (Sec. 3.3). This is important, as uncertain explanations might give an unwarranted trust in the model. Our framework is inspired by RISE [38]. 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 the Appendix.

3.1. RELAX

The central idea of RELAX is that when informative parts are masked out, the representation should change significantly. Let $\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$ represent an image consisting of $H \times W$ pixels, and f denote a feature extractor that transforms an image into a representation $\mathbf{h} = f(\mathbf{X})$. 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:

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

Eq. (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 and computing the sample mean:

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

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

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.

We choose to measure the similarity $s(\mathbf{h}, \bar{\mathbf{h}}_n)$ between the representations using the cosine kernel:

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

where $\|\cdot\|$ denotes the Euclidean norm of a vector. There are several motivations for this choice. First, a large portion of feature extractors trained using self-supervised learning use the cosine kernel in their loss function [8, 9]. Therefore, it is the natural choice for measuring similarities in their latent space. Second, [29] 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. Third, the cosine kernel does not rely on hyper parameters that must be selected, which may be beneficial in an unsupervised setting where we cannot do cross validation.

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 [38]. 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 demand of RELAX, we derive the following lower bound on the number of masks required for a certain estimation error.

²See Sec. A of the Appendix for evaluation of masking strategies.

Theorem 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 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 B for verification of bound. In all of our experiments, we generate 3000 masks, as this will ensure that the estimation error is 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).

Theorem 2. Suppose the similarity function $s(\cdot, \cdot)$ is a valid Mercer kernel [33]. \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 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 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 [31].

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

Theorem 3. Suppose $s(\cdot, \cdot)$ is a valid Parzen window [47]. Then:

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

where $p_{ij}(\cdot)$ is a weighted Parzen density estimate [35] of the density of the masked embeddings:

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

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

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

3.2. Uncertainty in Explanations

Trusting an explanation without knowing its uncertainty can lead to an unjustified faith in the model. Our intuition 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:

$$U_{ij} = \text{Var}_{\mathbf{M}}[s(\mathbf{h}, \bar{\mathbf{h}})M_{ij}]. \quad (8)$$

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

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

Eq. (9) estimates the uncertainty of the RELAX-score for pixel (i, j) by measuring the difference between the similarity score and the explanations. 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 [2]. Secondly, it avoid computationally intense sampling methods, for instance through Monte Carlo sampling [17, 45]. Lastly, the uncertainty estimation can be incorporated into the computation of Eq. (2), as explained in Sec. 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 reduce the complexity 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 approach as:

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

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 to choose ϵ as:

$$\epsilon = \frac{1}{HW} \sum_{ij} \bar{U}_{ij}, \quad (11)$$

that is, the average uncertainty for a particular image. This provides a simple and intuitive way of selecting the threshold, which is motivated by wanting to only consider pixels that have high importance and low uncertainty. We refer to this uncertainty-filtered version of RELAX as U-RELAX.

3.4. One-Pass Version of RELAX

Computing Eq. (9) requires first computing Eq. (2), which 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 [51]. Importance is computed as:

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

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 all mask elements (i, j) after n masks. Uncertainty is computed as:

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

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

```
# 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)
```

Algorithm 1. Pytorch-like pseudocode for RELAX.

4. Evaluation and Baseline

4.1. Evaluation of Explanations

Evaluation is a developing subfield of XAI, and a unifying score is not agreed upon [11], even more so for explanations of representations. To evaluate the explanations we use several recent explainability evaluation metrics. All metrics are computed using the Quantus toolbox ⁴.

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 two localisation metrics, the *pointing game* [57] and *top-k intersection* [46]. 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. 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 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 measures 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 *iterative removal of features (IROF)* metric [40]. IROF segments out regions as identify by an explanations, and measures the decrease in classification score as more segments are removed.

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 [1, 43], 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:

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

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

Here, 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.

5. Experiments

To evaluate RELAX, we consider a number of widely used feature extractors trained with and without labels. This allows us to investigate potential differences between supervised and unsupervised methods, and will reveal how unsupervised methods compare against each other. For the supervised model, we use the pretrained model from Pytorch [36]. For the models trained without labels but with self-supervision, we use the SimCLR [8] and SwAV [7] 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 [13]. We consider a ResNet50 [20] as the backbone for the feature extractors, and all models are trained on ImageNet [10].

Implementation details. Similarly as in previous works [14, 42], we use the test split of the PASCAL VOC07 (VOC) [12] and the validation split of MSCOCO2014 (COCO) [27] for evaluating the localisation metrics, since they contain information about the location of the objects in the images. For the faithfulness metric, we use the validation set of ImageNet [10]. For all datasets, we randomly sample 1000 images for evaluation and repeat all experiments 3 times. As explained at the end of Sec. 3.1, 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 [57]. For the IROF score, we use Alexnet [25] as the classifier, as suggested in prior works [41].

5.1. Qualitative Results

Fig. 2 and 3 displays the explanation and the uncertainty in the explanations provided by RELAX for an image from the VOC and COCO dataset, respectively. See Appendix D for additional results, including qualitative gradient-based and U-RELAX 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.

Fig. 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

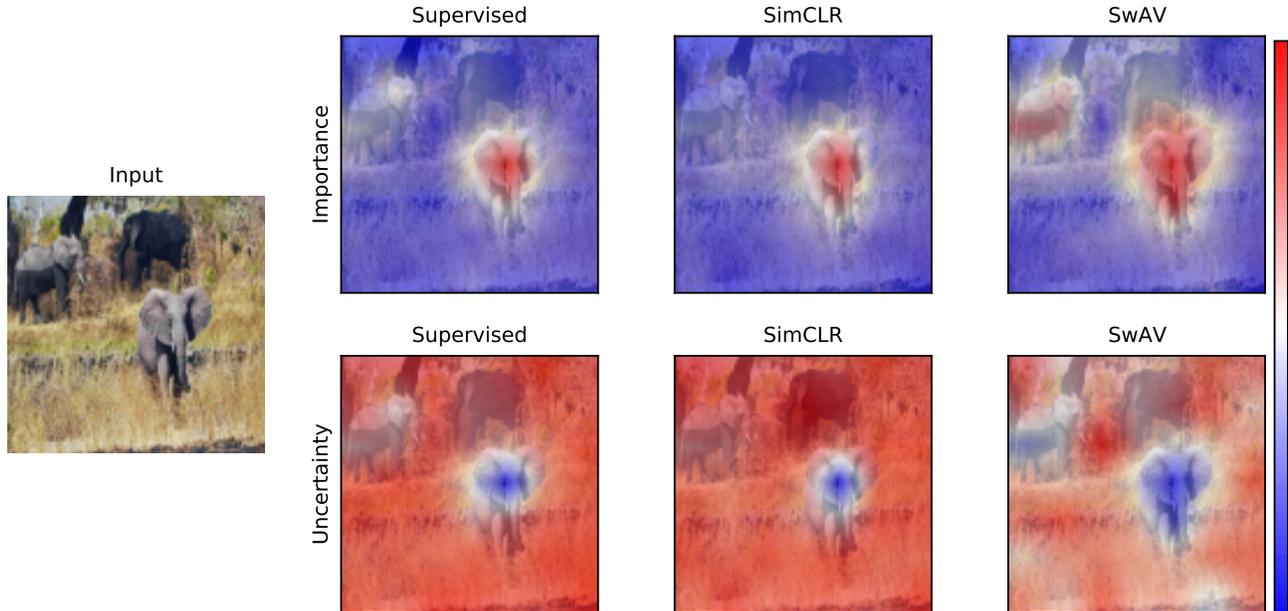


Figure 3. The figure shows the RELAX explanation and its uncertainty for the representation of the leftmost image for a number of widely used feature extractors. The first row displays the explanations for the representation and the second row shows the uncertainty associated with the different explanations. 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 [27].

Methods	Supervised		SimCLR		SwAV	
	VOC	COCO	VOC	COCO	VOC	COCO
Saliency	52.9±0.0 / 54.8±0.0	42.6±0.0 / 42.3±0.0	53.8±0.0 / 54.3±0.0	42.2±0.0 / 41.7±0.0	53.9±0.0 / 54.6±0.0	42.3±0.0 / 42.0±0.0
RELAX	87.1±0.3 / 86.4±0.1	73.1±0.3 / 72.2±0.1	84.5±0.2 / 85.1±0.2	68.8±1.2 / 67.9±0.1	85.6±0.2 / 84.9±0.2	67.6±0.3 / 67.0±0.2

Table 1. Pointing game / top-k intersection scores in percentages and averaged over 3 runs. Higher is better and bold numbers highlight the significantly best method. Results show that our method improves on the baseline across all scores.

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 explanations 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. 3 shows an image with 3 elephants, one in the foreground and two in the background, one of which is partially shaded. Again, RELAX enables investigation of interesting aspects of the representations, such as: are the models capable of recognizing all elephants and utilizing the information? All models highlight the elephant in the foreground as important with high certainty. However, there is little em-

phasis 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 Fig. 2 and 3, 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 this is, 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

Table 1 and 2 displays the quantitative evaluation of our proposed methodology using the evaluation metrics de-

scribed in Sec. 4.1, compared with the gradient-based baseline described in Sec. 4.2. The results show how the proposed method outperforms the baseline across all metrics. Note, the saliency explanation is deterministic, which is why the standard deviation is zero. Moreover, the low standard deviation scores 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 metrics.

Methods	Supervised	SimCLR	SwAV
Saliency	5408.3±0.0	5405.9±0.0	5407.9±0.0
RELAX	5408.1±0.2	5407.2±0.2	5408.0±0.4

Table 2. IROF scores averaged over 3 runs. Higher is better and bold numbers highlight the significantly best method. Results show that our method improves on the baseline.

5.3. U-RELAX for Low-Complexity Explanations

In some cases it can be desirable to remove parts of the input to produce explanations of lower complexity. We compare RELAX with U-RELAX using the same setup as for the faithfulness experiment, where complexity is measured as suggested by [5]. Results are reported in Table 3, which shows that U-RELAX can be utilized to reduce the complexity of explanations.

	Supervised	SimCLR	SwAV
RELAX	1082.3±0.0	1082.3±0.0	1082.3±0.0
U-RELAX	985.6±0.1	993.9±0.2	999.1±0.0

Table 3. Complexity scores averaged over 3 runs. Lower is better and bold numbers highlight the best method. Results show that U-RELAX can reduce the complexity of explanations.

5.4. Use Case for RELAX: Multi-View Clustering

To further illustrate RELAX’s ability to obtain insights into new tasks, we conduct an experiment on multi-view clustering. We learn a feature extractor using the Completer framework [28], which uses an information theoretic approach to allow for fusing several views into a unified representation. Clustering is performed through K-means on the learned representations. 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. Fig. 4 shows an example on Noisy MNIST [49], where one view is a digit and the other view is a noisy version of the same digit. The result

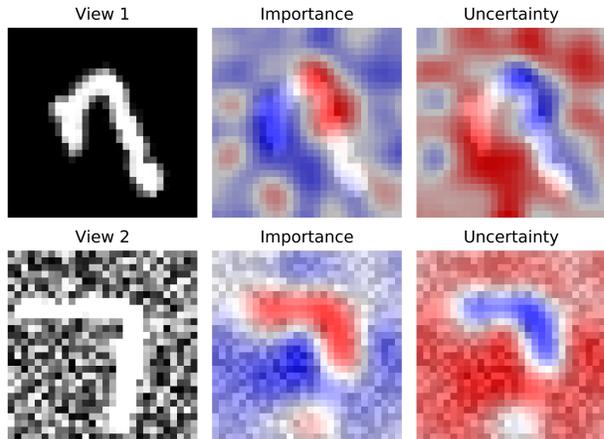


Figure 4. RELAX explanation and uncertainty for the representation of an example from Noisy MNIST image for a number of widely used feature extractors. The first row displays input, explanation, and uncertainty for view 1, and the second row for view 2. Red indicates high values and blue indicates low values. The Fig. shows that Completer is extracting complementary information from the two views for creating its unified representation.

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. See Appendix E for further details.

5.5. Limitations

We have demonstrated the utility and value of RELAX, but there are also limitations to consider. First, RELAX assumes that the masks can cover an object completely, partially, or not at all. If an object is too big or too small it can be difficult to mask them efficiently. Second, the limited number of available baselines weakens the evaluation, however, it also strongly highlights the need for development of more methods for explaining representations. Lastly, adopting RELAX to new non-image modalities is not trivial, as it requires a suitable masking procedure for the modality in questions. See Appendix F of the supplementary for discussion on the potential negative societal impact of RELAX.

6. Conclusion

In this work, we presented RELAX, a framework for explaining representations produced by deep learning-based feature extractors. RELAX is based on masking out parts of an image and measuring 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 reduce the complexity of explanations, and that RELAX provides new insights into multi-view clustering. We believe that RELAX can be an important addition in the intersection between XAI and representation learning.

Appendices

A. Masking Strategies

Fig. 5 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 [44] 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 [18]. 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.

B. Proofs

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

B.1. Proof of Theorem 1

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

Definition 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 $i = 1, \dots, N$,

$$\sup_{x_1, \dots, x_N, x'_i \in A} |f(x_1, \dots, x_N, x'_i) - f(x_1, \dots, x_N, x_i)| \leq c_i \quad (15)$$

We then have the following lemmas:

Lemma 3.1 (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, for all $t > 0$*

$$\begin{aligned} P(|f(X_1, \dots, X_N) - \mathbb{E}[f(X_1, \dots, X_N)]| \geq t) \\ \leq 2e^{-\frac{2t^2}{\sum_{n=1}^N c_n^2}} \end{aligned} \quad (16)$$

Proof. See [32]. □

Lemma 3.2. *Let X_1, \dots, X_N and f be defined as in Lemma 3.1, 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 [32]. □

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

$$X_n = \frac{s(\mathbf{h}, \bar{\mathbf{h}}_n) M_{ij}(n)}{N}, \quad (17)$$

and

$$f(X_1, \dots, X_N) = \sum_{n=1}^N X_n. \quad (18)$$

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 3.2.

Now, observe that

$$f(X_1, \dots, X_N) = \frac{1}{N} \sum_{n=1}^N s(\mathbf{h}, \bar{\mathbf{h}}_n) M_{ij}(n) = \bar{R}_{ij}. \quad (19)$$

Combining Lemmas 3.1 and 3.2 then gives

$$P(|\bar{R}_{ij} - R_{ij}| \geq t) \leq 2e^{-\frac{2t^2}{\sum_{n=1}^N (1/N)^2}} \quad (20)$$

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

$$P(|\bar{R}_{ij} - R_{ij}| \geq t) \leq 2e^{-\frac{2t^2}{\sum_{n=1}^N (1/N)^2}} \quad (21)$$

$$= 2e^{-2t^2 \left(-\frac{\ln(\delta/2)}{2t^2}\right)} \quad (22)$$

$$= 2e^{\ln(\delta/2)} \quad (23)$$

$$= \delta, \quad (24)$$

which concludes our proof. □

In Fig. 6 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 PASCAL VOC dataset. To obtain a value for R_{ij} , we use 10000 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 it is possible to obtain a similarity of 0, it is highly unlikely since our masking strategy never removes all information in an image.



Figure 5. Comparison of different masking strategies for RELAX. 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.

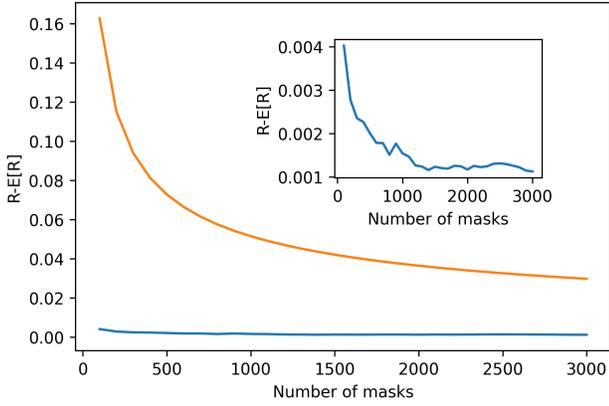


Figure 6. Empirical evaluation of the derived bound for the number of masks necessary for low estimation error. We calculate the absolute error as the number of masks increase, average over 10 randomly samples images from the PASCAL VOC dataset. To obtain a value for R_{ij} , we use 10000 masks and average over 10 runs for a single sample. Results indicate that the estimation error is much lower than the predicted bound.

B.2. Proof of Theorem 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

$$\bar{R}_{ij} = \frac{1}{N} \sum_{n=1}^N \langle \phi(\mathbf{h}), \phi(\bar{\mathbf{h}}_n) \rangle_{\mathcal{H}} M_{ij}(n) \quad (25)$$

$$= \langle \phi(\mathbf{h}), \frac{1}{N} \sum_{n=1}^N \phi(\bar{\mathbf{h}}_n) M_{ij}(n) \rangle_{\mathcal{H}} \quad (26)$$

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

B.3. Proof of Theorem 3

Proof. Observe that

$$\bar{R}_{ij} \cdot \frac{N}{\sum_{n'=1}^N M_{ij}(n')} \quad (27)$$

$$= \frac{N}{\sum_{n'=1}^N M_{ij}(n')} \cdot \frac{1}{N} \sum_{n=1}^N s(\cdot, \bar{\mathbf{h}}_n) M_{ij}(n) \quad (28)$$

$$= \frac{1}{\sum_{n'=1}^N M_{ij}(n')} \sum_{n=1}^N s(\cdot, \bar{\mathbf{h}}_n) M_{ij}(n) \quad (29)$$

$$= p_{ij}(\mathbf{h}) \quad (30)$$

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

C. One-Pass Versus Two-Pass RELAX

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 Figure 7 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.

D. Qualitative Results

This section presents additional qualitative results.

D.1. Additional Qualitative Results for RELAX

Fig. 8 to 17 displays examples of explanations and their associated uncertainty, provided by RELAX, for images from the VOC [12] and COCO [27] dataset. Fig. 8 displays an example where all feature extractors agree in terms of importance, but the degree of uncertainty varies. Fig. 9 shows an example where only SwAV highlight both objects

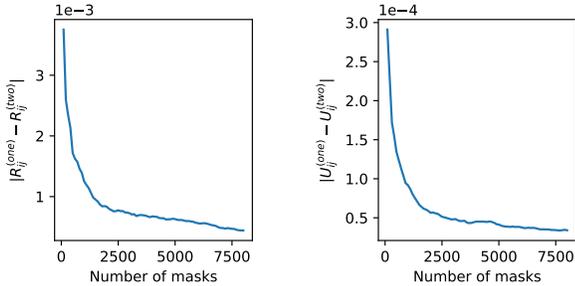


Figure 7. Absolute error of one-pass versus two-pass version of RELAX for importance (leftmost figure) and uncertainty (rightmost figure), averaged over 50 images from the VOC dataset. The figure shows how the difference between the versions is small for both the importance and uncertainty estimates.

as important for the representation. Similarly, Fig. 10 displays an example where only SwAV is considering both the person and the car as important for the representation. Fig. 10 to 17 shows similar examples where RELAX provides insights into the different feature extractors.

D.2. Qualitative Results for Saliency Explanation

Fig. 18 and Fig. 19 show a qualitative comparison between the RELAX and saliency explanation for a representation of an image. Both Figs. 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.

D.3. Qualitative Results for U-RELAX

Fig. 20 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.

E. Multi-View Clustering

We employ the state-of-the-art multi-view clustering approach in Completer [28] to illustrate the usability of RELAX in multi-view clustering. Completer uses an information theoretic objective to incorporate information from several views of the same data sample into a unified representation. It uses individual encoders for each view of the data, and concatenates the representation from each encoder to produce a unified representation. 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.

F. Negative Societal Impact

An important purpose of explainability is to enhance trust in deep learning models. However, there is also a risk for negative consequences with an increased focus on explainability. An undue trust in a model based on its explanation might lead to erroneous or unfair decisions. And users with little expertise in deep learning might misunderstand the explanations, which can lead to false conclusions. Such issues emphasise the need for explanations with uncertainty, as we have proposed in RELAX.

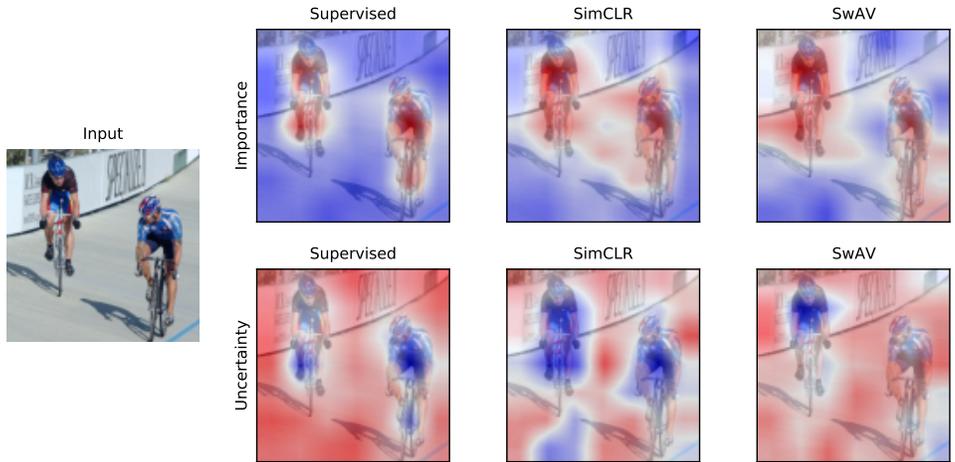


Figure 8. Example from the VOC dataset.

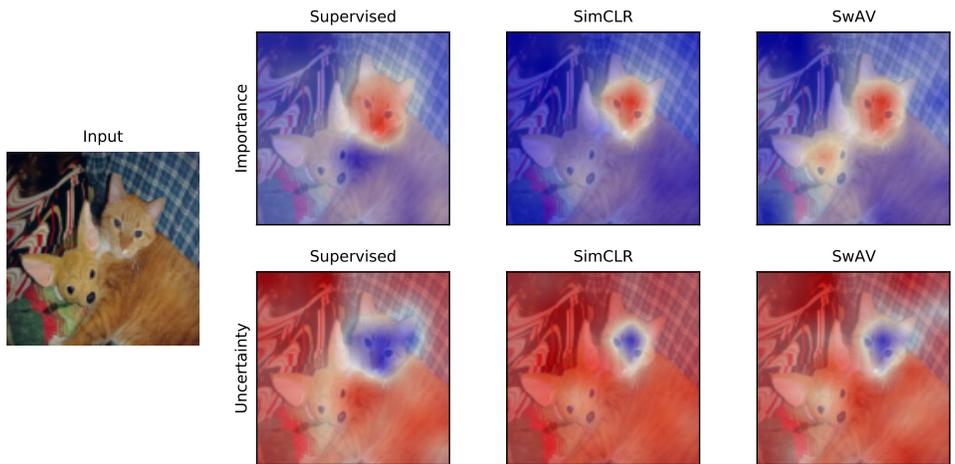


Figure 9. Example from the COCO dataset.

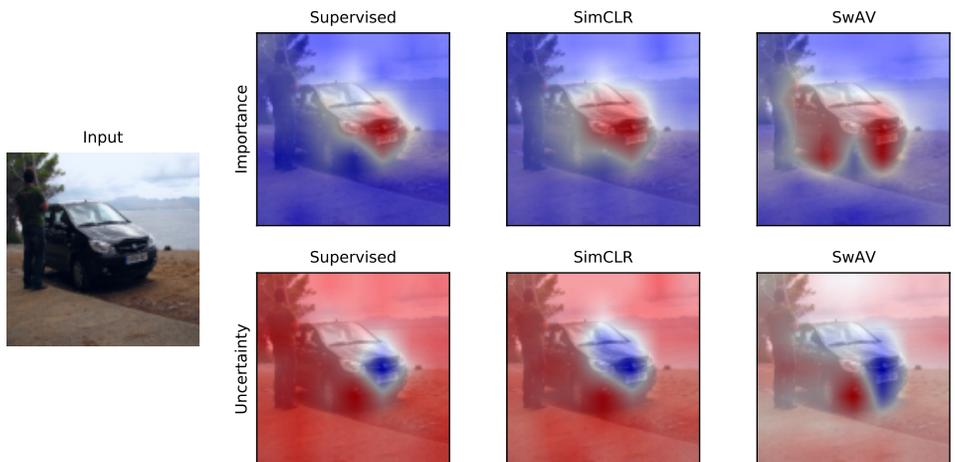


Figure 10. Example from the VOC dataset.

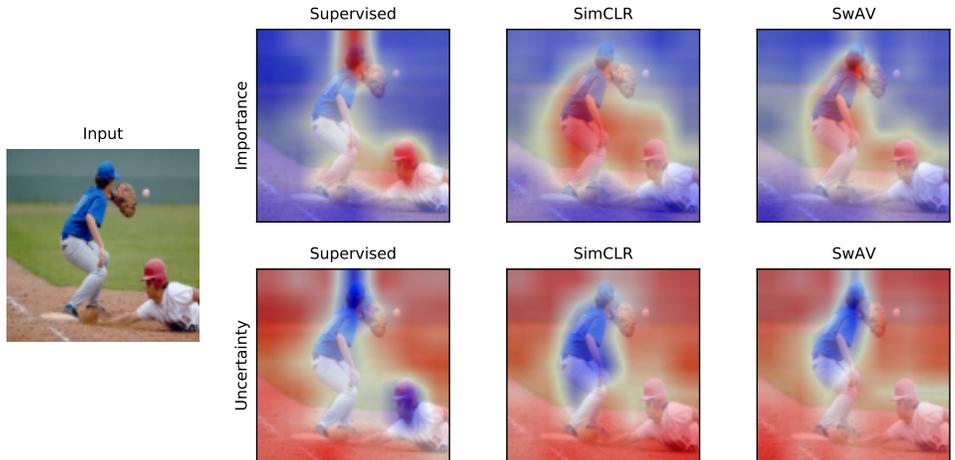


Figure 11. Example from the COCO dataset.

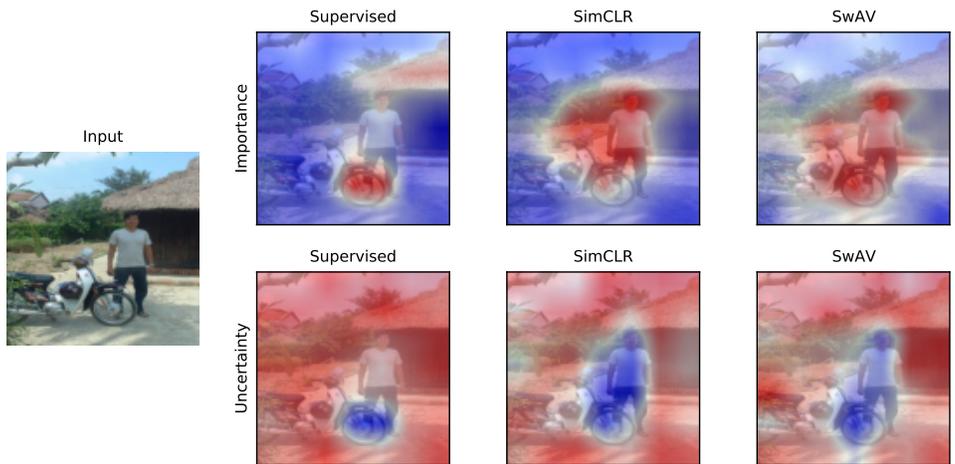


Figure 12. Example from the VOC dataset.

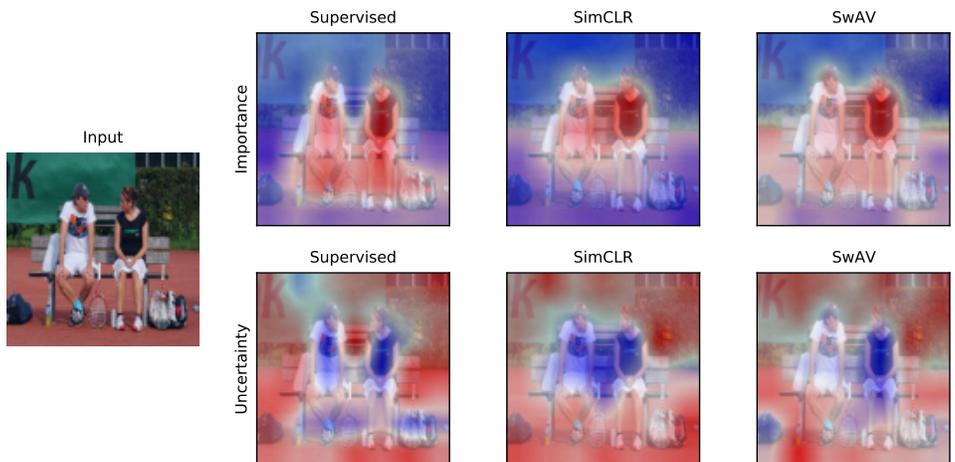


Figure 13. Example from the COCO dataset.

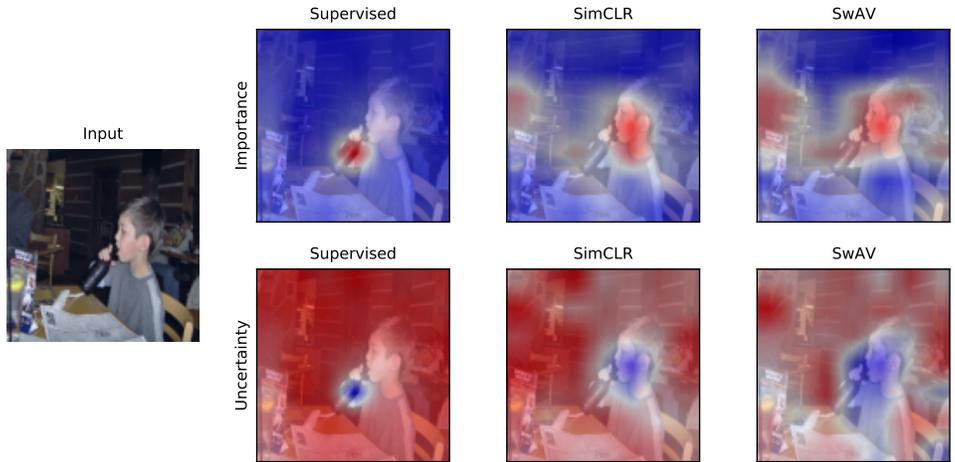


Figure 14. Example from the VOC dataset.

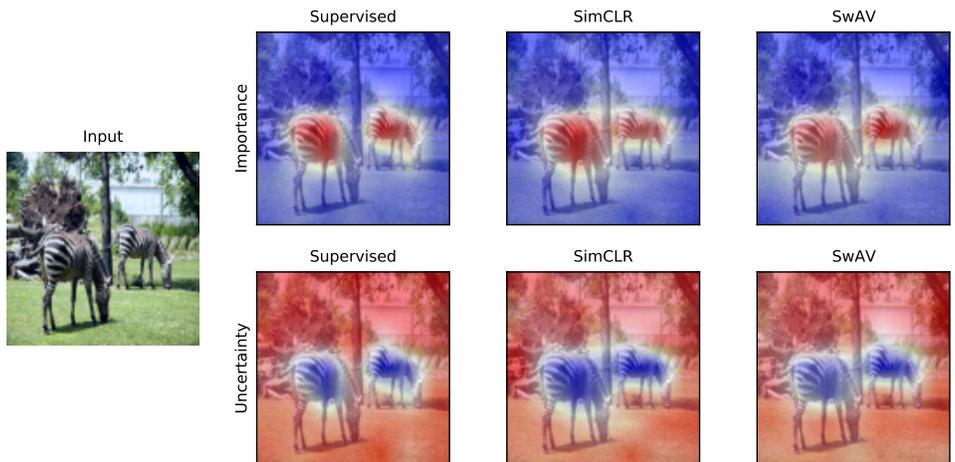


Figure 15. Example from the COCO dataset.

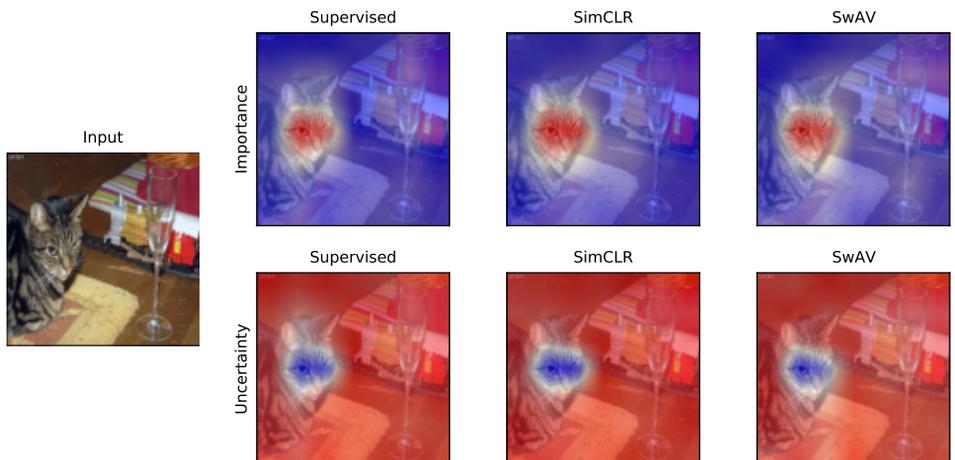


Figure 16. Example from the VOC dataset.

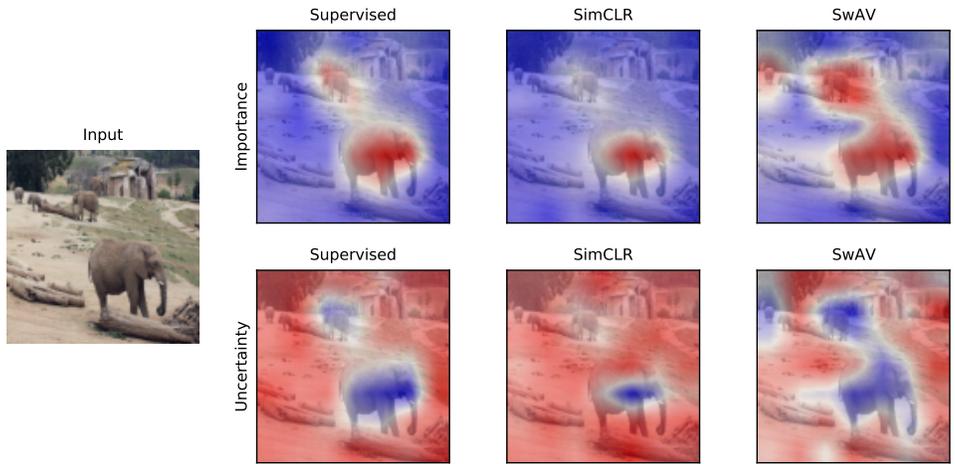


Figure 17. Example from the COCO dataset.



Figure 18. Comparison of RELAX and saliency explanation for representation of 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.



Figure 19. Comparison of RELAX and Saliency explanation for representation of image from PASCAL VOC. The example shows how RELAX captures information about both objects in the image, while the saliency explanation is focused on the gap in between the two objects.

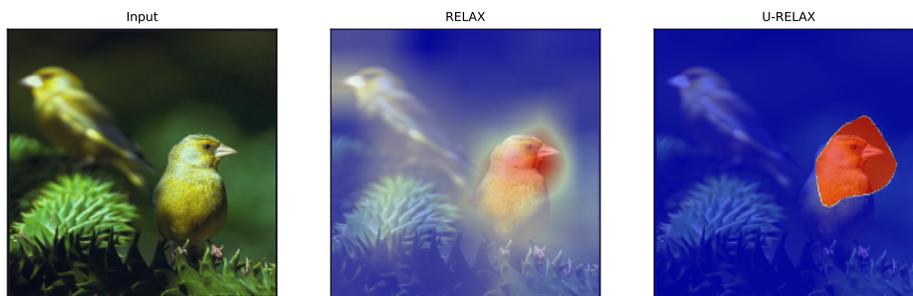


Figure 20. Comparison of RELAX and U-RELAX on an image taken from PASCAL VOC. In this case, the emphasis on the bird in the background is removed as the uncertainty was too high for this part of the explanation.

References

- [1] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2018. 6
- [2] Javier Antoran, Umang Bhatt, Tameem Adel, Adrian Weller, and José Miguel Hernández-Lobato. Getting a {clue}: A method for explaining uncertainty estimates. In *International Conference on Learning Representations*, 2021. 5
- [3] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE*, 10(7):e0130140, 07 2015. 3
- [4] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *IEEE Computer Vision and Pattern Recognition*, 2017. 3
- [5] Umang Bhatt, Adrian Weller, and José M. F. Moura. Evaluating and aggregating feature-based model explanations. In *International Joint Conference on Artificial Intelligence*, pages 3016–3022, 2020. 8
- [6] Kirill Bykov, Marina M.-C. Höhne, Klaus-Robert Müller, Shinichi Nakajima, and Marius Kloft. How much can I trust you? - quantifying uncertainties in explaining neural networks. *CoRR*, abs/2006.09000, 2020. 3
- [7] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *Advances in Neural Information Processing Systems*, pages 9912–9924, 2020. 1, 6
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pages 1597–1607, 2020. 1, 4, 6
- [9] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *IEEE Computer Vision and Pattern Recognition*, pages 15750–15758, June 2021. 4
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Computer Vision and Pattern Recognition*, pages 248–255. Ieee, 2009. 6
- [11] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning, 2017. 6
- [12] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (VOC) challenge. *International Journal of Computer Vision*, pages 303–338, 2009. 2, 6, 10
- [13] William Falcon and Kyunghyun Cho. A framework for contrastive self-supervised learning and designing a new approach. *arXiv preprint arXiv:2009.00104*, 2020. 6
- [14] Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Understanding deep networks via extremal perturbations and smooth masks. In *IEEE International Conference on Computer Vision*, October 2019. 3, 6
- [15] Ruth Fong and Andrea Vedaldi. Net2vec: Quantifying and explaining how concepts are encoded by filters in deep neural networks. In *IEEE Computer Vision and Pattern Recognition*, pages 8730–8738, 2018. 3
- [16] Ruth C. Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *IEEE International Conference on Computer Vision*, pages 3449–3457, 2017. 3
- [17] Yarín Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning*, pages 1050–1059, 2016. 3, 5
- [18] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V. Le. Dropblock: A regularization method for convolutional networks. In *International Conference on Neural Information Processing Systems*, page 10750–10760, 2018. 9
- [19] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *IEEE Computer Vision and Pattern Recognition*, June 2020. 1
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 CVPR*, pages 770–778, 2016. 6
- [21] A.-H. Karimi, G. Barthe, B. Balle, and I. Valera. Model-agnostic counterfactual explanations for consequential decisions. In *International Conference on Artificial Intelligence and Statistics*, pages 895–905, 2020. 3
- [22] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie J. Cai, James Wexler, Fernanda B. Viégas, and Rory Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). In *International Conference on Machine Learning*, pages 2673–2682, 2018. 3
- [23] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *In-*

- ternational Conference on Machine Learning*, page 1885–1894, 2017. 3
- [24] Stefan Kolek, Duc Anh Nguyen, Ron Levie, Joan Bruna, and Gitta Kutyniok. A rate-distortion framework for explaining black-box model decisions, 2021. 3
- [25] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. *Commun. ACM*, 60(6):84–90, May 2017. 6
- [26] Iro Laina, Ruth C. Fong, and Andrea Vedaldi. Quantifying learnability and descriptibility of visual concepts emerging in representation learning. In *Advances in Neural Information Processing Systems*, 2020. 3
- [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *Computer Vision – ECCV 2014*, pages 740–755. Springer International Publishing, 2014. 6, 7, 10
- [28] Yijie Lin, Yuanbiao Gou, Zitao Liu, Boyun Li, Jiancheng Lv, and Xi Peng. Completer: Incomplete multi-view clustering via contrastive prediction. In *IEEE Computer Vision and Pattern Recognition*, pages 11174–11183, June 2021. 1, 8, 11
- [29] W. Liu, R. Lin, Z. Liu, L. Xiong, B. Schölkopf, and A. Weller. Learning with hyperspherical uniformity. In *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 1180–1188. PMLR, Apr. 2021. 4
- [30] Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, page 4768–4777, 2017. 1
- [31] P. McCullagh and J.A. Nelder. *Generalized Linear Models, Second Edition*. Chapman & Hall, 1989. 4
- [32] Colin McDiarmid. *On the method of bounded differences*, page 148–188. Cambridge University Press, 1989. 9
- [33] J. Mercer. 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, 1909. 4
- [34] Alexander Mordvintsev, Christopher Olah, and Mike Tyka. Inceptionism: Going deeper into neural networks, 2015. 3
- [35] Emanuel Parzen. On estimation of a probability density function and mode. *The Annals of Mathematical Statistics*, 33(3):1065–1076, Sept. 1962. 4
- [36] Adam Paszke, Sam Gross, Francisco Massa, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019. 6
- [37] Dino Pedreschi, Fosca Giannotti, Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, and Franco Turini. Meaningful explanations of black box AI decision systems. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33:9780–9784, July 2019. 1
- [38] Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. In *Proceedings of the British Machine Vision Conference*, 2018. 3, 4
- [39] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ”why should I trust you?”: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pages 1135–1144, 2016. 3
- [40] Laura Rieger and Lars Kai Hansen. Irof: a low resource evaluation metric for explanation methods. In *Proceedings of the Workshop AI for Affordable Healthcare at ICLR 2020*, 2020. Workshop AI for Affordable Healthcare at ICLR 2020 ; Conference date: 26-04-2020 Through 26-04-2020. 6
- [41] Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE TNNLS*, 28(11):2660–2673, 2017. 6
- [42] Karl Schulz, Leon Sixt, Federico Tombari, and Tim Landgraf. Restricting the flow: Information bottlenecks for attribution. In *International Conference on Learning Representations*, 2020. 1, 3, 6
- [43] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. In *ICLR Workshop*, 2015. 6
- [44] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, pages 1929–1958, 2014. 9
- [45] Mattias Teye, Hossein Azizpour, and Kevin Smith. Bayesian uncertainty estimation for batch normalized deep networks. In *International Conference on Machine Learning*, pages 4907–4916, 2018. 5
- [46] Jonas Theiner, Eric Müller-Budack, and Ralph Ewerth. Interpretable semantic photo geolocalization. *CoRR*, abs/2104.14995, 2021. 6

- [47] Sergios Theodoridis and Konstantinos Koutroumbas. *Pattern Recognition, Fourth Edition*. Academic Press, 2009. 4
- [48] Sana Tonekaboni, Shalmali Joshi, Melissa D. McCradden, and Anna Goldenberg. What clinicians want: Contextualizing explainable machine learning for clinical end use. In *Machine Learning for Healthcare Conference*, pages 359–380, 2019. 1, 2
- [49] Weiran Wang, Raman Arora, Karen Livescu, and Jeff Bilmes. On deep multi-view representation learning. In *International Conference on Machine Learning*, page 1083–1092, 2015. 8
- [50] Jie Wen, Zheng Zhang, Yong Xu, Bob Zhang, Lunke Fei, and Guo-Sen Xie. Cdimc-net: Cognitive deep incomplete multi-view clustering network. In *International Joint Conference on Artificial Intelligence*, 2020. 1
- [51] D. H. D. West. Updating mean and variance estimates: An improved method. *Commun. ACM*, 22(9):532–535, sep 1979. 5
- [52] Kristoffer Wickstrøm, Michael Kampffmeyer, and Robert Jenssen. Uncertainty modeling and interpretability in convolutional neural networks for polyp segmentation. In *IEEE International Workshop on Machine Learning for Signal Processing*, pages 1–6, 2018. 3
- [53] Kristoffer Wickstrøm, Michael Kampffmeyer, and Robert Jenssen. Uncertainty and interpretability in convolutional neural networks for semantic segmentation of colorectal polyps. *Medical Image Analysis*, 60:101619, 2020. 3
- [54] Kristoffer Wickstrøm, KØ Mikalsen, Michael Kampffmeyer, Arthur Revhaug, and Robert Jenssen. Uncertainty-aware deep ensembles for reliable and explainable predictions of clinical time series. *IEEE Journal of Biomedical and Health Informatics*, 25(7):2435–2444, 2021. 3
- [55] Bo Yang, Xiao Fu, Nicholas D. Sidiropoulos, and Mingyi Hong. Towards k-means-friendly spaces: Simultaneous deep learning and clustering. In *International Conference on Machine Learning*, page 3861–3870, 2017. 1
- [56] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *European Conference on Computer Vision*, pages 818–833, 2014. 3
- [57] Jianming Zhang, Sarah Adel Bargal, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. Top-down neural attention by excitation back-prop. *International Journal of Computer Vision*, 126(10):1084–1102, Dec. 2017. 6
- [58] Yujia Zhang, Kuangyan Song, Yiming Sun, Sarah Tan, and Madeleine Udell. “Why Should You Trust My Explanation?” Understanding Uncertainty in LIME Explanations. In *Workshop on AI for Social Good*, 2019. 3