

---

# Concept Activation Regions: A Generalized Framework For Concept-Based Explanations

---

**Jonathan Crabbé**  
 University of Cambridge  
 jc2133@cam.ac.uk

**Mihaela van der Schaar**  
 University of Cambridge  
 The Alan Turing Institute  
 UCLA  
 mv472@cam.ac.uk

## Abstract

Concept-based explanations permit to understand the predictions of a deep neural network (DNN) through the lens of concepts specified by users. Existing methods assume that the examples illustrating a concept are mapped in a fixed direction of the DNN’s latent space. When this holds true, the concept can be represented by a concept activation vector (CAV) pointing in that direction. In this work, we propose to relax this assumption by allowing concept examples to be scattered across different clusters in the DNN’s latent space. Each concept is then represented by a region of the DNN’s latent space that includes these clusters and that we call concept activation region (CAR). To formalize this idea, we introduce an extension of the CAV formalism that is based on the kernel trick and support vector classifiers. This CAR formalism yields global concept-based explanations and local concept-based feature importance. We prove that CAR explanations built with radial kernels are invariant under latent space isometries. In this way, CAR assigns the same explanations to latent spaces that have the same geometry. We further demonstrate empirically that CARs offer (1) more accurate descriptions of how concepts are scattered in the DNN’s latent space; (2) global explanations that are closer to human concept annotations and (3) concept-based feature importance that meaningfully relate concepts with each other. Finally, we use CARs to show that DNNs can autonomously rediscover known scientific concepts, such as the prostate cancer grading system.

## 1 Introduction

Deep learning models are both useful and challenging. Their utility is reflected in their increasing contributions to sophisticated tasks such as natural language processing [1, 2], computer vision [3, 4] and scientific discovery [5–8]. Their challenging nature can be attributed to their inherent complexity. State of the art deep models typically contain millions to billions parameters and, hence, appear as *black-boxes* to human users. The opacity of black-box models make it difficult to: anticipate how models will perform at deployment [9]; reliably distil knowledge from the models [10] and earn the trust of stakeholders in high-stakes domains [11–13]. With the aim of increasing the transparency of black-box models, the field of explainable AI (XAI) developed [14–17]. We can broadly divide XAI methods in 2 categories: ① Methods that restrict the model’s architecture to enable explanations. Examples include attention models that motivate their predictions by highlighting features they pay attention to [18] and prototype-based models that motivate their predictions by highlighting relevant examples from their training set [19]. ② *Post-hoc* methods that can be used in a plug-in fashion to provide explanation for a pre-trained model. Examples include *feature importance methods* (also known as feature attribution or saliency methods) that highlight features the model is sensitive to [20–26]; *example importance methods* that identify influential training examples [27–29] and *hybrid*

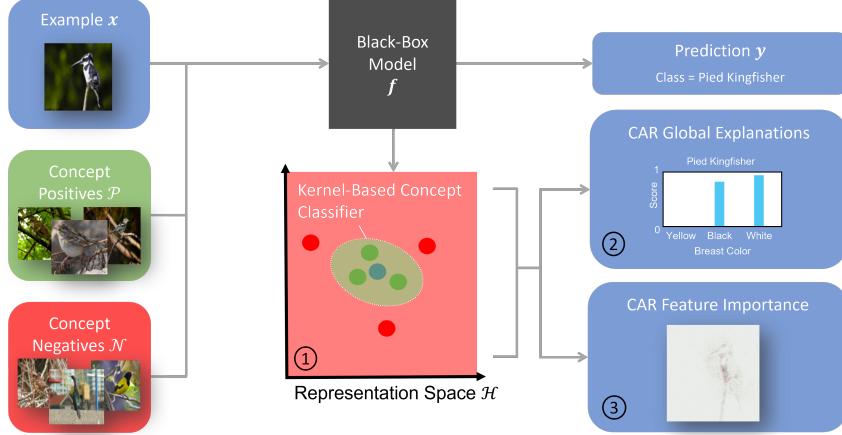


Figure 1: Illustration of the CAR framework. To define a concept, a set of positive and negative examples are fed to the model. In this illustration, we consider a neural network predicting the species of a bird based on a picture. The concept we illustrate is *white breast*. Concept positive and negative images respectively exhibit birds with and without white breasts. CARs aim at explaining the prediction for a given example  $x$ . ① All the examples are mapped in the model’s representation space  $\mathcal{H}$ . Unlike CAVs [31], CARs rely on kernel-based support vector classifiers (SVC) and do not require the positive and negative sets to be linearly separable in  $\mathcal{H}$ . ② CAR uses the SVC to output TCAR scores that indicate how classes are related to concepts. ③ The CAR formalism can be used in combination with any feature importance method for deep neural networks. This permits to create concept-specific saliency maps.

*methods* combining the two previous approaches [30]. In this work, we focus on a different type of explanation methods known as *concept-based* explanations. Let us now summarize the relevant literature to contextualize our own contribution.

**Related work.** Concept-based explanations were first formalized with the *concept activation vector* (CAV) formalism [31, 32]. Given a concept specified by the user (e.g. stripes in an image), linear classifiers are used as probes [33] to assess whether a deep model’s representation space separates examples where a concept is present (*concept positive examples*) from examples where a concept is absent (*concept negative examples*). A CAV is then extracted from the linear classifier’s weights. With this CAV, it is possible to provide post-hoc explanations such as the sensitivity of a model’s prediction to the presence/absence of a concept. It goes without saying that several concepts are needed to explain the prediction of a deep model. To formalize this idea, existing works have used concept basis decomposition [34] and sufficient statistics [35]. Although most applications of CAV involve image data, we note that the formalism has been successfully applied to time-series data [36, 37]. While concepts are typically specified by the user, early works in computer vision have been undertaken to discover concepts in the form of meaningful image segmentations [38]. The main criticism against the CAV formalism is that it requires concept positive examples to be linearly separable from concept negative examples [39, 40]. This is because linear separability of concept sets is a restrictive criterion that the model is not explicitly trained to fulfil. To address this issue, works like *concept whitening transformations* [39] and *concept bottleneck models* [41–43] propose to do away with the post-hoc nature of concept-based explanations. These methods introduce new neural network architectures that permit to train the models with concept labels. We stress that this requires the set of concepts to be specified *before* training the model. This assumes that we know what concepts are relevant for a model to solve a downstream task *a-priori*. This is not the case whenever we train a model to solve a task for which little or no knowledge is available. In this setup, it seems more appropriate to train a model for the task first and, then, perform a post-hoc analysis of the model to determine the concepts that were relevant in providing a solution. Furthermore, recent concerns have emerged regarding the reliability of interpretations provided by these altered model architectures [44, 45].

**Contributions.** In this work, our purpose is to retain the flexible post-hoc nature of concept-based explanations without assuming that the concept sets are linearly separable. To that aim, we introduce *concept activation regions* (CARs), an extension of the CAV framework illustrated in Figure 1.

**① Generalized formalism.** As a substitute to linear separability, we propose in Section 2.1 to adapt the smoothness assumption from semi-supervised learning [46]. Intuitively, this more general assumption only requires positive and negative examples to be scattered across distinct clusters in the model’s representation space. In practice, this generalization is implemented by substituting CAV’s linear classifiers by kernel-based support vector classifiers. We demonstrate that choosing radial kernels lead to CAR classifiers that are invariant under isometries of the latent space. This permits to assign identical explanations to latent spaces characterized by the same geometry. Moreover, we show in Section 3.1.1 that our CAR classifiers yield a substantially more accurate description of how concepts are distributed in deep model’s representation spaces. **② Better global explanations.** With the nonlinear decision boundaries of our support vector classifiers, there is no obvious way to adapt the notion of concept activation vectors. Since CAV’s concept importance (TCAV score) is computed with these vectors, we need an alternative approach. In Section 2.2, we propose to define concept importance by building on our smoothness assumption. Concretely, a concept is important for a given example if the model’s representation for this example lies in a cluster of concept positive representations. With this characterization, we define TCAR scores that are the CAR equivalent of TCAV scores. In Section 3.1.2, we demonstrate that TCAR scores lead to global explanations that are more consistent with concept annotations provided by humans. **③ Concept-based feature importance.** In Section 2.3, we argue that our CAR formalism permits to assign concept-specific feature importance scores for each example fed to the neural network. We verify empirically in Section 3.1.3 that those feature importance scores reflect meaningful concept associations. Finally, we illustrate in Section 3.2 how these contributions permit to establish that deep models implicitly discover known scientific concepts.

## 2 Concept Activation Regions (CARs)

### 2.1 Preliminaries

We assume a typical supervised setting where each sample is represented by a couple  $(\mathbf{x}, \mathbf{y})$  with input features  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^{d_X}$  and a label  $\mathbf{y} \in \mathcal{Y} \subseteq \mathbb{R}^{d_Y}$ , where  $d_X, d_Y \in \mathbb{N}^*$  are respectively the dimensions of the feature (or input) and label (or output) spaces<sup>1</sup>. In order to predict the labels from the features, we are given a deep neural network (DNN)  $f = l \circ g : \mathcal{X} \rightarrow \mathcal{Y}$ , where  $g : \mathcal{X} \rightarrow \mathcal{H}$  is a feature extractor that maps features  $\mathbf{x} \in \mathcal{X}$  to latent representations  $\mathbf{h} = g(\mathbf{x}) \in \mathcal{H} \subseteq \mathbb{R}^{d_H}$  and  $l : \mathcal{H} \rightarrow \mathcal{Y}$  maps latent representations  $\mathbf{h} \in \mathcal{H}$  to labels  $\mathbf{y} = l(\mathbf{h}) \in \mathcal{Y}$ . The representation  $\mathbf{h} = g(\mathbf{x})$  typically corresponds to the output from one of the DNN’s hidden layer. We assume that the model  $f$  was obtained by fitting a training set  $\mathcal{D}_{\text{train}} \subset \mathcal{X} \times \mathcal{Y}$ . Our purpose is to *understand* how the neural network representation induced by  $g$  allows the model  $f$  to solve the prediction task. To that end, we shall use concept-based explanations. Those explanations rely on a set of  $C \in \mathbb{N}^*$  concepts indexed<sup>2</sup> by  $c \in [C]$ , where  $[C]$  denotes the set of positive integers between 1 and  $C$ . Each concept  $c \in [C]$  is defined by the user through a set of  $N^c \in \mathbb{N}^*$  positive examples  $\mathcal{P}^c = \{\mathbf{x}^{c,n} \mid n \in [N^c]\}$  and a set of  $N^{-c}$  negative examples  $\mathcal{N}^c = \{\mathbf{x}^{-c,n} \mid n \in [N^{-c}]\}$ . The concept  $c \in [C]$  is present in the positive examples  $\mathcal{P}^c$  and absent from the negative examples  $\mathcal{N}^c$ . For instance, if we work in a computer vision context and the concept of interest is  $c = \text{Stripes}$ , then  $\mathcal{P}^c$  contains images with stripes (e.g. zebra images) and  $\mathcal{N}^c$  contains images without stripes (e.g. cow images). We note that non-binary concepts (e.g. a colour) can be one-hot encoded to fit this formulation. In this paper, we use balanced sets:  $N^c = N^{-c}$ .

**CAV formalism.** Let us start by summarizing the CAV formalism [31]. This formalism relies on the crucial assumption that the sets  $g(\mathcal{P}^c) \subset \mathcal{H}$  and  $g(\mathcal{N}^c) \subset \mathcal{H}$  are linearly separable in  $\mathcal{H}$ . Hence, it is possible to find a vector  $\mathbf{w}^c \in \mathbb{R}^{d_H}$  and a bias term  $b^c \in \mathbb{R}$  such that  $(\mathbf{w}^c)^\top \mathbf{h} + b^c > 0$  for all  $\mathbf{h} \in g(\mathcal{P}^c)$  and, conversely,  $(\mathbf{w}^c)^\top \mathbf{h} + b^c < 0$  for all  $\mathbf{h} \in g(\mathcal{N}^c)$ . The vector  $\mathbf{w}^c$  is called the *concept activation vector* (abbreviated CAV) associated to the concept  $c \in [C]$ . Geometrically, this vector is normal to the hyperplane separating  $g(\mathcal{P}^c)$  from  $g(\mathcal{N}^c)$  in the latent space  $\mathcal{H}$ . Intuitively, this vector points in a direction of the latent space where the presence of the concept  $c \in [C]$  increases. When  $f_k(\mathbf{x}) = l_k \circ g(\mathbf{x})$  corresponds to the predicted probability of class  $k \in [d_Y]$  for  $\mathbf{x} \in \mathcal{X}$ , this insight allows us to define a class-wise conceptual sensitivity. Indeed, the directional derivative

---

<sup>1</sup>Note that we use bold symbols for vectors.

<sup>2</sup>It is understood that there exists a dictionary between the integer concept identifier and the concept’s name. In this way, if the first concept of interest is *stripes*, we can write  $c = 1$  and  $c = \text{Stripes}$ .

$S_k^c(\mathbf{x}) \equiv (\mathbf{w}^c)^\top \nabla_{\mathbf{h}} l_k[\mathbf{g}(\mathbf{x})]$  measures the extent to which the probability of class  $k$  varies in the direction of the CAV. If the sensitivity is positive  $S_k^c(\mathbf{x}) > 0$ , this indicates that the presence of concept  $c$  increases the belief of the model  $\mathbf{f}$  that  $k$  is the correct class for example  $\mathbf{x}$ . Note that CAV sensitivities  $S_k^c(\mathbf{x})$  can be generalized to aggregate the gradients between  $\mathbf{g}(\mathbf{x})$  and a baseline  $\bar{\mathbf{h}} \in \mathcal{H}$  [32]. Given a dataset of examples split by class  $\mathcal{D} = \bigsqcup_{k=1}^{d_Y} \mathcal{D}_k$ , where  $\mathcal{D}_k$  contains only examples of class  $k \in [d_Y]$ , it is possible to summarize the overall sensitivity of class  $k$  to concept  $c$  by computing the score  $\text{TCAV}_k^c = |\{\mathbf{x} \in \mathcal{D}_k | S_k^c(\mathbf{x}) > 0\}| / |\mathcal{D}_k|$ , where  $|\cdot|$  denotes the set cardinality.

**The limitations of linear separability.** The linear separability of concept negatives and positives is central in the aforementioned formalism. Indeed, the existence of a separating hyperplane in  $\mathcal{H}$  is necessary to define CAVs, which in turn are required to define concept sensitivity and TCAV scores. This assumption has been criticized in the literature [39, 40]. The main criticism is the following: there is no reason to expect the representation map  $\mathbf{g}$  to linearly separate concept positive and negatives since these concept-related labels are not known by the model. Let us consider a concrete example to stress that generic classifiers should not be expected to linearly separate concepts even if they linearly separate classes. Consider a classifier  $\mathbf{f} = \mathbf{l} \circ \mathbf{g}$ , where for each class  $k \in [d_Y]$ :  $l_k(\mathbf{h}) = \text{Softmax}(\alpha_k^\top \mathbf{h} + \beta_k)$  with  $\alpha_k \in \mathbb{R}^{d_H}$  and  $\beta_k \in \mathbb{R}$ . This parametrization is typically realized when the representation space  $\mathcal{H}$  corresponds to the penultimate layer of a neural network. For concreteness, we consider a computer vision setting with the classes  $k \in \{\text{Tiger}, \text{Lion}, \text{Cow}, \text{Zebra}\}$ . We note that the decision boundary in  $\mathcal{H}$  for any pair of class is a hyperplane. For instance, the decision boundary between the classes lion and tiger is parametrized by  $l_{\text{Lion}}(\mathbf{h}) = l_{\text{Tiger}}(\mathbf{h})$ , which corresponds to the equation of a hyperplane in  $\mathcal{H}$ :  $(\alpha_{\text{Lion}} - \alpha_{\text{Tiger}})^\top \mathbf{h} + (\beta_{\text{Lion}} - \beta_{\text{Tiger}}) = 0$ . In this setting, any successful classifier  $\mathbf{f}$  requires a representation map  $\mathbf{g}$  that linearly separates the classes in the latent space  $\mathcal{H}$ . We now consider the concept  $c = \text{Stripes}$ . In this case, we can expect examples from the classes tiger and zebra in  $\mathcal{P}^{\text{Stripes}}$  (as both tigers and zebras have stripes) and examples from the classes lion and cow in  $\mathcal{N}^{\text{Stripes}}$  (as neither lions nor cows have stripes). As illustrated in Figure 2, it is therefore perfectly possible to have  $\mathbf{g}(\mathcal{P}^{\text{Stripes}})$  and  $\mathbf{g}(\mathcal{N}^{\text{Stripes}})$  that are not linearly separable in spite of the classes linear separability. With this example, we emphasize the subtle distinction between classes and concept linear separability. While the former is expected and holds experimentally [33], this is not the case for the latter.

**A better assumption.** Although the representations from Figure 2 do not linearly separate the concept sets, we note that concept positive and negative examples are scattered across distinct clusters. In this way, a model that produces these representations appears to make a difference between presence and absence of the concept. From this angle, we could consider that the concept is well encoded in the representation space geometry. To formalize this more general notion of concept set separability, we adapt the smoothness assumption originally formulated in semi-supervised learning [46].

**Assumption 2.1** (Concept Smoothness). A concept  $c \in [C]$  is encoded in the latent space  $\mathcal{H}$  if  $\mathcal{H}$  is smooth with respect to the concept. This means that we can separate  $\mathcal{H} = \mathcal{H}^c \sqcup \mathcal{H}^{\neg c}$  into a *concept activation region* (CAR)  $\mathcal{H}^c$  where the concept  $c$  is mostly present (i.e.  $|\mathbf{g}(\mathcal{P}^c) \cap \mathcal{H}^c| \gg |\mathbf{g}(\mathcal{N}^c) \cap \mathcal{H}^c|$ ) and a region  $\mathcal{H}^{\neg c}$  where the concept  $c$  is mostly absent (i.e.  $|\mathbf{g}(\mathcal{N}^c) \cap \mathcal{H}^{\neg c}| \gg |\mathbf{g}(\mathcal{P}^c) \cap \mathcal{H}^{\neg c}|$ ). If two points  $\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}$  in a high-density region of the latent space are close to each other, then we should have  $\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}^c$  or  $\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}^{\neg c}$ .

*Remark 2.1.* We can easily see that linear separability is trivially included in this assumption: it corresponds to the case where the CAR  $\mathcal{H}^c$  and  $\mathcal{H}^{\neg c}$  are separated by a hyperplane. Conversely, our assumption does not require the CAR  $\mathcal{H}^c$  and  $\mathcal{H}^{\neg c}$  to be separated by a hyperplane for a concept to be relevant, as illustrated in Figure 3.

There are two crucial components in the previous assumption that need to be detailed: what do we mean by *density* and how do we extract a CAR  $\mathcal{H}^c$  from  $\mathcal{H}$ .

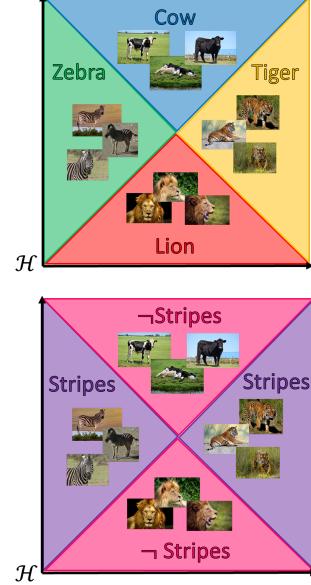


Figure 2: Classes are linearly separable (top) but concept sets are not (bottom).

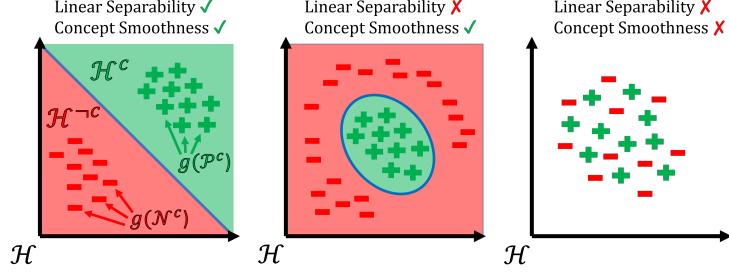


Figure 3: Linear separability implies concept smoothness but not the opposite.

## 2.2 Detecting Concepts

**Concept Density.** Let us start by detailing the notion of density that we use. We first note that the notion of proximity in  $\mathcal{H}$  is crucial in Assumption 2.1. It is conventionally formalized through a kernel function  $\kappa : \mathcal{H}^2 \rightarrow \mathbb{R}^+$  [47]. These functions are such that the proximity between  $\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}$  increases with  $\kappa(\mathbf{h}_1, \mathbf{h}_2)$ . Similarly, the proximity between  $\mathbf{h}$  and a discrete set  $\mathcal{S} \subset \mathcal{H}$  of representations increases with  $\sum_{\mathbf{h}' \in \mathcal{S}} \kappa(\mathbf{h}, \mathbf{h}')$ . This last sum can be interpreted as the density of examples from  $\mathcal{S}$  at a point  $\mathbf{h} \in \mathcal{H}$ . In our setup, it is natural to define a density relative to each concept  $c \in [C]$  by using the two sets  $\mathcal{P}^c$  and  $\mathcal{N}^c$ . This motivates the following definition.

**Definition 2.1** (Concept Density). Let  $\kappa : \mathcal{H}^2 \rightarrow \mathbb{R}^+$  be a kernel function. For each concept  $c \in [C]$ , we assume that we have a positive set  $\mathcal{P}^c = \{\mathbf{x}^{c,n} \mid n \in [N^c]\}$  and a negative set  $\mathcal{N}^c = \{\mathbf{x}^{-c,n} \mid n \in [N^c]\}$ . We define the concept density as a function  $\rho^c : \mathcal{H} \rightarrow \mathbb{R}$  such that

$$\begin{aligned} \rho^c(\mathbf{h}) &= \rho^{\mathcal{P}^c}(\mathbf{h}) - \rho^{\mathcal{N}^c}(\mathbf{h}), \\ \rho^{\mathcal{P}^c}(\mathbf{h}) &= \frac{1}{N^c} \sum_{n=1}^{N^c} \kappa[\mathbf{h}, \mathbf{g}(\mathbf{x}^{c,n})], \quad \rho^{\mathcal{N}^c}(\mathbf{h}) = \frac{1}{N^c} \sum_{n=1}^{N^c} \kappa[\mathbf{h}, \mathbf{g}(\mathbf{x}^{-c,n})]. \end{aligned}$$

The concept density for an example  $\mathbf{x} \in \mathcal{X}$  can similarly be defined as  $\rho^c[\mathbf{g}(\mathbf{x})]$ .

*Remark 2.2.* This density function is not necessarily positive. Indeed, we have assigned a positive contribution for examples from  $\mathcal{P}^c$  and a negative one for those of  $\mathcal{N}^c$ . The idea is that  $\rho^c(\mathbf{h}) > 0$  whenever the density of  $\mathbf{g}(\mathcal{P}^c)$  is higher around  $\mathbf{h}$ . Conversely,  $\rho^c(\mathbf{h}) < 0$  whenever the density of  $\mathbf{g}(\mathcal{N}^c)$  is higher around  $\mathbf{h}$ . Finally,  $\rho^c(\mathbf{h}) \approx 0$  if  $\mathbf{h}$  is isolated from positive and negative examples or if the density of  $\mathbf{g}(\mathcal{P}^c)$  balances the density of  $\mathbf{g}(\mathcal{N}^c)$  around  $\mathbf{h}$ .

**Concept Activation Regions.** With this definition, it is tempting to use the density  $\rho^c$  to define the regions  $\mathcal{H}^c$  and  $\mathcal{H}^{-c}$ . A natural choice would be to define the CAR as the positive density region  $\mathcal{H}^c = (\rho^c)^{-1}(\mathbb{R}^+)$  and its complementary as the negative density region  $\mathcal{H}^{-c} = (\rho^c)^{-1}(\mathbb{R}^-)$ . This corresponds to using a Parzen window classifier for the concept [48]. An obvious limitation of this approach is that each evaluation of the density  $\rho^c$  scales linearly with the number  $N^c$  of concept examples. If the size of concept sets is large, it is possible to obtain a sparse version of this Parzen window classifier with a *support vector classifier* (SVC) [49, 50]  $s_\kappa^c : \mathcal{H} \rightarrow \{0, 1\}$ . This SVC is fitted to discriminate the concept sets  $\mathbf{g}(\mathcal{P}^c)$  and  $\mathbf{g}(\mathcal{N}^c)$ . We can then define the CARs as  $\mathcal{H}^c = (s_\kappa^c)^{-1}(1)$  and  $\mathcal{H}^{-c} = (s_\kappa^c)^{-1}(0)$ . More details can be found in Appendix A.

**Global explanations.** Our CAR formalism permits to extend those local (i.e. sample-wise) considerations globally. Let us start by defining the equivalent of TCAV scores described in Section 2.1. This score aims at understanding how the model relates classes with concepts. We define the TCAR score as the fraction of examples that have class  $k$  and whose representation lies in the CAR  $\mathcal{H}^c$ :  $\text{TCAR}_k^c = |\mathbf{g}(\mathcal{D}_k) \cap \mathcal{H}^c| / |\mathcal{D}_k|$ . Note that  $\text{TCAR}_k^c \in [0, 1]$ , where 0 corresponds to no overlap and 1 to a full overlap. In Appendix B, we extend approach to measure the overlap between two concepts.

**Latent space isometries invariance.** In many applications such as clustering [51] or data visualization [52], the only relevant geometrical information of the representation space  $\mathcal{H}$  is the distance  $\|\mathbf{h}_1 - \mathbf{h}_2\|_{\mathcal{H}}$  between every pair of points  $(\mathbf{h}_1, \mathbf{h}_2) \in \mathcal{H}^2$ . Informally, we say that two representation spaces are isometric if they assign the same distance to each pair of points. In the aforementioned applications, two isometric representation spaces are therefore indistinguishable from one another.

Since concept-based explanations similarly describe the representation space geometry, one might require similar invariance to hold in this context. In Appendix D, we show that our CAR formalism provides such guarantee if  $\kappa$  is a radial kernel [53]. To the best of our knowledge, this type of analysis has not been performed in the context of concept-based explanation methods. We believe that future works in this domain would greatly benefit from this type of insight.

### 2.3 Concepts and Features

Not all features are relevant to identify a concept. As an example, let us consider the concept  $c = \text{Stripes}$  in a computer vision setting. If the image  $\mathbf{x}$  represents a zebra, only some small portion of the image will exhibit stripes (namely the body of the zebra). Therefore, if we build a saliency map for the concept  $c$ , we expect only this part of the image to be relevant for the identification of this concept. Existing works to produce concept-level saliency maps relying on the CAV formalism exist in the literature [34, 54]. In the absence of CAVs, we need an alternative approach. We now describe how any feature importance method can be used in conjunction with CARs. A generic feature importance method assigns a score  $a_i(f, \mathbf{x}) \in \mathbb{R}$  to each feature  $i \in [d_X]$  for a *scalar* model  $f : \mathcal{X} \rightarrow \mathbb{R}$  to make a prediction  $f(\mathbf{x})$  [55]. Since our purpose is to measure the relevance of each feature in identifying a concept  $c \in [C]$ , we compute the feature importance for the concept density:  $a_i(\rho^c \circ \mathbf{g}, \mathbf{x})$ . Note that some specific feature importance methods, such as Integrated Gradients [25] or Gradient Shap [24], require the explained model to be differentiable. This is the case whenever the kernel  $\kappa$  and the latent representation  $\mathbf{g}$  are differentiable. We note that the support vector classifier  $s_\kappa^c$  cannot be used in this context due to its non-differentiability. In Appendix C, we show that our CAR-based feature importance can be endowed with a useful completeness property that relates the importance scores  $a_i(\rho^c \circ \mathbf{g}, \mathbf{x})$  to the concept density  $\rho^c \circ \mathbf{g}(\mathbf{x})$ .

## 3 Experiments

The code to reproduce all the experiments from this section is available at <https://github.com/JonathanCrabbe/CARs> and <https://github.com/vanderschaarlab/CARs>.

### 3.1 Empirical Evaluation

Our purpose is to empirically validate the formalism described in the previous section. We have several independent components to evaluate: ① the concept classifier used to detect the CARs  $\mathcal{H}^c$ , ② the global explanations induced by the TCAR values and ③ the feature importance scores induced by the concept densities  $\rho^c$ .

**Datasets.** We perform our experiments on 3 datasets. ① The MNIST dataset [56] consists of  $28 \times 28$  grayscale images, each representing a digit. We train a convolutional neural network (CNN) with 2 layers to identify the digit of each image. ② The MIT-BIH Electrocardiogram (ECG) dataset [57, 58] consists of univariate time series with 187 time steps, each representing a heartbeat cycle. We train a CNN with 3 layers to determine whether each heartbeat is normal or abnormal. ③ The Caltech-UCSD Birds-200 (CUB) dataset [59] consists of coloured images of various sizes, each representing a bird from one of the 200 species present in the dataset. We fine-tune an Inceptionv3 neural network [60] to identify the species each bird belongs to among the 200 possible choices.

**Concepts.** For each dataset, we study the models through the lens of several well-defined concepts that are provided by human annotations. ① For MNIST: we use  $C = 4$  concepts that correspond to simple geometrical attributes of the images: *Loop* (positive images include a loop), *Vertical/Horizontal Line* (positive images include a vertical/horizontal line) and *Curvature* (positive images contain segments that are not straight lines). ② For ECG: we use  $C = 4$  concepts defined by cardiologists to better characterize abnormal heartbeats [61]: *Premature Ventricular*, *Supraventricular*, *Fusion Beats* and *Unknown*. The exact definition for each of these concepts is beyond the scope of this paper. We simply note that annotations for those concepts are available in the ECG dataset. ③ For CUB: we use  $C = 112$  concepts that correspond to visual attributes of the birds (e.g. their size, the colour of their wings, etc.). We use the same procedure as [41] to extract those concepts from the CUB dataset. We stress that, in each case, we selected concepts that can be unambiguously associated to the classes that are predicted by the models. Hence, it is reasonable to expect those concepts to be salient for the

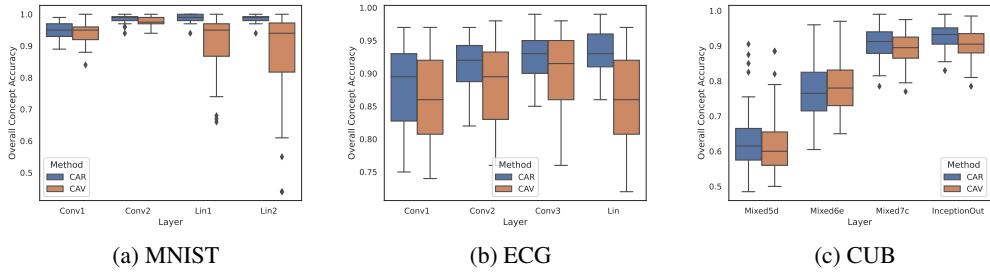


Figure 4: Overall accuracy of concept classifiers.

models. In each case, we sample the positive and negative sets  $\mathcal{P}^c$  and  $\mathcal{N}^c$  from the model’s training sets. For more details on the concepts and the models, please refer to Appendix E.

### 3.1.1 Accuracy of concept activation regions

**Methodology.** The purpose of this experiment is to assess if the concept regions  $\mathcal{H}^c$  identified by our CAR classifier generalize well to unseen examples. Each of the models described above are endowed with several representation spaces (one per hidden layer). For several of those latent spaces, we fit our CAR classifier (SVC with radial basis function kernel) to discriminate the concept sets  $\mathcal{P}^c, \mathcal{N}^c$  for each concept  $c \in [C]$ . These two sets have a size  $N^c = 200$  and are sampled from the model’s training set. The classifier is then evaluated by computing its accuracy on a holdout balanced concept set  $\mathcal{T}^c$  of size 100 sampled from the model’s testing set. For MNIST and ECG, we repeat this experiment 10 times for each concept and let the sets  $\mathcal{P}^c, \mathcal{N}^c, \mathcal{T}^c$  vary on each run. For comparison, we perform the same experiment with a linear CAV classifier as a benchmark. We report the overall (all the concepts together) accuracy in Figure 4.

**Analysis.** The CAR classifier substantially outperforms the CAV classifier. Note that this advantage is even more striking in the representation spaces associated to the deeper (last) DNN layers. This can be better understood through the lens of Cover’s theorem [62]: the linear separation underlying CAV is usually easier to achieve in higher dimensional spaces, which corresponds to the shallower (first) DNN layers in this case. When the dimension  $d_H$  of the latent space becomes comparable to the size  $N^c$  of the concept sets, linear separation often fails to maintain a high accuracy. By contrast, the SVC underlying CARs manage to maintain high accuracy through the more flexible notion of concept smoothness defined in Assumption 2.1. This suggests that concepts can be well encoded in the geometry of the latent space even when accurate linear separability is not possible. We also note that the accuracy CAR classifiers seems to increase with the representation’s depth. This is consistent with the behaviour of class probes [33].

**Statistical significance** The statistical significance of the concept classifiers is evaluated with the permutation test from [32]. All of them are statistically significant with p-value < .05 except for some concepts classifiers that are fitted with the layers *Mixed5d* and *Mixed6e* of the CUB Inceptionv3 model. We note that those classifiers do not generalize well in Figure 4c. This suggests that deeper networks are required to identify more challenging concepts correctly.

**Take-away 1:** CAR classifiers better capture how concepts are spread across representation spaces.

### 3.1.2 Consistency of global explanations

**Methodology.** The purpose of this experiment is to assess if the aforementioned concept classifiers create meaningful global association between classes and concepts. For each model, we now focus our study on the penultimate layer. We compute the TCAV and TCAR scores for each (class, concept) pair over the whole testing set. Ideally, these scores should be correlated to the true proportion (computed with the human concept annotations) of examples from the class that exhibit the concept. To that aim, we report the Pearson correlation  $r$  between each score and this true proportion in Table 1. To illustrate those global explanations, we also provide the scores for one class of each dataset in Figure 5 (for CUB, we restrict to wing colour concepts, other examples are reported in Appendix K).

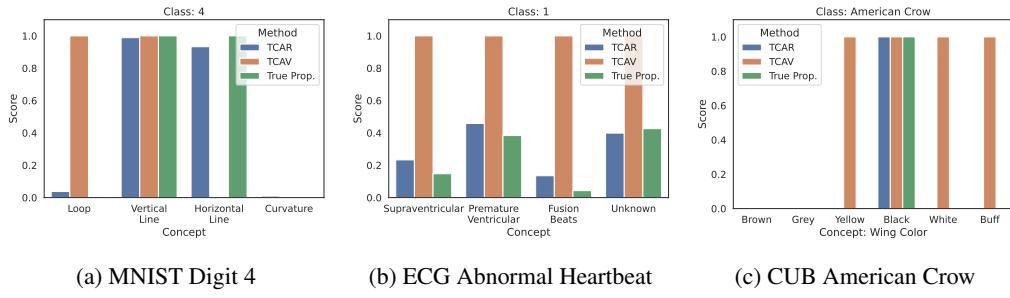


Figure 5: Examples of global concept-based explanations.

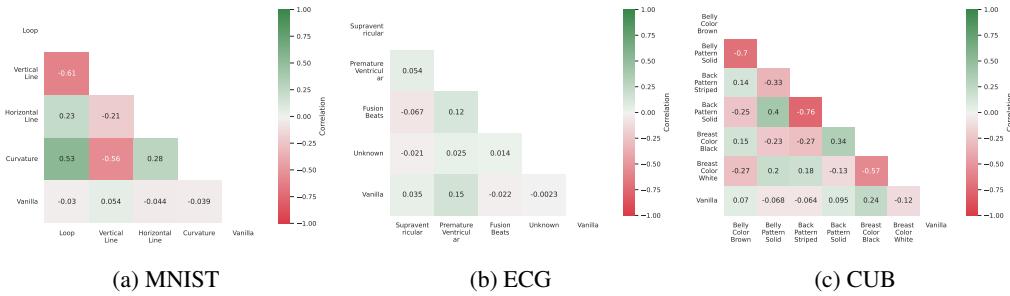


Figure 6: Correlations between concept feature importance.

**Analysis.** The TCAR scores better correlate with the true presence of concepts. This difference can be understood by looking at the examples from Figure 5. We note that TCAV scores tend to predict nonexistent associations (e.g. yellow wings for American crows) and miss existing associations (e.g. horizontal lines for digit 4). For ECG, we note that TCAR does capture the fact that some concepts (like fusion beats) are less represented within the class, while TCAV does not. In all of these cases, TCAV explanations might give the impression that the model did not learn meaningful class-concept associations. The TCAR analysis leads to the opposite conclusion. Since we have established that TCAR is built upon more accurate concept classifiers, it seems that models indeed learn concepts as intended, in spite of what TCAV explanations suggest.

**Take-away 2:** TCAR scores more faithfully reflect the true association between classes and concepts.

### 3.1.3 Coherency of concept-based feature importance

**Methodology.** Focusing again on the model’s penultimate layer, we are now interested in feature importance. The evaluation of feature importance methods is a notoriously difficult problem [63]. To perform an analysis in line with our concept-based explanations, we propose to check that our concept-based feature importance meaningfully captures features that are concept-specific. We analyse this through 2 desiderata: ① concept-based feature importance is not equivalent to the model (vanilla) feature importance and ② only concepts that can be identified with similar features should yield correlated feature importance. To evaluate these desiderata empirically, we compute our CAR-based Integrated Gradients  $a_i(\rho^c \circ g, x)$  for each concept  $c \in [C]$  and for each example  $x \in \mathcal{D}_{\text{test}}$  from the test set. Similarly, we compute the vanilla Integrated Gradients  $a_i(f, x)$  for the model<sup>3</sup>. We quantitatively compare these various feature importance scores by computing their Pearson correlation  $r$  as in [64, 65]. We report the results in Figure 6 (again, we selected a subset of 6 concepts for CUB).

Table 1: Evaluation of global explanations.

Dataset	$r(\text{TCAR}, \text{TrueProp.})$	$r(\text{TCAV}, \text{TrueProp.})$
MNIST	<b>1.00</b>	.70
ECG	.97	.74
CUB	<b>.86</b>	.52

<sup>3</sup>The vanilla Integrated Gradients are computed for the model estimation of the true class probability.

**Analysis.** First, we observe that CAR-based feature importance correlates weakly with vanilla feature importance ( $|r| < .25$  for all datasets). This confirms that desideratum ① is fulfilled. Then, we note that most of the CAR-based feature importance scores are decorrelated or weakly correlated with each other. The counterexamples that we observe indeed correspond to concepts that can be identified with similar features. A first example is the positive correlation between the loop and the curvature concepts for MNIST ( $r = .53$ ), both concepts are generally associated to curved symbols. Another example is the negative correlation between striped and solid back patterns for CUB birds ( $r = -.76$ ), those concepts are mutually exclusive and identified by inspecting the bird’s back. Those examples support that CAR-based feature importance fulfils desideratum ②.

**Take-away 3:** CAR-based feature importance is concept-specific and captures concept associations.

### 3.2 Use Case: Machine Learning Model Rediscovering Known Medical Concepts

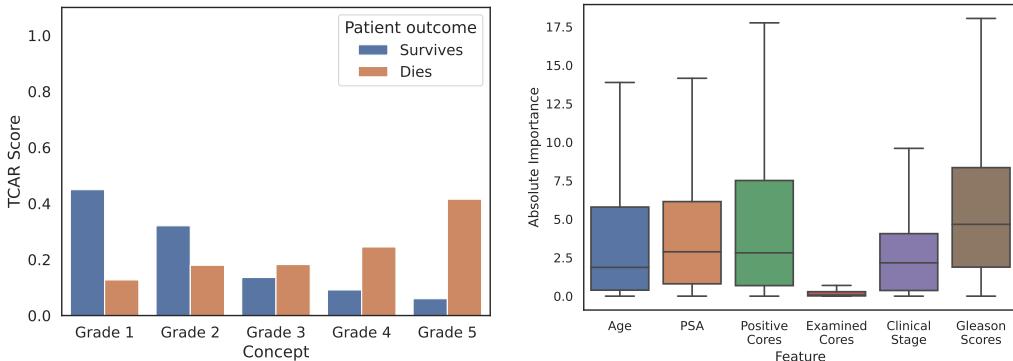
We will now describe a use case of the CAR formalism introduced in this paper. We stress that the literature already contains numerous use cases of CAV concept-based explanations, especially in the medical setting [66, 67, 37]. Since CARs generalize CAVs, it goes without saying that they apply to these use cases. Rather than repeating existing usage of concept-based explanations, we discuss an alternative use case motivated by recent trends in machine learning. With the successes of deep models in various scientific domains [5–8], we witness an increasing overlap between scientific discovery and machine learning. While this new trend opens up fascinating opportunities, it comes with a set of new challenges. The evaluation of machine learning models is arguably one of the most important of these challenges. In a scientific context, the canonical machine learning approach to validate models (out-of-sample generalization) is likely to be insufficient. Beyond generalization on unseen data, the scientific validity of a model requires consistency with established scientific knowledge [68]. We propose to illustrate how our CARs can be used in this context.

**Dataset.** We use the data collected with the Surveillance, Epidemiology, and End Results (SEER) Program. The dataset [69] contains a US population-based cohort of 171,942 men diagnosed with non-metastatic prostate cancer between Jan 1, 2000, and Dec 31, 2016. Each patient is described by age, the results of a prostate-specific antigen blood test (PSA), the clinical stage of its tumour and two Gleason scores (primary and secondary). Each patient also has a label that indicates if they died because of their prostate cancer. We train a multilayer perceptron (MLP) to predict the patient’s mortality on 90% of the data and test on the remaining 10%. For a more detailed description of the data and the model, please refer to Appendix F.

**Concepts.** Doctors use an established grading system to predict how likely the cancer is to spread [70]. This system assigns to each patient a grade between 1 and 5. The probability that the cancer spreads increases with this grade. It can be computed from the two Gleason scores (more details in Appendix F). We can consider each of these 5 grades as a concept. We note that this grade is not explicitly part of the input features of the MLP that we trained. Our purpose is to assess if the MLP implicitly discovered those grades in order to predict the patient’s mortality. If this happens to be the case, this would demonstrate that the MLP is in-line with existing medical knowledge.

**Methodology.** As in the previous section, we are going to use the output of the MLP’s penultimate layer as a representation space. We fit a CAR classifier (SVC with linear kernel) to discriminate the concepts sets  $\mathcal{P}^c, \mathcal{N}^c$  for each grade  $c \in [5]$ . Both of those sets contain  $N^c = 250$  patients sampled from the training set. We then proceed to several verifications. ① We determine if the grades have been learned by measuring the accuracy of each CAR classifier on a holdout balanced concept set  $\mathcal{T}^c$  of size 100 sampled from the model’s testing set. ② We check that each grade is associated to the appropriate outcome by reporting the TCAR scores between each pair (mortality, grade) in Figure 7a. ③ We verify that grades are identified with the right features by plotting a summary of the absolute feature importance scores  $|a_i(\rho^c \circ \mathbf{g}, \mathbf{x})|$  on the test set  $\mathbf{x} \in \mathcal{D}_{\text{test}}$  in Figure 7b.

**Analysis.** Let us summarize the findings for each of the above points. ① All of the CAR classifiers generalize well on the test set (their accuracy ranges from 90% to 100%). This strongly suggests that the MLP implicitly separates patients with different grades in representation space. ② Figure 7a demonstrates that the model associates higher grades with higher mortality. This is in line with the clinical interpretation of the grades. ③ Figure 7b suggests that the Gleason scores constitute the most important features overall (highest quantiles) for the model to discriminate between grades. This is consistent with the fact that grades are computed with the Gleason scores. We conclude that the MLP



(a) TCAR grades vs mortality.

(b) Grades feature importance.

Figure 7: CAR explanations with the prostate cancer grades as concepts.

implicitly identifies the patient’s grades with high accuracy and in a way that is consistent with the medical literature.

**Take-away 4:** The CAR formalism can reliably support scientific evaluation of a model.

## 4 Conclusion

We introduced Concept Activation Regions, a new framework to relax the linear separability assumption underlying the Concept Activation Vector formalism. We showed that our framework guarantees crucial properties, such as invariance with respect to latent symmetries. Through extensive validation on several datasets, we verified that ① Concept Activation Regions better capture the distribution of concepts across the model’s representation space, ② The resulting global explanations are more consistent with human annotations and ③ Concept Activation Regions permit to define concept-specific feature importance that is consistent with human intuition. Finally, through a use case involving prostate cancer data, we show the neural network can implicitly rediscover known scientific concepts, such as the prostate cancer grading system.

Many important points that were not covered in the main paper can be found in the appendices. In Appendix A, we discuss how to tune the various hyperparameters of the CAR classifiers. In Appendix B, we explain how to generalize concept activation vectors to nonlinear decision boundaries. In Appendix G, we show that the CAR explanations are robust to adversarial perturbations and background shifts. In Appendix H, we demonstrate that CAR explanations can be used to relate abstract concepts discovered by self-explaining neural networks with human concepts. Finally, Appendix I illustrates how CAR explanations allow us to probe language models.

We believe that our extended concept explainability framework opens up many interesting avenues for future work. A first one would be to probe state of the art neural networks with an approach similar to Section 3. In particular, it would be interesting to analyse if improving model performance is associated with a better encoding of human concepts. A second one would be to analyze how concept discovery [38] can benefit from our generalized notion of concept activation. A more fine-grained characterization of a model’s latent space is likely to improve the surfaced concept. A third one, as suggested by Section 3.2, would be to use concept-based explanations to make a better scientific assessment of neural networks. Indeed, consistency with well-established knowledge is crucial for a scientific model to be accepted.

## Acknowledgments and Disclosure of Funding

The authors are grateful to Fergus Imrie, Yangming Li and the 3 anonymous NeurIPS reviewers for their useful comments on an earlier version of the manuscript. Jonathan Crabbé is funded by Aviva and Mihaela van der Schaar by the Office of Naval Research (ONR), NSF 172251.

## References

- [1] Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in deep learning based natural language processing. *IEEE Computational Intelligence Magazine*, 13(3):55–75, 2018.
- [2] K. R. Chowdhary. Natural Language Processing. *Fundamentals of Artificial Intelligence*, pages 603–649, 2020.
- [3] Richard Szeliski. *Computer vision: algorithms and applications*. Springer Science & Business Media, 2010.
- [4] Athanasios Voulodimos, Nikolaos Doulamis, Anastasios Doulamis, and Eftychios Protopapadakis. Deep Learning for Computer Vision: A Brief Review. *Computational Intelligence and Neuroscience*, 2018.
- [5] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with AlphaFold. *Nature* 2021 596:7873, 596(7873):583–589, 2021.
- [6] Alex Davies, Petar Veličković, Lars Buesing, Sam Blackwell, Daniel Zheng, Nenad Tomašev, Richard Tanburn, Peter Battaglia, Charles Blundell, András Juhász, Marc Lackenby, Geordie Williamson, Demis Hassabis, and Pushmeet Kohli. Advancing mathematics by guiding human intuition with AI. *Nature* 2021 600:7887, 600(7887):70–74, 2021.
- [7] James Kirkpatrick, Brendan McMorrow, David H.P. Turban, Alexander L. Gaunt, James S. Spencer, Alexander G.D.G. Matthews, Annette Obika, Louis Thiry, Meire Fortunato, David Pfau, Lara Román Castellanos, Stig Petersen, Alexander W.R. Nelson, Pushmeet Kohli, Paula Mori-Sánchez, Demis Hassabis, and Aron J. Cohen. Pushing the frontiers of density functionals by solving the fractional electron problem. *Science*, 374(6573):1385–1389, 2021.
- [8] Jonas Degrave, Federico Felici, Jonas Buchli, Michael Neunert, Brendan Tracey, Francesco Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego de las Casas, Craig Donner, Leslie Fritz, Cristian Galperti, Andrea Huber, James Keeling, Maria Tsipoukelli, Jackie Kay, Antoine Merle, Jean Marc Moret, Sébastien Noury, Federico Pesamosca, David Pfau, Olivier Sauter, Cristian Sommariva, Stefano Coda, Basil Duval, Ambrogio Fasoli, Pushmeet Kohli, Koray Kavukcuoglu, Demis Hassabis, and Martin Riedmiller. Magnetic control of tokamak plasmas through deep reinforcement learning. *Nature* 2022 602:7897, 602(7897):414–419, 2022.
- [9] Zachary C Lipton. The Mythos of Model Interpretability. *Communications of the ACM*, 61(10):35–43, 2016.
- [10] Finale Doshi-Velez and Been Kim. Towards A Rigorous Science of Interpretable Machine Learning. *arXiv*, 2017.
- [11] Travers Ching, Daniel S. Himmelstein, Brett K. Beaulieu-Jones, Alexandr A. Kalinin, Brian T. Do, Gregory P. Way, Enrico Ferrero, Paul-Michael Agapow, Michael Zietz, Michael M. Hoffman, Wei Xie, Gail L. Rosen, Benjamin J. Lengerich, Johnny Israeli, Jack Lanchantin, Stephen Woloszynek, Anne E. Carpenter, Avanti Shrikumar, Jinbo Xu, Evan M. Cofer, Christopher A. Lavender, Srinivas C. Turaga, Amr M. Alexandari, Zhiyong Lu, David J. Harris, Dave DeCaprio, Yanjun Qi, Anshul Kundaje, Yifan Peng, Laura K. Wiley, Marwin H. S. Segler, Simina M. Boca, S. Joshua Swamidass, Austin Huang, Anthony Gitter, and Casey S. Greene. Opportunities and obstacles for deep learning in biology and medicine. *Journal of The Royal Society Interface*, 15(141):20170387, 2018.
- [12] Sana Tonekaboni, Shalmali Joshi, Melissa D Mccradden, and Anna Goldenberg. What Clinicians Want: Contextualizing Explainable Machine Learning for Clinical End Use. *Proceedings of Machine Learning Research*, 2019.

- [13] Markus Langer, Daniel Oster, Timo Speith, Holger Hermanns, Lena Kästner, Eva Schmidt, Andreas Sesing, and Kevin Baum. What do we want from Explainable Artificial Intelligence (XAI)? – A stakeholder perspective on XAI and a conceptual model guiding interdisciplinary XAI research. *Artificial Intelligence*, 296:103473, 2021.
- [14] Amina Adadi and Mohammed Berrada. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access*, 6:52138–52160, 2018.
- [15] Arun Das and Paul Rad. Opportunities and Challenges in Explainable Artificial Intelligence (XAI): A Survey. *arXiv*, 2020.
- [16] Erico Tjoa and Cuntai Guan. A Survey on Explainable Artificial Intelligence (XAI): Toward Medical XAI. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–21, 2020.
- [17] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador Garcia, Sergio Gil-Lopez, Daniel Molina, Richard Benjamins, Raja Chatila, and Francisco Herrera. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58:82–115, 2020.
- [18] Edward Choi, Mohammad Taha Bahadori, Joshua A. Kulas, Andy Schuetz, Walter F. Stewart, and Jimeng Sun. RETAIN: An Interpretable Predictive Model for Healthcare using Reverse Time Attention Mechanism. *Advances in Neural Information Processing Systems*, pages 3512–3520, 2016.
- [19] Chaofan Chen, Oscar Li, Chaofan Tao, Alina Jade Barnett, Jonathan Su, and Cynthia Rudin. This Looks Like That: Deep Learning for Interpretable Image Recognition. *Advances in Neural Information Processing Systems*, 32, 2018.
- [20] Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, Klaus Robert Müller, and Wojciech Samek. Layer-wise Relevance Propagation for Neural Networks with Local Renormalization Layers. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 9887 LNCS:63–71, 2016.
- [21] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. Association for Computing Machinery, 2016.
- [22] Ruth C Fong and Andrea Vedaldi. Interpretable Explanations of Black Boxes by Meaningful Perturbation. *Proceedings of the IEEE International Conference on Computer Vision*, pages 3449–3457, 2017.
- [23] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning Important Features Through Propagating Activation Differences. *34th International Conference on Machine Learning, ICML 2017*, 7:4844–4866, 2017.
- [24] Scott Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 4768–4777. Neural information processing systems foundation, 2017.
- [25] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic Attribution for Deep Networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 3319–3328. International Machine Learning Society (IMLS), 2017.
- [26] Jonathan Crabbé and Mihaela van der Schaar. Explaining Time Series Predictions with Dynamic Masks. In *International Conference on Machine Learning*, 2021.
- [27] Pang Wei Koh and Percy Liang. Understanding Black-box Predictions via Influence Functions. In *International Conference on Machine Learning*, volume 34, pages 2976–2987. International Machine Learning Society (IMLS), 2017.

- [28] Amirata Ghorbani and James Zou. Data Shapley: Equitable Valuation of Data for Machine Learning. *36th International Conference on Machine Learning, ICML 2019*, pages 4053–4065, 2019.
- [29] Garima Pruthi, Frederick Liu, Mukund Sundararajan, and Satyen Kale. Estimating Training Data Influence by Tracing Gradient Descent. In *Advances in Neural Information Processing Systems*, pages 19920–19930. Curran Associates, Inc., 2020.
- [30] Jonathan Crabbé, Zhaozhi Qian, Fergus Imrie, and Mihaela van der Schaar. Explaining Latent Representations with a Corpus of Examples. In *Advances in Neural Information Processing Systems*, 2021.
- [31] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, and Rory Sayres. Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV). *35th International Conference on Machine Learning, ICML 2018*, 6:4186–4195, 2017.
- [32] Jessica Schrouff, Sébastien Baur, Shaobo Hou, Diana Mincu, Eric Lureau, Ralph Blanes, James Wexler, Alan Karthikesalingam, and Been Kim. Best of both worlds: local and global explanations with human-understandable concepts. *arXiv*, 2021.
- [33] Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. In *International Conference on Learning Representations*, 2017.
- [34] Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Interpretable basis decomposition for visual explanation. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 11212 LNCS, pages 122–138, 2018.
- [35] Chih-Kuan Yeh, Been Kim, Sercan Ö Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar. On Completeness-aware Concept-Based Explanations in Deep Neural Networks. *Advances in Neural Information Processing Systems*, 33:20554–20565, 2020.
- [36] Ferdinand Kusters, Peter Schichtel, Sheraz Ahmed, and Andreas Dengel. Conceptual Explanations of Neural Network Prediction for Time Series. *Proceedings of the International Joint Conference on Neural Networks*, 2020.
- [37] Diana Mincu, Eric Lureau, Shaobo Hou, Sébastien Baur, Ivan Protsyuk, Martin G Seneviratne, Anne Mottram, Nenad Tomasev, Alan Karthikesalingam, and Jessica Schrouff. Concept-based model explanations for Electronic Health Records. *ACM CHIL 2021 - Proceedings of the 2021 ACM Conference on Health, Inference, and Learning*, pages 36–46, 2020.
- [38] Amirata Ghorbani, James Wexler, James Zou, and Been Kim. Towards Automatic Concept-based Explanations. *Advances in Neural Information Processing Systems*, 32, 2019.
- [39] Zhi Chen, Yijie Bei, and Cynthia Rudin. Concept whitening for interpretable image recognition. *Nature Machine Intelligence 2020 2:12*, 2(12):772–782, 2020.
- [40] Rahul Soni, Naresh Shah, Chua Tat Seng, and Jimmy D Moore. Adversarial TCAV – Robust and Effective Interpretation of Intermediate Layers in Neural Networks. *arXiv*, 2020.
- [41] Pang Wei Koh, Thao Nguyê, Yew Siang Tang, Stephen Mussmann, Emma Pierso, Been Kim, and Percy Liang. Concept Bottleneck Models. In *37th International Conference on Machine Learning, ICML 2020*, pages 5294–5304. PMLR, 2020.
- [42] Pietro Barbiero, Gabriele Ciravegna, Francesco Giannini, Pietro Liò, Marco Gori, and Stefano Melacci. Entropy-based Logic Explanations of Neural Networks. *arXiv*, 2021.
- [43] Dobrik Georgiev, Pietro Barbiero, Dmitry Kazhdan, Petar Veličković, and Pietro Liò. Algorithmic Concept-based Explainable Reasoning. *arXiv*, 2021.
- [44] Andrei Margeloiu, Matthew Ashman, Umang Bhatt, Yanzhi Chen, Mateja Jamnik, and Adrian Weller. Do Concept Bottleneck Models Learn as Intended? In *ICLR Workshop on Responsible AI*, 2021.

- [45] Anita Mahinpei, Justin Clark, Isaac Lage, Finale Doshi-Velez, and Weiwei Pan. Promises and Pitfalls of Black-Box Concept Learning Models. *arXiv*, 2021.
- [46] Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. *Semi-Supervised Learning*. The MIT Press, 2006.
- [47] Thomas Hofmann, Bernhard Schölkopf, and Alexander J Smola. Kernel methods in machine learning. *The Annals of Statistics*, 36(3):1171–1220, 2008.
- [48] Emanuel Parzen. On Estimation of a Probability Density Function and Mode. *The Annals of Mathematical Statistics*, 33(3):1065–1076, 1962.
- [49] Kristin P Bennett and O L Mangasarian. Robust linear programming discrimination of two linearly inseparable sets. *Optimization Methods and Software*, 1(1):23–34, 1992.
- [50] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [51] Lior Rokach and Oded Maimon. Clustering methods. In *Data mining and knowledge discovery handbook*, pages 321–352. Springer, 2005.
- [52] Laurens van der Maaten and Geoffrey Hinton. Visualizing Data using t-SNE. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008.
- [53] Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- [54] Lennart Brocki and Neo Christopher Chung. Concept saliency maps to visualize relevant features in deep generative models. *Proceedings - 18th IEEE International Conference on Machine Learning and Applications, ICMLA 2019*, pages 1771–1778, 2019.
- [55] Ian Covert, Scott Lundberg, and Su-In Lee. Explaining by Removing: A Unified Framework for Model Explanation. *Journal of Machine Learning Research*, 22(209):1–90, 2021.
- [56] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2323, 1998.
- [57] A. L. Goldberger, L. A. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C. K. Peng, and H. E. Stanley. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. *Circulation*, 101(23), 2000.
- [58] G.B. Moody and R.G. Mark. The impact of the mit-bih arrhythmia database. *IEEE Engineering in Medicine and Biology Magazine*, 20(3):45–50, 2001.
- [59] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.
- [60] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 07-12-June-2015:1–9, 2015.
- [61] Mohammad Kachuee, Shayan Fazeli, and Majid Sarrafzadeh. ECG Heartbeat Classification: A Deep Transferable Representation. *Proceedings - 2018 IEEE International Conference on Healthcare Informatics, ICHI 2018*, pages 443–444, 2018.
- [62] Thomas M. Cover. Geometrical and Statistical Properties of Systems of Linear Inequalities with Applications in Pattern Recognition. *IEEE Transactions on Electronic Computers*, EC-14(3):326–334, 1965.
- [63] Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. The (un)reliability of saliency methods. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pages 267–280. Springer, 2019.

- [64] Olivier Le Meur and Thierry Baccino. Methods for comparing scanpaths and saliency maps: Strengths and weaknesses. *Behavior Research Methods*, 45(1):251–266, 2013.
- [65] Jonathan Crabbé and Mihaela van der Schaar. Label-Free Explainability for Unsupervised Models. *arXiv*, 2022.
- [66] Mara Graziani, Vincent Andrearzyk, and Henning Müller. Regression Concept Vectors for Bidirectional Explanations in Histopathology. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 11038 LNCS:124–132, 2019.
- [67] James R. Clough, Ilkay Oksuz, Esther Puyol-Antón, Bram Ruijsink, Andrew P. King, and Julia A. Schnabel. Global and Local Interpretability for Cardiac MRI Classification. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 11767 LNCS:656–664, 2019.
- [68] Ribana Roscher, Bastian Bohn, Marco F. Duarte, and Jochen Garske. Explainable Machine Learning for Scientific Insights and Discoveries. *IEEE Access*, 8:42200–42216, 2020.
- [69] Margaret Peggy Adamo, Jessica A. Boten, Linda M. Coyle, Kathleen A. Cronin, Clara J.K. Lam, Serban Negoita, Lynne Penberthy, Jennifer L. Stevens, and Kevin C. Ward. Validation of prostate-specific antigen laboratory values recorded in Surveillance, Epidemiology, and End Results registries. *Cancer*, 123(4):697–703, 2017.
- [70] Jennifer Gordetsky and Jonathan Epstein. Grading of prostatic adenocarcinoma: current state and prognostic implications. *Diagnostic Pathology*, 11(1), 2016.
- [71] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [72] Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. Captum: A unified and generic model interpretability library for PyTorch. *arXiv*, 2020.
- [73] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.
- [74] Pascal Sturmels, Scott Lundberg, and Su-In Lee. Visualizing the Impact of Feature Attribution Baselines. *Distill*, 5(1):e22, 2020.
- [75] Christopher M Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006.
- [76] Guido Van Rossum and Fred L Drake Jr. *Python reference manual*. Centrum voor Wiskunde en Informatica Amsterdam, 1995.
- [77] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [78] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. SMOTE: Synthetic Minority Over-sampling Technique. *Journal Of Artificial Intelligence Research*, 16:321–357, 2011.
- [79] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. *37th International Conference on Machine Learning*, pages 1575–1585, 2020.

- [80] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1452–1464, 2018.
- [81] David Alvarez-Melis and Tommi S. Jaakkola. Towards Robust Interpretability with Self-Explaining Neural Networks. *Advances in Neural Information Processing Systems*, pages 7775–7784, 2018.
- [82] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. GloVe: Global Vectors for Word Representation. *EMNLP 2014 - 2014 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference*, pages 1532–1543, 2014.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes] All the claims from the abstract and Section 1 are verified empirically in Section 3. The theoretical claims are demonstrated in Appendices C and D.
  - (b) Did you describe the limitations of your work? [Yes] The assumption for our concept-based explanations to be valid are clearly stated in Assumption 2.1. Our method only applies to neural networks, which is stated in Section 2.
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes] Our work improves the transparency of deep neural networks. This has many beneficial societal impacts, as described in Section 1. We do not see any potential negative societal impact to this approach. We have carefully reviewed the points from the ethical guideline and none of the mentioned negative impact seems to apply to our work.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes] We have read the ethics review guidelines and confirm that our paper conforms to them.
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [Yes] The assumptions in Propositions C.1 and D.1 are clearly stated in Appendices C and D respectively.
  - (b) Did you include complete proofs of all theoretical results? [Yes] The proofs for Propositions C.1 and D.1 are given in Appendices C and D respectively.
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] The full code is available at <https://github.com/JonathanCrabbe/CARs> and <https://github.com/vanderschaarl/CARs>. The implementation of our method closely follows Algorithms 1, 2 and 4 in the appendices. All the details to reproduce the experimental results are given in Appendices E and F.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] All the training details are specified in Appendices E and F.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] Figure 4 measures the accuracy of our method over several runs. All the runs are aggregated together in the form of a box-plot.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] Our computing resources are described in Appendices E and F.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [Yes] We have included a citation for the MNIST dataset [56], for the ECG dataset [57], for the CUB dataset [59], for the SEER dataset [69] and for the InceptionV3 model [60].
  - (b) Did you mention the license of the assets? [Yes] The licenses are mentioned in Appendices E and F.
  - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A] We have used only existing assets.
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A] All the datasets that we use are publicly available.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A] The medical datasets we use are public and have been de-identified.
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

## A CAR Classifiers

In this appendix, we provide some details about our CAR classifiers.

**Implementation.** To determine the concept activation regions, we fit a SVC  $s_\kappa^c$  for each concept  $c \in [C]$ . This process is described in Algorithm 1.

---

**Algorithm 1:** Fit CAR Classifier

---

**Input:** Neural network  $f = l \circ g : \mathcal{X} \rightarrow \mathcal{H} \rightarrow \mathcal{Y}$ , kernel function  $\kappa : \mathcal{H}^2 \rightarrow \mathbb{R}^+$ , concept positives  $\mathcal{P}^c \subset \mathcal{X}$ , concept negatives  $\mathcal{N}^c \subset \mathcal{X}$

**Output:** CAR SVC Classifier  $s_\kappa^c : \mathcal{H} \rightarrow \{0, 1\}$

```

 $s_\kappa^c \leftarrow \text{SVC}(\kappa);$  /* Initialize SVC classifier */
 $\mathcal{D} \leftarrow \{(g(x), \mathbb{I}_{\mathcal{P}^c}(x)) \mid x \in \mathcal{P}^c \sqcup \mathcal{N}^c\};$  /* Assemble SVC training set */
 $(s_\kappa^c).\text{fit}(\mathcal{D});$  /* Fit SVC classifier */
return  $s_\kappa^c$ 

```

---

where  $\mathbb{I}_{\mathcal{P}^c}$  is the indicator function on the set  $\mathcal{P}^c$ . Our implementation leverages the SVC class from scikit-learn [71]. We use the default hyperparameters for each CAR classifier that appears in our experiments. The resulting classifier  $s_\kappa^c$  can then be used to assess if a test point  $x \in \mathcal{X}$  has a representation  $g(x) \in \mathcal{H}$  that falls in the CAR  $\mathcal{H}^c$  or not:

$$g(x) \in \begin{cases} \mathcal{H}^c & \text{if } s_\kappa^c \circ g(x) = 1 \\ \mathcal{H}^{\neg c} & \text{if } s_\kappa^c \circ g(x) = 0. \end{cases}$$

**Choosing the kernel.** Algorithm 1 requires the user to specify a kernel function  $\kappa$ . To select this kernel, a good approach is to train several SVCs  $s_\kappa^c$  with different kernels  $\kappa$  and see how each SVC generalizes on a validation set. We illustrate this process with MNIST in Figure 8. We note that the Gaussian RBF kernel outperforms the other kernels on the validation set, although the Matern kernel achieves perfect accuracy on the training set. This highlights the importance of evaluating a CAR classifier on a held-out dataset to make sure that the related CAR offers a good description of how concepts are distributed in the latent space  $\mathcal{H}$ . In our experiments, we found that Gaussian RBF kernels are often the most interesting option.

**Concept set size.** Algorithm 1 requires the user to specify concepts sets  $\mathcal{P}^c$  and  $\mathcal{N}^c$ . What size  $N^c = |\mathcal{P}^c| = |\mathcal{N}^c|$  should we choose for these concept sets? It seems logical that larger concepts sets are more likely to yield more accurate CARs. To study this experimentally, we propose to fit several concept classifiers for MNIST by varying the size  $N^c$  of their training concept sets  $\mathcal{P}^c$  and

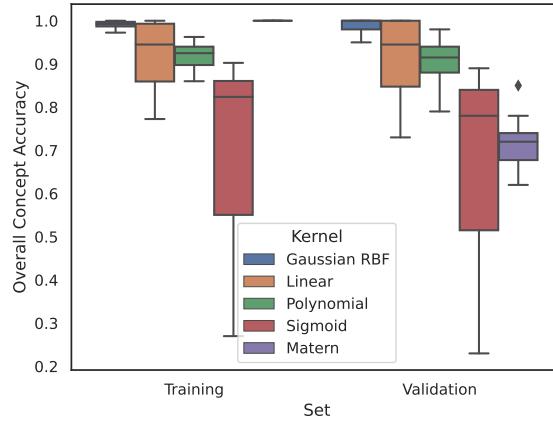


Figure 8: Accuracy of CAR classifiers for MNIST concepts with different kernels.

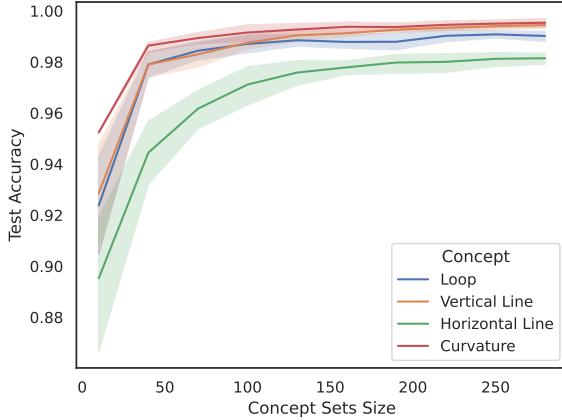


Figure 9: Accuracy of CAR classifiers for MNIST concepts for various concept set size  $N^c$  (average and 95% confidence intervals based on 10 runs).

$\mathcal{N}^c$ . We report the results in Figure 9. As we can see, the curves flatten above  $N^c > 200$ . Increasing the concept sets size beyond this point does not improve the accuracy of the resulting CAR classifier. We recommend to acquire concept examples until the performance of the CAR classifier stabilizes. In our experiments, we found that  $N^c = 200$  examples is often sufficient to obtain accurate CAR classifiers.

**Tuning hyperparameters.** In the case where the user desires a CAR classifier that generalizes as well as possible, tuning these hyperparameters might be useful. We propose to tune the kernel type, kernel width and error penalty of our CAR classifiers  $s_\kappa^c$  for each concept  $c \in [C]$  by using Bayesian optimization and a validation concept set:

1. Randomly sample the hyperparameters from an initial prior distribution  $\theta_h \sim P_{\text{prior}}$ .
2. Split the concept sets  $\mathcal{P}^c, \mathcal{N}^c$  into training concept sets  $\mathcal{P}_{\text{train}}^c, \mathcal{N}_{\text{train}}^c$  and validation concept sets  $\mathcal{P}_{\text{val}}^c, \mathcal{N}_{\text{val}}^c$ .
3. For the current value  $\theta_h$  of the hyperparameters, fit a model  $s_\kappa^c$  to discriminate the training concept sets  $\mathcal{P}_{\text{train}}^c, \mathcal{N}_{\text{train}}^c$ .
4. Measure the accuracy  $\text{ACC}_{\text{val}} = \frac{\sum_{x \in \mathcal{P}_{\text{val}}^c} \mathbf{1}(s_\kappa^c \circ g(x)=1) + \sum_{x \in \mathcal{N}_{\text{val}}^c} \mathbf{1}(s_\kappa^c \circ g(x)=0)}{|\mathcal{P}_{\text{val}}^c \cup \mathcal{N}_{\text{val}}^c|}$ .
5. Update the current hyperparameters  $\theta_h$  based on  $\text{ACC}_{\text{val}}$  using Bayesian optimization (Optuna in our case).
6. Repeat 3-5 for a predetermined number of trials.

We applied this process to the CAR accuracy experiment (same setup as in Section 3.1.1 of the main paper) to tune the CAR classifiers for the CUB concepts. Interestingly, we noticed no improvement with respect to the CAR classifiers reported in the main paper: tuned and standard CAR classifier have an average accuracy of  $(93 \pm .2)\%$  for the penultimate Inception layer. This suggests that the accuracy of CAR classifiers is not heavily dependant on hyperparameters in this case.

## B TCAR Global Explanations

In this appendix, we provide some details about the TCAR scores.

**Implementation.** When the CAR classifiers are available, they permit to compute TCAR scores through Algorithm 2.

**TCAR between concepts.** Up until now, we have discussed TCAR scores that indicate how models relate classes to concepts. It is possible to define a similar score to estimate how models relate two

---

**Algorithm 2:** Compute Class-Concept TCAR Score

---

**Input:** Neural network  $f = l \circ g : \mathcal{X} \rightarrow \mathcal{H} \rightarrow \mathcal{Y}$ , CAR classifier  $s_\kappa^c : \mathcal{H} \rightarrow \{0, 1\}$ , set of examples  $\mathcal{D}_k \subset \mathcal{X}$  of a given class  $k \in [d_Y]$

**Output:** TCAR score  $\text{TCAR}_k^c \in [0, 1]$

```

count ← 0 ;                                /* Initialize concept positive count */
for  $x \in \mathcal{D}_k$  do
  if  $s_\kappa^c \circ g(x) = 1$  then
    | count ← count + 1 ;           /* Increment count if example in the CAR */
    | end
end
TCAR $k$  $c$  ← count/ $|\mathcal{D}_k|$  ;          /* Normalize count to get TCAR */
return  $\text{TCAR}_k^c$ 

```

---

concepts with each other. Given a set  $\mathcal{D} \subset \mathcal{X}$  of examples, we define the TCAR score associated to the concepts  $c_1, c_2 \in [C]$  as the ratio

$$\text{TCAR}^{c_1, c_2} = \frac{|\mathbf{g}(\mathcal{D}) \cap \mathcal{H}^{c_1} \cap \mathcal{H}^{c_2}|}{|\mathbf{g}(\mathcal{D}) \cap (\mathcal{H}^{c_1} \cup \mathcal{H}^{c_2})|}.$$

Again,  $\text{TCAR} = 0$  corresponds to no overlap and  $\text{TCAR} = 1$  describes a perfect overlap. We note that the concept-concept TCAR score can be interpreted as a Jaccard index between the sets  $\mathbf{g}(\mathcal{D}) \cap \mathcal{H}^{c_1}$  and  $\mathbf{g}(\mathcal{D}) \cap \mathcal{H}^{c_2}$ . This score is symmetric with respect to the concepts:  $\text{TCAR}^{c_1, c_2} = \text{TCAR}^{c_2, c_1}$ . The computation of this score is done as in Algorithm 3.

---

**Algorithm 3:** Compute Concept-Concept TCAR Score

---

**Input:** Neural network  $f = l \circ g : \mathcal{X} \rightarrow \mathcal{H} \rightarrow \mathcal{Y}$ , CAR classifiers  $s_\kappa^{c_1} : \mathcal{H} \rightarrow \{0, 1\}$  and  $s_\kappa^{c_2} : \mathcal{H} \rightarrow \{0, 1\}$ , set of examples  $\mathcal{D} \subset \mathcal{X}$

**Output:** TCAR score  $\text{TCAR}^{c_1, c_2} \in [0, 1]$

```

numerator ← 0 ;                                /* Initialize score numerator */
denominator ← 0 ;                             /* Initialize score denominator */
for  $x \in \mathcal{D}$  do
  if  $s_\kappa^{c_1} \circ g(x) = 1$  and  $s_\kappa^{c_2} \circ g(x) = 1$  then
    | numerator ← numerator + 1 ;           /* Increment if example in both CARs */
    | end
  if  $s_\kappa^{c_1} \circ g(x) = 1$  or  $s_\kappa^{c_2} \circ g(x) = 1$  then
    | denominator ← denominator + 1 ;       /* Increment if example in one CAR */
    | end
end
TCAR $c_1, c_2$  ← numerator/denominator ;          /* Assemble the TCAR score */
return  $\text{TCAR}^{c_1, c_2}$ 

```

---

**TCAR between MNIST concepts.** As an illustration, we compute concept-concept TCAR scores for the MNIST concepts and report the results in Figure 10. We see that the model relates concepts that tend to appear together (e.g. curvature and loop). Hence, concept-concept TCAR scores can serve as a proxy for the concept semantics encoded in a model’s representation space. We note the similarity with the correlation between concept-based feature importance illustrated in Figure 6. The main difference is that concept-concept TCAR scores do not explicitly refer to input features.

**Generalizing concept sensitivity.** In our formalism, it is perfectly possible to define a *local* concept activation vector through the concept density  $\rho^c : \mathcal{H} \rightarrow \mathbb{R}^+$  defined in Definition 2.1. Indeed, the vector  $\nabla_{\mathbf{h}}\rho^c[\mathbf{h}] \in \mathcal{H}$  points in the direction of the representation space  $\mathcal{H}$  where the concept density (and hence the presence of the concept) increases. Hence, this vector can be interpreted as a *local* concept activation vector. Note that this vector becomes global whenever we parametrize the concept density  $\rho^c$  with a linear kernel  $\kappa(\mathbf{h}_1, \mathbf{h}_2) = \mathbf{h}_1^\top \mathbf{h}_2$ . Equipped with this generalized notion of concept activation vector, we can also generalize the CAV concept sensitivity  $S_k^c$  by replacing the CAV  $w^c$  by  $\nabla_{\mathbf{h}}\rho^c[\mathbf{h}]$  for the representation  $\mathbf{h} = \mathbf{g}(\mathbf{x})$  of the input  $\mathbf{x} \in \mathcal{X}$ :

$$S_k^c(\mathbf{x}) \equiv (\nabla_{\mathbf{h}}\rho^c[\mathbf{g}(\mathbf{x})])^\top (\nabla_{\mathbf{h}}l_k[\mathbf{g}(\mathbf{x})]).$$

In this way, all the interpretation provided by the CAV formalism are also available in the CAR formalism.

## C CAR Feature Importance

In this appendix, we provide some details about our concept-based feature importance.

**Implementation.** Concept-based feature importance uses the concept densities  $\rho^c$  to attribute an importance score to each feature in order to confirm/reject the presence of a concept  $c \in [C]$  for a given example  $\mathbf{x} \in \mathcal{X}$ . This process is described in Algorithm 4. We stress that features importance scores are computed with *concept densities* and not with SVC classifiers. The reason for this is that SVCs are non-differentiable functions of the input and are therefore incompatible with gradient-based attribution methods such as Integrated-Gradients [25]. One could argue that a sparse version of the density  $\rho^c$  could be used by only allowing support vectors from the SVC to contribute. Our first implementation was relying on this approach. Unfortunately, this often led to vanishing importance scores. A possible explanation is the following: Gaussian radial basis function kernels  $\kappa(\mathbf{h}_1, \mathbf{h}_2) = \exp[-(\gamma \|\mathbf{h}_1 - \mathbf{h}_2\|_{\mathcal{H}})^2]$  decay very quickly as  $\|\mathbf{h}_1 - \mathbf{h}_2\|_{\mathcal{H}}$  increases. Hence, unless the example  $\mathbf{x}$  we wish to explain has a representation  $\mathbf{g}(\mathbf{x})$  located near a support vector of the SVC classifier, the sparse density (and its gradients) vanishes. On the other hand, using the density from Definition 2.1 allows us to incorporate the representations of all concept examples. This makes it more likely that the representation  $\mathbf{g}(\mathbf{x})$  is located near (at least) one of the representation that contributes to the density. This permits to solve the vanishing gradient problem that led to vanishing attributions. It goes without saying that this solution comes at the expense of having to compute a density that scales linearly with the size  $N^c$  of the concept sets. In our implementation, we make this tractable by restricting to small concept sets ( $N^c \leq 250$ ) and by saving the concept sets representations  $\mathbf{g}(\mathcal{P}^c)$  and  $\mathbf{g}(\mathcal{N}^c)$  to limit the number of queries to the model. We use Captum’s implementation of feature importance methods [72]. When using radial basis function kernels, we tune the kernel width  $\gamma$

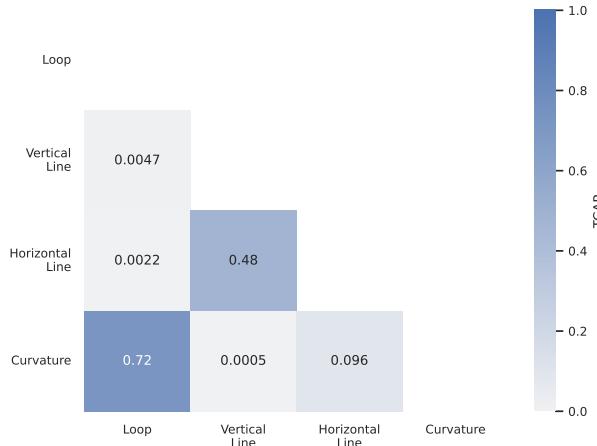


Figure 10: TCAR score between MNIST concepts.

with Optuna [73] with 1,000 trials so that the Parzen window classifier  $\mathbb{I}_{\mathbb{R}^+} \circ \rho^c$ , where  $\mathbb{I}_{\mathbb{R}^+}$  denotes the indicator function on  $\mathbb{R}^+$ , accurately discriminates between positives representations  $\mathbf{g}(\mathcal{P}^c)$  and negative representations  $\mathbf{g}(\mathcal{N}^c)$ .

---

**Algorithm 4:** Compute concept-based feature importance

---

**Input:** Neural network  $\mathbf{f} = \mathbf{l} \circ \mathbf{g} : \mathcal{X} \rightarrow \mathcal{H} \rightarrow \mathcal{Y}$ , kernel function  $\kappa : \mathcal{H}^2 \rightarrow \mathbb{R}^+$ , concept positives  $\mathcal{P}^c \subset \mathcal{X}$ , concept negatives  $\mathcal{N}^c \subset \mathcal{X}$ , feature importance method  $\mathbf{a} : \mathcal{A}(\mathbb{R}^{\mathcal{X}}) \times \mathcal{X} \rightarrow \mathbb{R}^{d_X}$ , example  $\mathbf{x} \in \mathcal{X}$

**Output:** Feature importance vector  $\mathbf{a} \in \mathbb{R}^{d_X}$  for the example  $\mathbf{x}$

```

 $\rho^c \leftarrow \text{Density}(\mathbf{g}, \kappa, \mathcal{P}^c, \mathcal{N}^c); \quad /* \text{ Initialize concept density as Def. 2.1 */}$ 
 $\mathbf{a} \leftarrow \mathbf{a}(\rho^c \circ \mathbf{g}, \mathbf{x}); \quad /* \text{ Compute feature importance vector */}$ 
return  $\mathbf{a}$ 

```

---

where  $\mathcal{A}(\mathbb{R}^{\mathcal{X}})$  denotes the hypothesis set of scalar functions on the input space  $\mathcal{X}$ . In this case, this corresponds to the set of neural networks with input space  $\mathcal{X}$  and scalar output.

**Completeness.** Completeness is a crucial property of some feature importance methods. It guarantees that the feature importance scores can be used to reconstruct the model prediction. To make this more precise, let  $f : \mathcal{X} \rightarrow \mathbb{R}$  be a model with *scalar* output and let  $\mathbf{x} \in \mathcal{X}$  be an example we wish to explain. Feature importance methods assign a score  $a_i(f, \mathbf{x})$  to each feature  $i \in [d_X]$ . This score reflects the sensitivity of the model prediction  $f(\mathbf{x})$  with respect to the component  $x_i$  of  $\mathbf{x}$ . Completeness is fulfilled whenever the sum of this importance scores equals the model prediction up to a constant baseline  $b \in \mathbb{R}$ :  $\sum_{i=1}^{d_X} a_i(f, \mathbf{x}) = f(\mathbf{x}) - b$ . The baseline varies from one method to the other. For instance, Lime [21] uses a vanishing baseline  $b = 0$ , SHAP [24] uses the average prediction  $b = \mathbb{E}_{\mathbf{X}}[f(\mathbf{X})]$  and Integrated Gradients [25] use a baseline prediction  $b = f(\bar{\mathbf{x}})$  for some baseline input  $\bar{\mathbf{x}}$ . With completeness, the importance scores  $a_i(f, \mathbf{x})$  are given a natural interpretation: their sign indicates if the features tend to increase/decrease the prediction  $f(\mathbf{x})$  and their absolute value indicates how important this effect is. When combined with CAR concept densities  $\rho^c$ , this interpretation is even more insightful.

**Proposition C.1** (CAR Completeness). *Consider a neural network decomposed as  $\mathbf{f} = \mathbf{l} \circ \mathbf{g}$ , where  $\mathbf{g} : \mathcal{X} \rightarrow \mathcal{H}$  is a feature extractor mapping the input space  $\mathcal{X}$  to the representation space  $\mathcal{H}$  and  $\mathbf{l}$  is a label function that maps the representation space  $\mathcal{H}$  to the label set  $\mathcal{Y}$ . Let  $\rho^c : \mathcal{H} \rightarrow \mathbb{R}^+$  be a concept density for some concept  $c \in [C]$ , defined as in Definition 2.1. Let  $\mathbf{a} : \mathcal{A}(\mathbb{R}^{\mathcal{X}}) \times \mathcal{X} \rightarrow \mathbb{R}$  be a feature importance method satisfying the completeness property:  $\sum_{i=1}^{d_X} a_i(\mathbf{f}, \mathbf{x}) = f(\mathbf{x}) - b$  for some baseline  $b \in \mathbb{R}$  and for all  $\mathbf{x} \in \mathcal{X}$ . Then, we have the following completeness property for the concept-based density:*

$$\sum_{i=1}^{d_X} a_i(\rho^c \circ \mathbf{g}, \mathbf{x}) = \rho^c \circ \mathbf{g}(\mathbf{x}) - b$$

*Remark C.1.* With this property, we can interpret features  $i \in [d_X]$  with  $a_i(\rho^c \circ \mathbf{g}, \mathbf{x}) > 0$  as those that tend to increase the concept density. This means that those features are important for the feature extractor  $\mathbf{g}$  to map the example in a region of the representation space  $\mathcal{H}$  where the concept is present. Hence, those are features that are important to identify a given concept  $c \in [C]$ . Conversely, features  $i \in [d_X]$  with  $a_i(\rho^c \circ \mathbf{g}, \mathbf{x}) < 0$  tend to decrease the concept density and therefore brings the example  $\mathbf{x}$  in a region of the representation space  $\mathcal{H}$  where the concept is absent. These features can therefore be interpreted as important to reject the presence of a given concept  $c \in [C]$ .

*Proof.* We simply note that  $\rho^c \circ \mathbf{g}$  is a scalar function as  $\rho^c \circ \mathbf{g}(\mathbf{x}) \in \mathbb{R}^+$  for all  $\mathbf{x} \in \mathcal{X}$ . Hence, the proposition immediately follows by applying the completeness property to  $f = \rho^c \circ \mathbf{g}$ .  $\square$

**Input baseline choice.** The choice of baseline input  $\bar{\mathbf{x}}$  has a notable effect on feature importance methods [74]. What constitutes a good input baseline is problem dependant. Intuitively,  $\bar{\mathbf{x}}$  should correspond to an input  $\mathbf{x} \in \mathcal{X}$  where no information is present [55]. Let us now explain how this information removal is achieved with the datasets that we use in our experiments. ① For MNIST,

we chose a black image  $\bar{x} = \mathbf{0}$  as an input baseline. This is because MNIST images have a black background and the information comes from white pixels that represent the digits. ② For ECG, we chose a constant  $\bar{x} = \mathbf{0}$  time series as an input baseline. This is because ECG time series are normalized ( $x_t \in [0, 1]$  for all  $t \in [d_X]$ ) and the pulse information comes from time steps where the time series is non-vanishing  $x_t \neq 0$ . ③ For CUB, choosing a baseline is more complicated. The reason for this is that the background colour changes from one image to the other and rarely corresponds to a single colour. To address this, we proceed as in the literature [22] and select a baseline  $\bar{x}$  that corresponds to a blurred version of the image we wish to explain:  $\bar{x}(x) = \mathbf{G}_\sigma \otimes x$ , where  $\mathbf{G}_\sigma$  is a Gaussian filter of width  $\sigma$  and  $\otimes$  denotes the convolution operation. We note that this baseline depends on which input  $x \in \mathcal{X}$  we want to explain. In our implementation, we use  $\sigma = 50$  to have images that are significantly blurred. ④ For SEER, we chose a constant vector  $\bar{x} = \mathbf{0}$  as an input baseline. This is because all continuous features are standardized and all categorical features are one-hot encoded.

## D CAR Latent Isometry Invariance

In this appendix, we prove that CAR explanations are invariant under isometries of the latent space when built with a radial kernel. Let us first rigorously define the notion of isometry between two vector spaces.

**Definition D.1** (Isometry). Let  $(\mathcal{H}, \|\cdot\|_{\mathcal{H}})$  and  $(\mathcal{H}', \|\cdot\|_{\mathcal{H}'})$  be two normed vector spaces. An *isometry* from  $\mathcal{H}$  to  $\mathcal{H}'$  is a map  $\tau : \mathcal{H} \rightarrow \mathcal{H}'$  such that for all  $\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}$ :

$$\|\tau(\mathbf{h}_1) - \tau(\mathbf{h}_2)\|_{\mathcal{H}'} = \|\mathbf{h}_1 - \mathbf{h}_2\|_{\mathcal{H}}.$$

We say that the two spaces  $(\mathcal{H}, \|\cdot\|_{\mathcal{H}})$  and  $(\mathcal{H}', \|\cdot\|_{\mathcal{H}'})$  are *isometric* if there exists a bijective isometry  $\tau$  from  $\mathcal{H}$  to  $\mathcal{H}'$ .

An explanation method is invariant to latent space isometries if applying a bijective isometry  $\tau$  to the model's latent space  $\mathcal{H}$  does not affect the explanations produced by the method. To make this more formal, we write the model as  $f = l \circ g = l \circ \tau^{-1} \circ \tau \circ g$ , where  $\tau^{-1}$  is the inverse of the bijective isometry  $\tau$ . In this setup, we could produce explanations with CARs by making the following replacements for the feature extractor and the label map:  $g \xrightarrow{\tau} \tau \circ g$  and  $l \xrightarrow{\tau} l \circ \tau^{-1}$ . The explanations are defined to be invariant to latent space isometries if they are unaffected by this replacement. It is legitimate to expect this since the previous replacement leads to the same model  $f \xrightarrow{\tau} f$  and substitutes the latent space by an isometric latent space  $\mathcal{H} \xrightarrow{\tau} \mathcal{H}' = \tau(\mathcal{H})$ .

Let us now discuss the isometry invariance of CAR explanations. First, we recall that CARs are defined through a kernel  $\kappa$ . Not all kernels lead to isometry invariant CAR explanations. We will show that it holds for a family of kernels known as radial kernels.

**Definition D.2** (Radial Kernel). A *radial kernel* is a kernel function  $\kappa : \mathcal{H}^2 \rightarrow \mathbb{R}^+$  that can be written as

$$\kappa(\mathbf{h}_1, \mathbf{h}_2) = \chi(\|\mathbf{h}_1 - \mathbf{h}_2\|_{\mathcal{H}}), \quad (1)$$

with a function  $\chi : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ .

A typical example of radial kernel is the Gaussian radial basis function kernel (RBF) that corresponds to  $\chi(x) = \exp[-(\gamma x)^2]$  for some  $\gamma \in \mathbb{R}^+$ . We are now ready to state and prove the isometry invariance property of CAR explanations.

**Proposition D.1** (CAR Isometry Invariance). Consider a neural network decomposed as  $f = l \circ g$ , where  $g : \mathcal{X} \rightarrow \mathcal{H}$  is a feature extractor mapping the input space  $\mathcal{X}$  to the representation space  $\mathcal{H}$  and  $l$  is a label function that maps the representation space  $\mathcal{H}$  to the label set  $\mathcal{Y}$ . Let  $\kappa : \mathcal{H}^2 \rightarrow \mathbb{R}^+$  be a radial kernel  $\kappa(\mathbf{h}_1, \mathbf{h}_2) = \chi(\|\mathbf{h}_1 - \mathbf{h}_2\|_{\mathcal{H}})$  that we use to define CAR's concept density in Definition 2.1. Let  $\tau : \mathcal{H} \rightarrow \mathcal{H}'$  be a bijective isometry between the normed spaces  $(\mathcal{H}, \|\cdot\|_{\mathcal{H}})$  and  $(\mathcal{H}', \|\cdot\|_{\mathcal{H}'})$ . All the explanations outputted by CAR remain the same if we transform the latent space with the isometry  $\tau$  by making the following replacements:  $g \xrightarrow{\tau} \tau \circ g$  and  $l \xrightarrow{\tau} l \circ \tau^{-1}$ .

*Proof.* For each concept  $c \in [C]$ , we note that CAR explanations exclusively rely on the concept density  $\rho^c$  from Definition 2.1 and the associated support vector classifier  $s_\kappa^c$ . Hence, it is sufficient to show that these two functions are invariant under isometry.

**Concept Density.** We start with the concept density  $\rho^c$ . We note that the kernel function  $\kappa$  can be applied to vectors from  $\mathcal{H}'$  as  $\kappa'(\mathbf{h}_1', \mathbf{h}_2') \equiv \chi(\|\mathbf{h}_1' - \mathbf{h}_2'\|_{\mathcal{H}'})$  for all  $(\mathbf{h}_1', \mathbf{h}_2') \in \mathcal{H}'$ . Applying CAR in the latent space  $\mathcal{H}'$  isometric to  $\mathcal{H}$  corresponds to using this kernel to compute an alternative density  $\rho'^c$ . Let us fix an example  $\mathbf{x} \in \mathcal{X}$ . Under the isometry, the concept density for this example transforms as  $\rho^c \circ \mathbf{g}(\mathbf{x}) \xrightarrow{\tau} \rho'^c \circ \tau \circ \mathbf{g}(\mathbf{x})$ . Let us show that this is in fact an invariance:

$$\begin{aligned} \rho'^c \circ \tau \circ \mathbf{g}(\mathbf{x}) &\stackrel{\text{Def. 2.1}}{=} \sum_{n=1}^{N^c} \kappa'[\tau \circ \mathbf{g}(\mathbf{x}), \tau \circ \mathbf{g}(\mathbf{x}^{c,n})] - \kappa'[\tau \circ \mathbf{g}(\mathbf{x}), \tau \circ \mathbf{g}(\mathbf{x}^{\neg c,n})] \\ &\stackrel{\text{Def. D.2}}{=} \sum_{n=1}^{N^c} \chi(\|\tau \circ \mathbf{g}(\mathbf{x}) - \tau \circ \mathbf{g}(\mathbf{x}^{c,n})\|_{\mathcal{H}'}) - \chi(\|\tau \circ \mathbf{g}(\mathbf{x}) - \tau \circ \mathbf{g}(\mathbf{x}^{\neg c,n})\|_{\mathcal{H}'}) \\ &\stackrel{\text{Def. D.1}}{=} \sum_{n=1}^{N^c} \chi(\|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}^{c,n})\|_{\mathcal{H}}) - \chi(\|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}^{\neg c,n})\|_{\mathcal{H}}) \\ &\stackrel{\text{Def. D.2}}{=} \sum_{n=1}^{N^c} \kappa[\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{x}^{c,n})] - \kappa[\mathbf{g}(\mathbf{x}), \mathbf{g}(\mathbf{x}^{\neg c,n})] \\ &\stackrel{\text{Def. 2.1}}{=} \rho^c \circ \mathbf{g}(\mathbf{x}). \end{aligned}$$

We deduce that the concept density is invariant under isometry:  $\rho^c \circ \mathbf{g}(\mathbf{x}) \xrightarrow{\tau} \rho^c \circ \mathbf{g}(\mathbf{x})$ .

**SVC.** The proof is more involved for the SVC  $s_\kappa^c$ . We assume that the reader is familiar with the standard theory of SVC. If this is not the case, please refer e.g. to Chapter 7 of [75]. For the sake of notation, we will abbreviate  $\mathbf{g}(\mathbf{x}^{c,n})$  and  $\mathbf{g}(\mathbf{x}^{\neg c,n})$  by  $\mathbf{h}^{c,n}$  and  $\mathbf{h}^{\neg c,n}$  respectively. Similarly, we abbreviate  $\tau \circ \mathbf{g}(\mathbf{x}^{c,n})$  and  $\tau \circ \mathbf{g}(\mathbf{x}^{\neg c,n})$  by  $\mathbf{h}'^{c,n}$  and  $\mathbf{h}'^{\neg c,n}$  respectively. The SVC concept classifier can be written as

$$s_\kappa^c(\mathbf{h}) = \mathbb{I}_{\mathbb{R}^+} \left( \sum_{n=1}^{N^c} \alpha^{c,n} \kappa[\mathbf{h}, \mathbf{h}^{c,n}] - \alpha^{\neg c,n} \kappa[\mathbf{h}, \mathbf{h}^{\neg c,n}] + \beta \right), \quad (2)$$

where  $\mathbb{I}_{\mathbb{R}^+}$  is the indicator function on  $\mathbb{R}^+$ , the real numbers  $\alpha^{c,n}, \alpha^{\neg c,n}$  maximize the objective

$$\begin{aligned} \mathcal{L}(\boldsymbol{\alpha}^c, \boldsymbol{\alpha}^{\neg c}) &= \sum_{n=1}^{N^c} \alpha^{c,n} + \alpha^{\neg c,n} - \frac{1}{2} \sum_{n=1}^{N^c} \sum_{m=1}^{N^c} \alpha^{c,n} \alpha^{c,m} \kappa[\mathbf{h}^{c,n}, \mathbf{h}^{c,m}] \\ &\quad + \alpha^{\neg c,n} \alpha^{\neg c,m} \kappa[\mathbf{h}^{\neg c,n}, \mathbf{h}^{\neg c,m}] - 2\alpha^{c,n} \alpha^{\neg c,m} \kappa[\mathbf{h}^{c,n}, \mathbf{h}^{\neg c,m}] \end{aligned} \quad (3)$$

under the constraints

$$\begin{aligned} \alpha^{c,n} &\geq 0 & \alpha^{\neg c,n} &\geq 0 & \forall n \in [N^c], \\ \sum_{n=1}^{N^c} \alpha^{c,n} - \alpha^{\neg c,n} &= 0. \end{aligned}$$

Finally, the bias term can be written as

$$\begin{aligned} \beta &= \frac{1}{|\mathcal{S}^c| + |\mathcal{S}^{\neg c}|} \left( |\mathcal{S}^c| - |\mathcal{S}^{\neg c}| - \sum_{n \in \mathcal{S}^c} \sum_{m \in \mathcal{S}^c} \alpha^{c,m} \kappa[\mathbf{h}^{c,n}, \mathbf{h}^{c,m}] \right. \\ &\quad + \sum_{n \in \mathcal{S}^c} \sum_{m \in \mathcal{S}^{\neg c}} \alpha^{\neg c,m} \kappa[\mathbf{h}^{c,n}, \mathbf{h}^{\neg c,m}] - \sum_{n \in \mathcal{S}^{\neg c}} \sum_{m \in \mathcal{S}^c} \alpha^{c,m} \kappa[\mathbf{h}^{\neg c,n}, \mathbf{h}^{c,m}] \\ &\quad \left. + \sum_{n \in \mathcal{S}^{\neg c}} \sum_{m \in \mathcal{S}^{\neg c}} \alpha^{\neg c,m} \kappa[\mathbf{h}^{\neg c,n}, \mathbf{h}^{\neg c,m}] \right), \end{aligned} \quad (4)$$

where  $\mathcal{S}^c = \{n \in [N^c] \mid \alpha^{c,n} \neq 0\}$  and  $\mathcal{S}^{\neg c} = \{n \in [N^c] \mid \alpha^{\neg c,n} \neq 0\}$  are the indices of the support vectors. Under isometry, the SVC classification for a latent vector  $\mathbf{h} \in \mathcal{H}$  transforms as  $s_\kappa^c(\mathbf{h}) \xrightarrow{\tau} s'_{\kappa'}^c(\mathbf{h}')$  with  $\mathbf{h}' = \tau(\mathbf{h})$  and

$$s'_{\kappa'}^c(\mathbf{h}') = \mathbb{I}_{\mathbb{R}^+} \left( \sum_{n=1}^{N^c} \underbrace{\alpha'^{c,n}}_{\textcircled{2}} \underbrace{\kappa'[\mathbf{h}', \mathbf{h}'^{c,n}]}_{\textcircled{1}} - \underbrace{\alpha'^{\neg c,n}}_{\textcircled{2}} \underbrace{\kappa'[\mathbf{h}', \mathbf{h}'^{\neg c,n}]}_{\textcircled{1}} + \underbrace{\beta'}_{\textcircled{3}} \right). \quad (5)$$

We will now show that  $s_\kappa^c(\mathbf{h}) = s_{\kappa'}^{c'}(\mathbf{h}')$  so that the SVC is invariant under isometries. We proceed in 3 steps.

① We show that  $\kappa'[\mathbf{h}'_1, \mathbf{h}'_2] = \kappa[\mathbf{h}_1, \mathbf{h}_2]$  for any  $(\mathbf{h}_1, \mathbf{h}_2) \in \mathcal{H}^2$  and  $\mathbf{h}'_1 = \boldsymbol{\tau}(\mathbf{h}_1), \mathbf{h}'_2 = \boldsymbol{\tau}(\mathbf{h}_2)$ :

$$\begin{aligned}\kappa'[\mathbf{h}'_1, \mathbf{h}'_2] &\stackrel{\text{Def. D.2}}{=} \chi(\|\boldsymbol{\tau}(\mathbf{h}_1) - \boldsymbol{\tau}(\mathbf{h}_2)\|_{\mathcal{H}'}) \\ &\stackrel{\text{Def. D.1}}{=} \chi(\|\mathbf{h}_1 - \mathbf{h}_2\|_{\mathcal{H}}) \\ &\stackrel{\text{Def. D.2}}{=} \kappa[\mathbf{h}_1, \mathbf{h}_2].\end{aligned}$$

By injecting this in (5), we are able to make the following replacements:  $\kappa'[\mathbf{h}', \mathbf{h}'^{c,n}] = \kappa[\mathbf{h}, \mathbf{h}^{c,n}]$  and  $\kappa'[\mathbf{h}', \mathbf{h}'^{\neg c,n}] = \kappa[\mathbf{h}, \mathbf{h}^{\neg c,n}]$  for all  $n \in [N^c]$ .

② We show that  $\alpha'^{c,n} = \alpha^{c,n}$  and  $\alpha'^{\neg c,n} = \alpha^{\neg c,n}$  for all  $n \in [N^c]$ . To that aim, we note that  $\boldsymbol{\alpha}'^c$  and  $\boldsymbol{\alpha}'^{\neg c}$  maximize the objective

$$\begin{aligned}\mathcal{L}'(\boldsymbol{\alpha}'^c, \boldsymbol{\alpha}'^{\neg c}) &= \sum_{n=1}^{N^c} \alpha'^{c,n} + \alpha'^{\neg c,n} - \frac{1}{2} \sum_{n=1}^{N^c} \sum_{m=1}^{N^c} \alpha'^{c,n} \alpha'^{c,m} \kappa'[\mathbf{h}'^{c,n}, \mathbf{h}'^{c,m}] \\ &\quad + \alpha'^{\neg c,n} \alpha'^{\neg c,m} \kappa'[\mathbf{h}'^{\neg c,n}, \mathbf{h}'^{\neg c,m}] - 2\alpha'^{c,n} \alpha'^{\neg c,m} \kappa'[\mathbf{h}'^{c,n}, \mathbf{h}'^{\neg c,m}] \\ &\stackrel{\text{①}}{=} \sum_{n=1}^{N^c} \alpha'^{c,n} + \alpha'^{\neg c,n} - \frac{1}{2} \sum_{n=1}^{N^c} \sum_{m=1}^{N^c} \alpha'^{c,n} \alpha'^{c,m} \kappa[\mathbf{h}^{c,n}, \mathbf{h}^{c,m}] \\ &\quad + \alpha'^{\neg c,n} \alpha'^{\neg c,m} \kappa[\mathbf{h}^{\neg c,n}, \mathbf{h}^{\neg c,m}] - 2\alpha'^{c,n} \alpha'^{\neg c,m} \kappa[\mathbf{h}^{c,n}, \mathbf{h}^{\neg c,m}] \\ &\stackrel{(3)}{=} \mathcal{L}(\boldsymbol{\alpha}'^c, \boldsymbol{\alpha}'^{\neg c}).\end{aligned}$$

Since this objective is identical to the one from the original SVC and the constraints are unaffected by the isometry, we deduce that the solution to this convex optimization problem is identical to the solution of (3). Hence, we have that  $\alpha'^{c,n} = \alpha^{c,n}$  and  $\alpha'^{\neg c,n} = \alpha^{\neg c,n}$  for all  $n \in [N^c]$ . Again, we can make these replacements in (5).

③ We show that  $\beta' = \beta$ . First, we note that the support vector indices are invariant under isometry:  $\mathcal{S}'^c = \{n \in [N^c] \mid \alpha'^{c,n} \neq 0\} = \{n \in [N^c] \mid \alpha^{c,n} \neq 0\} = \mathcal{S}^c$  and similarly  $\mathcal{S}'^{\neg c} = \mathcal{S}^{\neg c}$ . Hence we have

$$\begin{aligned}\beta' &= \frac{1}{|\mathcal{S}^c| + |\mathcal{S}^{\neg c}|} \left( |\mathcal{S}^c| - |\mathcal{S}^{\neg c}| - \sum_{n \in \mathcal{S}^c} \sum_{m \in \mathcal{S}^c} \alpha'^{c,m} \kappa'[\mathbf{h}'^{c,n}, \mathbf{h}'^{c,m}] \right. \\ &\quad + \sum_{n \in \mathcal{S}^c} \sum_{m \in \mathcal{S}^{\neg c}} \alpha'^{\neg c,m} \kappa'[\mathbf{h}'^{c,n}, \mathbf{h}'^{\neg c,m}] - \sum_{n \in \mathcal{S}^{\neg c}} \sum_{m \in \mathcal{S}^c} \alpha'^{c,m} \kappa'[\mathbf{h}'^{\neg c,n}, \mathbf{h}'^{c,m}] \\ &\quad \left. + \sum_{n \in \mathcal{S}^{\neg c}} \sum_{m \in \mathcal{S}^{\neg c}} \alpha'^{\neg c,m} \kappa'[\mathbf{h}'^{\neg c,n}, \mathbf{h}'^{\neg c,m}] \right) \\ &\stackrel{\text{①&②}}{=} \frac{1}{|\mathcal{S}^c| + |\mathcal{S}^{\neg c}|} \left( |\mathcal{S}^c| - |\mathcal{S}^{\neg c}| - \sum_{n \in \mathcal{S}^c} \sum_{m \in \mathcal{S}^c} \alpha^{c,m} \kappa[\mathbf{h}^{c,n}, \mathbf{h}^{c,m}] \right. \\ &\quad + \sum_{n \in \mathcal{S}^c} \sum_{m \in \mathcal{S}^{\neg c}} \alpha^{\neg c,m} \kappa[\mathbf{h}^{c,n}, \mathbf{h}^{\neg c,m}] - \sum_{n \in \mathcal{S}^{\neg c}} \sum_{m \in \mathcal{S}^c} \alpha^{c,m} \kappa[\mathbf{h}^{\neg c,n}, \mathbf{h}^{c,m}] \\ &\quad \left. + \sum_{n \in \mathcal{S}^{\neg c}} \sum_{m \in \mathcal{S}^{\neg c}} \alpha^{\neg c,m} \kappa[\mathbf{h}^{\neg c,n}, \mathbf{h}^{\neg c,m}] \right) \\ &= \beta.\end{aligned}$$

By making the replacements from points ①, ② and ③ in (5), we deduce that  $s_\kappa^c(\mathbf{h}) = s_{\kappa'}^{c'}(\mathbf{h}')$ .  $\square$

## E Empirical Evaluation

This appendix provides useful details to reproduce the empirical evaluation from Section 3.

**Computing Resources.** All the empirical evaluations were run on a single machine equipped with a 18-Core Intel Core i9-10980XE CPU and a NVIDIA RTX A4000 GPU. The machine runs on Python 3.9 [76] and Pytorch 1.10.2 [77].

**Dataset licenses.** The MNIST dataset dataset is made available under the terms of the Creative Commons Attribution-Share Alike 3.0 License. The ECG dataset dataset is made available under the terms of the Open Data Commons Attribution License v1.0. The CUB dataset is made available for non-commercial research and educational purposes.

**Models.** The detailed architecture of the models are provided in Tables 2, 3 and 4. The InceptionV3 architecture is the same as in the literature [60]. We use its official Pytorch implementation.

**Data Split.** All the datasets are naturally split in training and testing data. In the ECG dataset, the different types of abnormal heartbeats are imbalanced (e.g. the fusion beats constitute only 0.7% of the training set). Hence, we create a synthetic training set with balanced concepts using SMOTE [78]. As in [41], the CUB dataset is augmented by using random crops and random horizontal flips.

**Model Fitting.** In fitting each model, we use the test set as a validation set since our purpose is not to obtain the models with the best generalization but simply models that perform well on a set of examples we wish to explain (here the examples from the test set). All the models are trained to minimize the cross-entropy between their prediction and the true labels. The hyperparameters are as follows. ① For MNIST we use a Adam optimizer with batches of 120 examples, a learning rate of  $10^{-3}$ , a weight decay of  $10^{-5}$  for 50 epochs with patience 10. ② For ECG we use a Adam optimizer with batches of 300 examples, a learning rate of  $10^{-3}$ , a weight decay of  $10^{-5}$  for 50 epochs with patience 10. ③ For CUB, we use a stochastic gradient descent optimizer with batches of 64 examples, a learning rate of  $10^{-3}$ , a weight decay of  $4 \cdot 10^{-5}$  for 1,000 epochs with patience 50.

**Concepts.** The concept mapping between MNIST classes and concepts is provided in Table 5. For the ECG and the CUB datasets, the presence/absence of a concept for each example is readily available in the dataset.

**Concept classifiers.** All the concept classifiers are implemented with scikit-learn [71]. For CAR classifiers, we fit a SVC with Gaussian RBF kernel and default hyperparameters from scikit-learn. For CAV classifiers, we fit a linear classifier with a stochastic gradient descent optimizer with learning rate  $10^{-2}$  and a tolerance of  $10^{-3}$  for 1,000 epochs and the remaining default hyperparameters from scikit-learn.

**Statistical significance.** The statistical significance test from Section 3.1.1 is performed with the scikit-learn [71] implementation of the permutation test. For MNIST and ECG, we consider 100 permutations per concept. For CUB, this test is more expensive since the latent spaces are high-dimensional. We consider only 25 permutations per concept in that case.

**Concept-based feature importance.** In the experiment from Section 2.3, we use Captum’s implementation [72] of Integrated Gradients [25] with default parameters. In the case of CUB, storing the feature importance scores for each concept and for the whole test set requires a prohibitive amount of memory. To avoid this problem, we select  $C = 6$  concepts and subsample 50 positive and 50 negative

Table 2: MNIST Model

Block Name	Layer Type	Hyperparameters	Activation
Conv1	Conv2d	Input Channels = 1, Output Channels = 16, Kernel Size = 5, Stride = 1, Padding = 0	ReLU
	Dropout2d	p = 0.2	
	MaxPool2d	Kernel Size = 2, Stride = 2	
Conv2	Conv2d	Input Channels = 16, Output Channels = 32, Kernel Size = 5, Stride = 1, Padding = 0	ReLU
	Dropout2d	p = 0.2	
	MaxPool2d	Kernel Size = 2, Stride = 2	
Lin1	Flatten Linear	Input Features = 512, Output Features = 10	
Lin2	Dropout	p = 0.2	
	Linear	Input Features = 10, Output Features = 5	
	Dropout Linear	p = 0.2 Input Features = 5, Output Features = 10	

Table 3: ECG Model

Block Name	Layer Type	Hyperparameters	Activation
Conv1	Conv1d MaxPool1d	Input Channels = 1, Output Channels = 16, Kernel Size = 3, Stride = 1, Padding = 1 Kernel Size = 2	
Conv2	Conv1d MaxPool1d	Input Channels = 16, Output Channels = 64, Kernel Size = 3, Stride = 1, Padding = 1 Kernel Size = 2	
Conv3	Conv1d MaxPool1d	Input Channels = 64, Output Channels = 128, Kernel Size = 3, Stride = 1, Padding = 1 Kernel Size = 2	
Lin	Flatten Linear	Input Features = 2944, Output Features = 32	
	Leaky ReLU Linear	Negative Slope = $10^{-2}$ Input Features = 32, Output Features = 2	Leaky ReLU

Table 4: CUB Model

Block Name	Layer Type	Hyperparameters	Activation
InceptionOut	InceptionV3 [60]	Pretrained = True	
	Linear	Input Features = 2048, Output Features = 200	

examples per concept from the test set. This corresponds to a set of 600 examples. We compute the feature importance for these examples only.

**Alternative architecture.** We extended the analysis of Section 3.1.1 to a ResNet-50 architecture. We fine-tune the ResNet model on the CUB dataset and reproduced the experiment from Section 3.1.1 with this new architecture. In particular, we fit a CAR and a CAV classifier on the penultimate layer of the ResNet. We then measure the accuracy averaged over the CUB concepts. This results in  $(89 \pm 1)\%$  accuracy for CAR classifiers and  $(87 \pm 1)\%$  accuracy for CAV classifiers. We deduce that CAR classifiers are highly accurate to identify concepts in the penultimate ResNet layer. As in the main paper, we observe that CAR classifiers outperform CAV classifiers, although the gap is smaller than for the Inception-V3 neural network. We deduce that our CAR formalism extends beyond the architectures explored in the paper and we hope that CAR will become widely used to interpret any more architectures.

## F Use Case

This appendix provides useful details to reproduce the use case from Section 3.2.

**Computing Resources.** The use case was run on a single machine equipped with a 18-Core Intel Core i9-10980XE CPU and a NVIDIA RTX A4000 GPU. The machine runs on Python 3.9 [76] and Pytorch 1.10.2 [77].

Table 5: MNIST Concepts

Class \ Concept	Loop	Vertical Line	Horizontal Line	Curvature
0	✓	✗	✗	✓
1	✗	✓	✗	✗
2	✓	✗	✗	✓
3	✗	✗	✗	✓
4	✗	✓	✓	✗
5	✗	✗	✓	✓
6	✓	✗	✗	✓
7	✗	✓	✓	✗
8	✓	✗	✗	✓
9	✓	✗	✗	✓

Table 6: SEER Model

Layer Type	Hyperparameters	Activation
Linear	Input Features = 21, Output Features = 400	ReLU
Dropout	p = 0.3	
Linear	Input Features = 400, Output Features = 100	ReLU
Dropout	p = 0.3	
Linear	Input Features = 100, Output Features = 2	

**Dataset license.** The SEER dataset is made available under the terms of the SEER Research Data Use Agreement.

**Model.** The detailed architecture of the model is provided in Table 6.

**Data split.** We randomly split the whole SEER dataset into a training set (90% of the data) and a test set (the remaining 10%). Since patients with a death outcome are in minority (less than 3%), we oversample them to obtain a balanced training set.

**Model fitting.** In fitting the model, we use the test set as a validation set since our purpose is not to obtain the model with the best generalization but simply a model that performs well on a set of examples we wish to explain (here the examples from the test set). The model is trained to minimize the cross-entropy between its prediction and the true labels. We use a Adam optimizer with batches of 500 examples, a learning rate of  $10^{-3}$ , a weight decay of  $10^{-5}$  for 500 epochs with patience 50.

**Concepts.** The concepts correspond to prostate cancer grades. Those grades can be computed from the Gleason score as follows [70]:

$$\text{Grade}(\text{Gleason}_1, \text{Gleason}_2) = \begin{cases} 1 & \text{if } \text{Gleason}_1 + \text{Gleason}_2 \leq 6 \\ 2 & \text{if } \text{Gleason}_1 = 3 \wedge \text{Gleason}_2 = 4 \\ 3 & \text{if } \text{Gleason}_1 = 4 \wedge \text{Gleason}_2 = 3 \\ 4 & \text{if } \text{Gleason}_1 + \text{Gleason}_2 = 8 \\ 5 & \text{if } \text{Gleason}_1 + \text{Gleason}_2 \geq 9. \end{cases}$$

It goes without saying that the model is trained with the Gleason scores only, not with the grades.

**Concept classifiers.** We fit a SVC with linear kernel and default hyperparameters from scikit-learn.

**Concept-based feature importance.** We use Captum’s implementation [72] of Integrated Gradients [25] with default parameters.

## G Explanation Robustness

**Adversarial perturbations.** We perform an experiment to evaluate the robustness of CAR explanations with respect to adversarial perturbations. In this experiment, we work with the MNIST dataset in the same setting as the experiment from Section 3.1.2. We train a CAR concept classifier for each MNIST concept  $c \in [C]$ . We use the CAR classifier to output TCAR scores relating the concept  $c$  with each class  $k \in [d_Y]$ . As in the main paper, since the ground-truth association between concepts and classes is known (e.g. the class corresponding to digit 8 will always have the concept loop), we can compute the correlation  $r(\text{TCAR}, \text{TrueProp})$  between our TCAR score and the ground-truth proportion of examples that exhibit the concept. In this experiment, the correlation is evaluated on a test set  $\mathcal{D}_{\text{test}} = \mathcal{D}_{\text{adv}} \sqcup \mathcal{D}_{\text{orig}}$  that contains adversarial test examples  $\mathcal{D}_{\text{adv}}$  and original test examples  $\mathcal{D}_{\text{orig}}$ . Each adversarial MNIST image  $\mathbf{x}_{\text{adv}} \in \mathcal{D}_{\text{adv}}$  is constructed by finding a small (w.r.t. the  $\|\cdot\|_\infty$  norm) perturbation  $\epsilon \in \mathbb{R}^{d_X}$  around an original test image  $\mathbf{x} \in \mathcal{X}$  that maximizes the prediction shift for the model  $f : \mathcal{X} \rightarrow \mathcal{Y}$ :

$$\epsilon = \arg \max_{\tilde{\epsilon} \in \mathbb{R}^{d_X}} \text{CrossEntropy}[f(\mathbf{x}), f(\mathbf{x} + \tilde{\epsilon})] \text{ s.t. } \|\tilde{\epsilon}\|_\infty < .1$$

Table 7: MNIST Adversarial Perturbation Sensitivity

Adversarial %	$r(\text{TCAR}, \text{TrueProp})$
0	.99
5	.99
10	.99
20	.99
50	.97
70	.96
100	.92

The adversarial image is then defined as  $\mathbf{x}_{\text{adv}} \equiv \mathbf{x} + \epsilon$ . We measure the correlation  $r(\text{TCAR}, \text{TrueProp})$  by varying the proportion  $\frac{|\mathcal{D}_{\text{adv}}|}{|\mathcal{D}_{\text{test}}|}$  of adversarial examples in the test set. The results are reported in Table 7.

We observe that the TCAR scores keep a high correlation with the true proportion of examples that exhibit the concept even when all the test examples are adversarially perturbed. We conclude that TCAR explanations are robust with respect to adversarial perturbations in this setting.

**Background shift.** For completeness, we have also adapted the background shift robustness experiment in Section 7 from [27]. We use CAR to explain the predictions of our Inception-V3 model trained on the original CUB training set. The explanations are made on test images where the background has been replaced. As [27], we use the segmentation of the CUB dataset to isolate the bird on each image. The rest of the image is replaced by a random background sampled from the *Place365* dataset [80]. This results in a test set  $\mathcal{D}_{\text{test}}$  with a background shift with respect to the training set. By following the approach from Section 3.1.2 of our paper, we measure the correlation  $r(\text{TCAR}, \text{TrueProp})$  between the TCAR score and the true proportion of examples in the class that exhibit the concept for each (class, concept) pair. We measured a correlation of  $r(\text{TCAR}, \text{TrueProp}) = .82$  in the background-shifted test set. This is close to the correlation for the original test set reported in the main paper, which suggests that CAR explanations are robust with respect to background shifts. Note that this correlation is still better than the one obtained with TCAV on the original test set.

## H Using CAR to Understand Unsupervised Concepts

Our CAR formalism adapts to a wide variety of neural network architectures. In this appendix, we use CAR to analyze the concepts discovered by a self explaining neural network (SENN) trained on the MNIST dataset. As in [81], we use a SENN of the form

$$f(\mathbf{x}) = \sum_{s=1}^S \theta_s(\mathbf{x}) \cdot g_s(\mathbf{x}),$$

Where  $g_s(\mathbf{x})$  and  $\theta_s(\mathbf{x})$  are respectively the activation and the relevance of the synthetic concept  $s \in [S]$  discovered by the SENN model. We follow the same training process as [81]. This yields a set of  $S = 5$  concepts explaining the predictions made by the SENN  $f : \mathcal{X} \rightarrow \mathcal{Y}$ .

We use our CAR formalism to study how the synthetic concepts  $s \in [S]$  discovered by the SENN are related to the concepts  $c \in \{\text{Loop, Vertical Line, Horizontal Line, Curvature}\}$  introduced in our paper. With our formalism, the relevance of a concept  $c$  for a given prediction  $\mathbf{x} \mapsto f(\mathbf{x})$  is measured by the concept density  $\rho^c \circ g(\mathbf{x})$ . To analyze the relationship between the SENN concept  $s$  and the concept  $c$ , we can therefore compute the correlation of their relevance:

$$r(s, c) = \text{corr}_{\mathbf{X} \sim P_{\text{empirical}}(\mathcal{D}_{\text{test}})}[\theta_s(\mathbf{X}), \rho^c \circ g(\mathbf{X})].$$

When this correlation increases, the concepts  $s$  and  $c$  tend to be relevant together more often. We report the correlation between each pair  $(s, c)$  in Table 8.

We note the following:

Table 8: SENN Concepts Correspondence

Correlation r(s, c)	Loop	Vertical Line	Horizontal Line	Curvature
<b>SENN Concept 1</b>	-0.28	-0.12	0.26	0.11
<b>SENN Concept 2</b>	-0.50	0.71	-0.03	-0.69
<b>SENN Concept 3</b>	-0.47	0.10	0.71	-0.14
<b>SENN Concept 4</b>	-0.33	0.02	-0.06	-0.01
<b>SENN Concept 5</b>	0.57	-0.00	-0.63	0.07

1. SENN Concept 2 correlates well with the Vertical Line Concept.
2. SENN Concept 3 correlates well with the Horizontal Line Concept
3. SENN Concept 5 correlates well with the Loop Concept.
4. SENN Concepts 1 and 4 are not well covered by our concepts.

The above analysis shows the potential of our CAR explanations to better understand the abstract concepts discovered by SENN models. We believe that the community would greatly benefit from the ability to perform similar analyses for other interpretable architectures, such as disentangled VAEs.

## I Using CAR with NLP

CAR is a general framework and can be used in a wide variety of domains that involve neural networks. In the main paper, we show that CAR provides explanations for various modalities:

1. Large image dataset
2. Medical time series
3. Medical tabular data

We now perform a toy experiment to assess if those conclusions extend to the NLP setting. We train a small CNN on the IMDB Review dataset to predict whether a review is positive or negative. We use GloVe [82] to turn the word tokens into embeddings. We would like to assess whether the concept  $c = \text{Positive Adjective}$  is encoded in the model’s representations. Examples that exhibit the concept  $c$  are sentences containing positive adjectives. We collect a positive set  $\mathcal{P}^c$  of  $N^c = 90$  such sentences. The negative set  $\mathcal{N}^c$  is made of  $N^c$  sentences randomly sampled from the Gutenberg Poem Dataset. We verified that the sentences from  $\mathcal{N}^c$  did not contain positive adjectives. We then fit a CAR classifier on the representations obtained in the penultimate layer of the CNN.

We assess the generalization performance of the CAR classifier on a holdout concept set made of  $N^c = 30$  concept positive and negative sentences (60 sentences in total). The CAR classifier has an accuracy of 87% on this holdout dataset. This suggests that the concept  $c$  is smoothly encoded in the model’s representation space, which is consistent with the importance of positive adjectives to identify positive reviews. We deduce that our CAR formalism can be used in a NLP setting. We believe that using CARs to analyze large-scale language model would be an interesting study that we leave for future work.

## J Increasing Explainability at Training Time

Improving neural networks explainability at training time constitutes a very interesting area of research but is beyond the scope of our paper. That said, we believe that our paper indeed contains insights that might be the seed of future developments in neural network training. As an illustration, we consider an important insight from our paper: the fact that the accuracy of concept classifiers seems to increase with the depth of the layer for which we fit a classifier. In the main paper, this is mainly reflected in Figure 4. This observation has a crucial consequence: it is not possible to reliably characterize the shallow layers in terms of the concepts we use.

In order to improve the explainability of those shallow layers, one could leverage the recent developments in contrastive learning. The purpose of this approach would be to separate the concept set

$\mathcal{P}^c$  and  $\mathcal{N}^c$  in the representation space  $\mathcal{H}$  corresponding to a shallow layer of the neural network. A practical way to implement this would be to follow [79]. Assume that we want to separate concept positives and negatives in the representation space  $\mathcal{H}$  induced by the shallow feature extractor  $\mathbf{g} : \mathcal{X} \rightarrow \mathcal{H}$ . As [79], one can use a projection head  $\mathbf{p} : \mathcal{H} \rightarrow \mathcal{Z}$  and enforce the separation of the concept sets through the contrastive loss

$$\mathcal{L}_{\text{cont}}^c = \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in (\mathcal{P}^c)^2} -\log \frac{\exp(\tau^{-1} \cdot \cos[\mathbf{p} \circ \mathbf{g}(\mathbf{x}_i), \mathbf{p} \circ \mathbf{g}(\mathbf{x}_j)])}{\sum_{\mathbf{x}_k \in (\mathcal{P}^c \cup \mathcal{N}^c) \setminus \{\mathbf{x}_i\}} \exp(\tau^{-1} \cdot \cos[\mathbf{p} \circ \mathbf{g}(\mathbf{x}_i), \mathbf{p} \circ \mathbf{g}(\mathbf{x}_k)])), \quad (6)$$

where  $\cos(\mathbf{z}_1, \mathbf{z}_2) \equiv \frac{\mathbf{z}_1^\top \mathbf{z}_2}{\|\mathbf{z}_1\|_2 \cdot \|\mathbf{z}_2\|_2}$  and  $\tau \in \mathbb{R}^+$  is a temperature parameter. The effect of this loss is to group the concept positive examples from  $\mathcal{P}^c$  together and apart from the concept negatives  $\mathcal{N}^c$  in the representation space  $\mathcal{H}$ . To the best of our knowledge, concept-based contrastive learning has not been explored in the literature. We believe that it would constitute an interesting contribution to the field based on the insights from our paper.

## K Further Examples

In this appendix, we provide several examples to illustrate the experiments from Section 3.

**Concept classifiers.** The accuracy of the concept classifiers for all the MNIST and ECG concepts are given in Figures 11 and 12. Each box-plot is built with 10 random seeds where the concept sets  $\mathcal{P}^c$  and  $\mathcal{N}^c$  are allowed to vary. We observe that the CAR classifiers are more accurate for each of the observed concepts

**Global explanations.** The global concept explanations for all the MNIST concepts and some of the CUB concepts are given in Figures 13 and 14. As the examples from Section 3.1.2, we see that TCAR explanations are more consistent with the human concept annotations.

**Saliency maps.** Examples of concept-based vanilla saliency maps for MNIST, ECG and CUB examples are given in Figures 15, 16, 17, 18 and 19. As explained in the main paper, we observe that concept-based saliency maps are indeed distinct from vanilla saliency maps. Furthermore, saliency maps for different concepts are not interchangeable. In the CUB case, we note that concept are not always identified with the minimal amount of features (e.g. some of the saliency maps from Figure 19 highlight pixels that do not always belong to the bird's breast). This surprising observation seems to occur even for concept-bottleneck models that are explicitly trained to recognize concepts [44]. We believe that this could be improved by training concept classifiers with images that include a segmentation highlighting the concept of interest (e.g. the bird's breast). We leave this idea for future works.

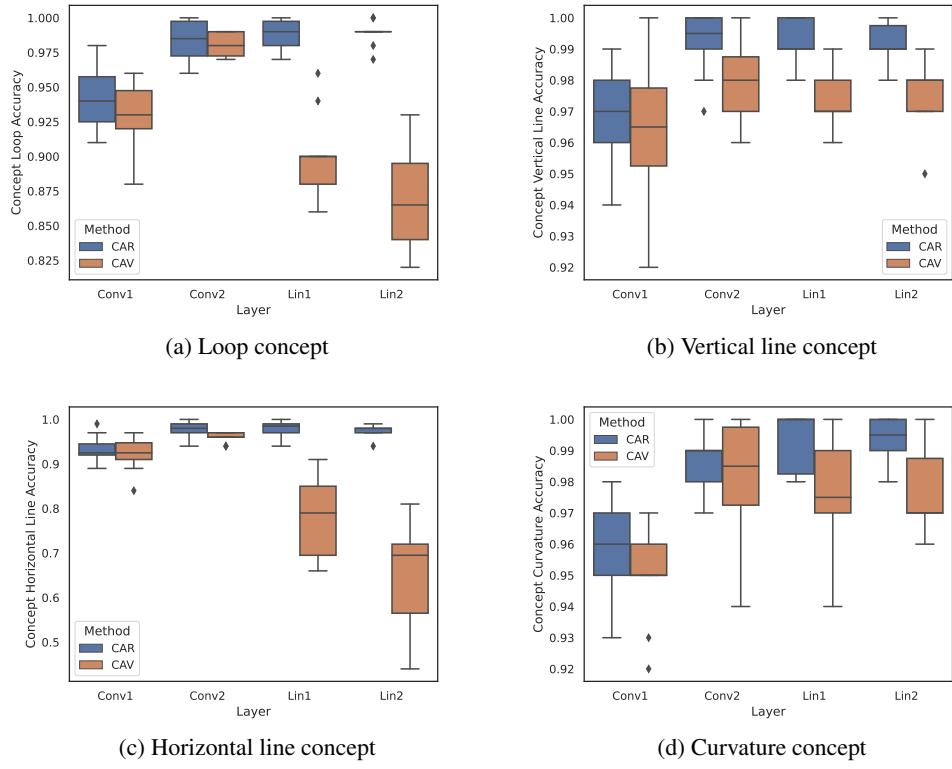


Figure 11: Concept accuracy for MNIST concepts

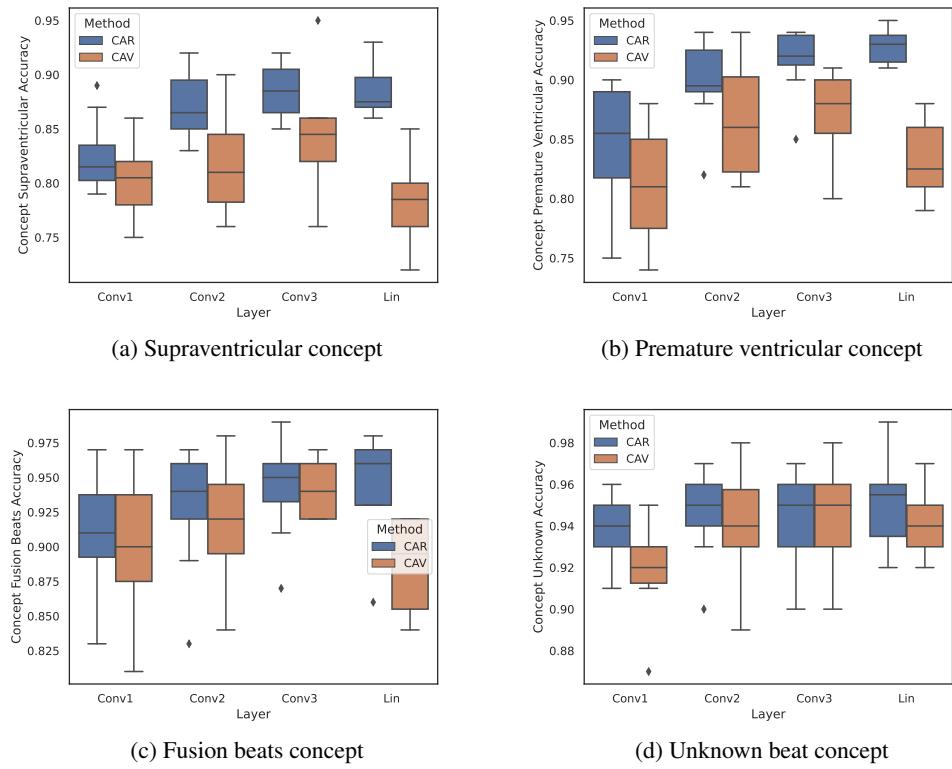


Figure 12: Concept accuracy for ECG concepts

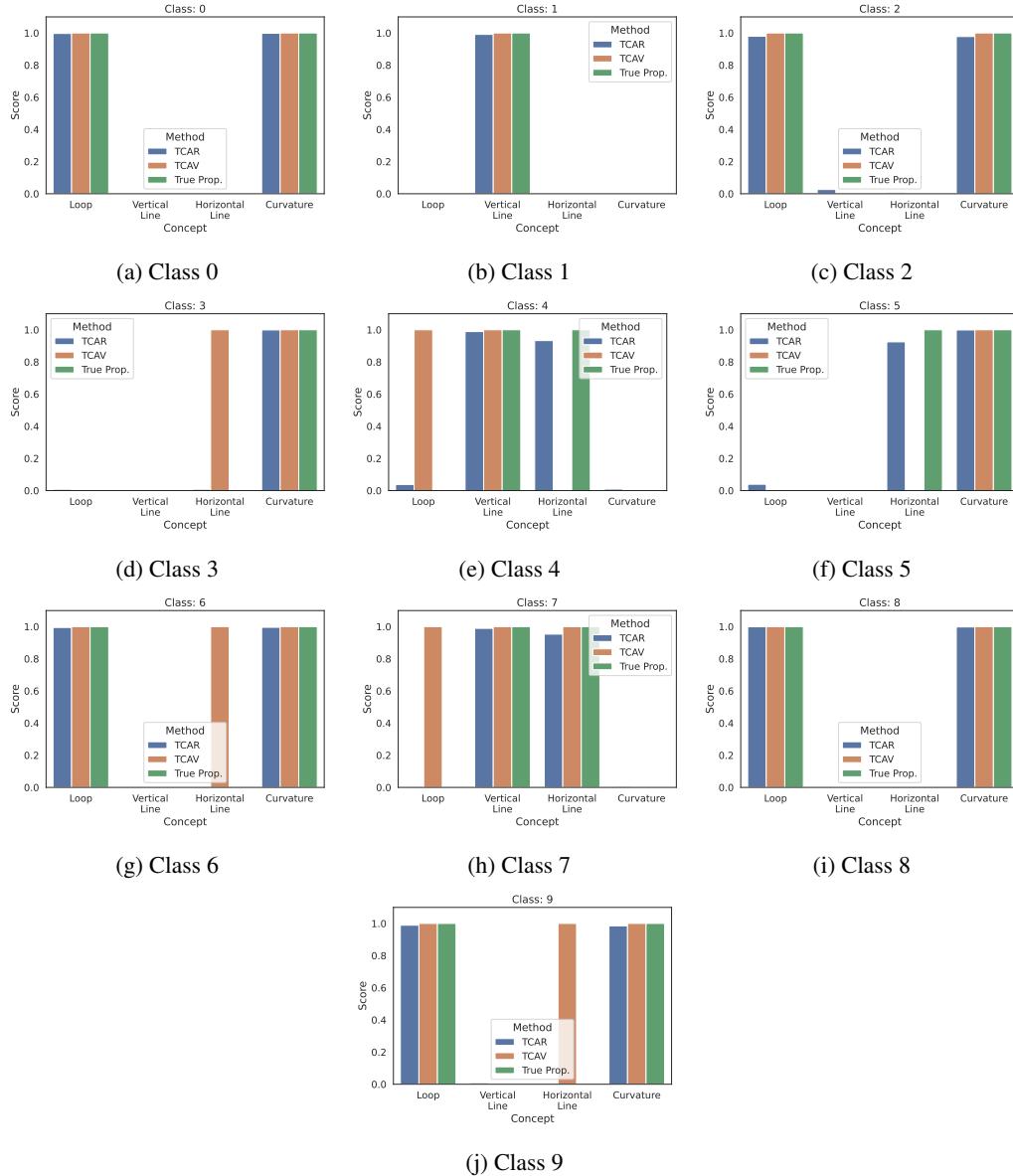


Figure 13: Global concept explanations for MNIST

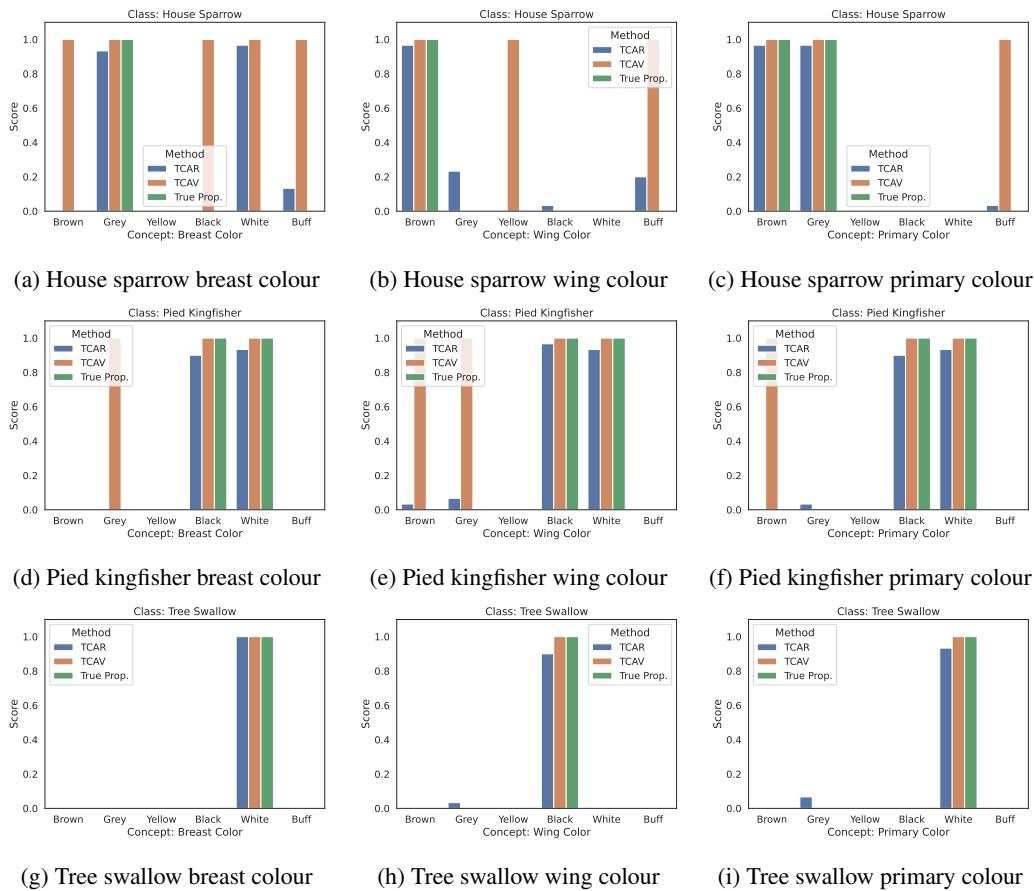
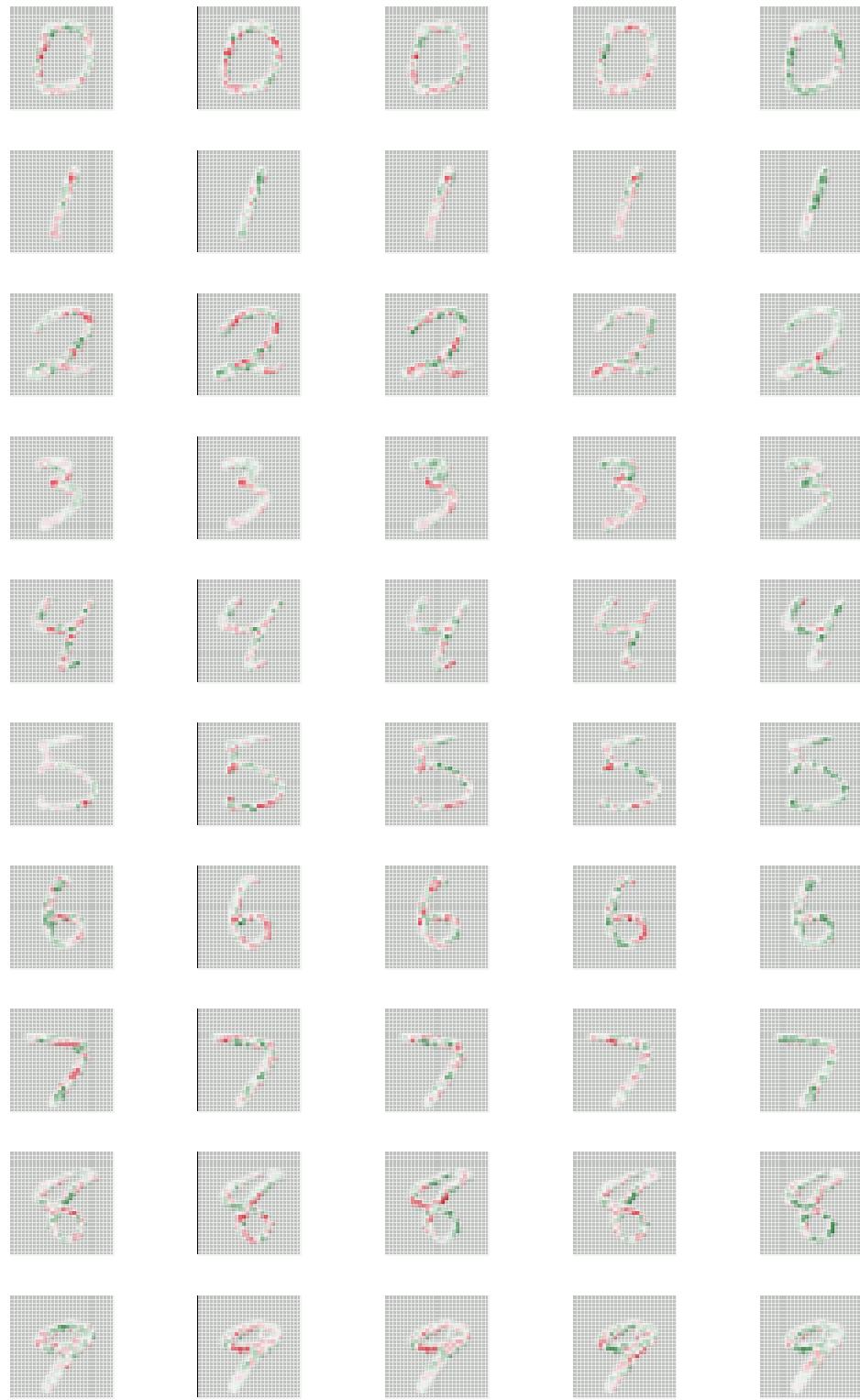


Figure 14: Global concept explanations for CUB



(a) Loop

(b) Vertical line

(c) Horiz. line

(d) Curvature

(e) Vanilla

Figure 15: Concept-based saliency maps for MNIST

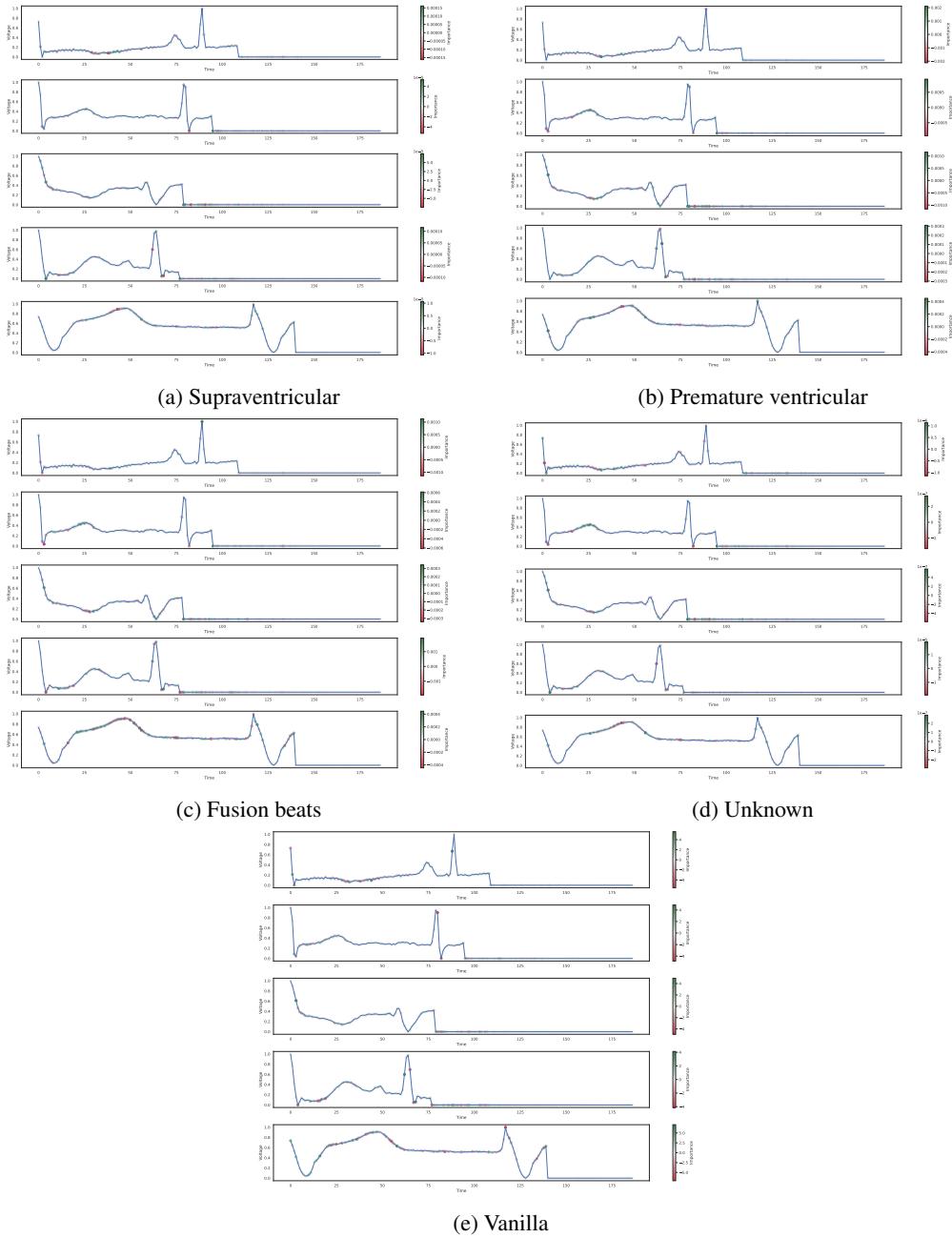
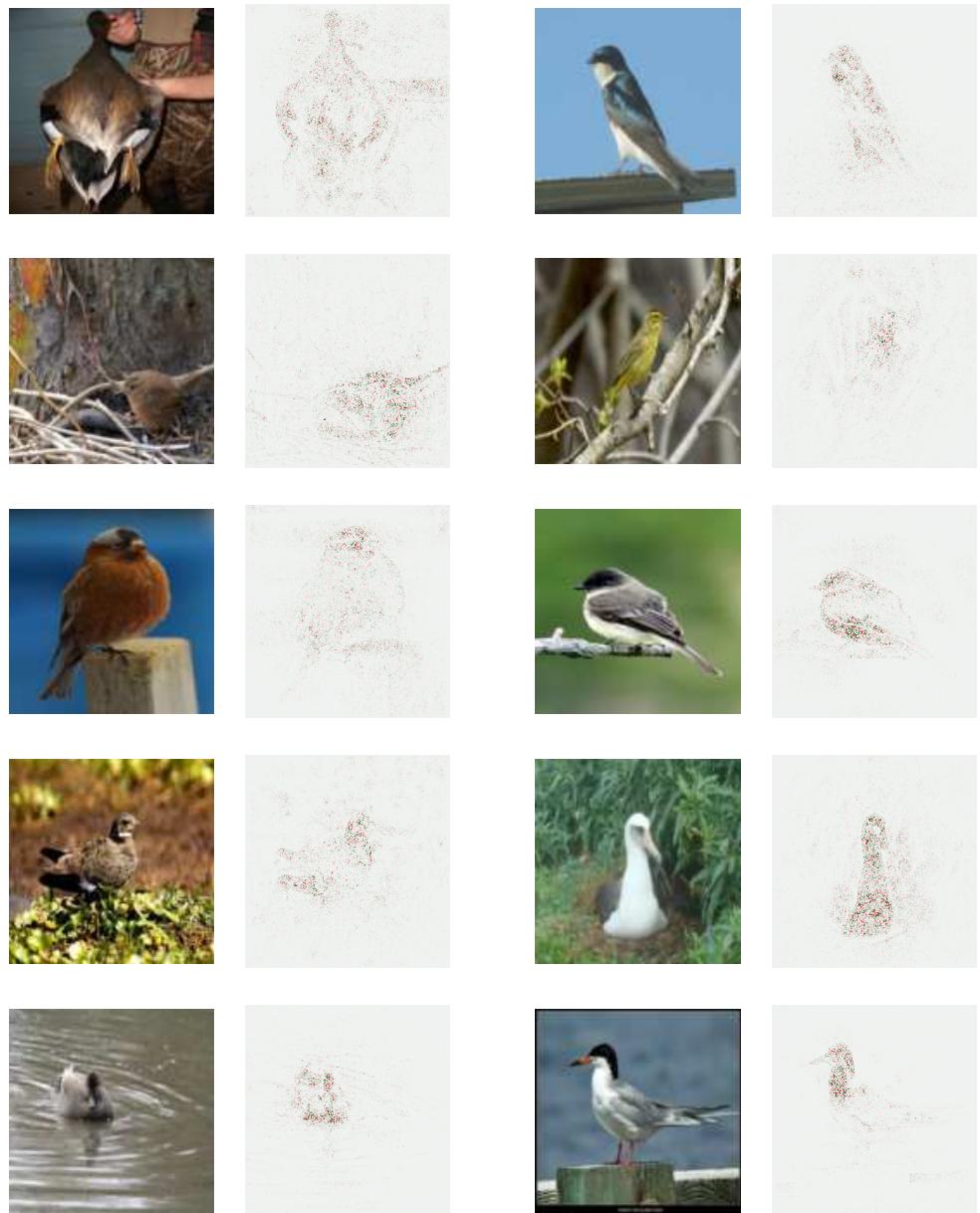


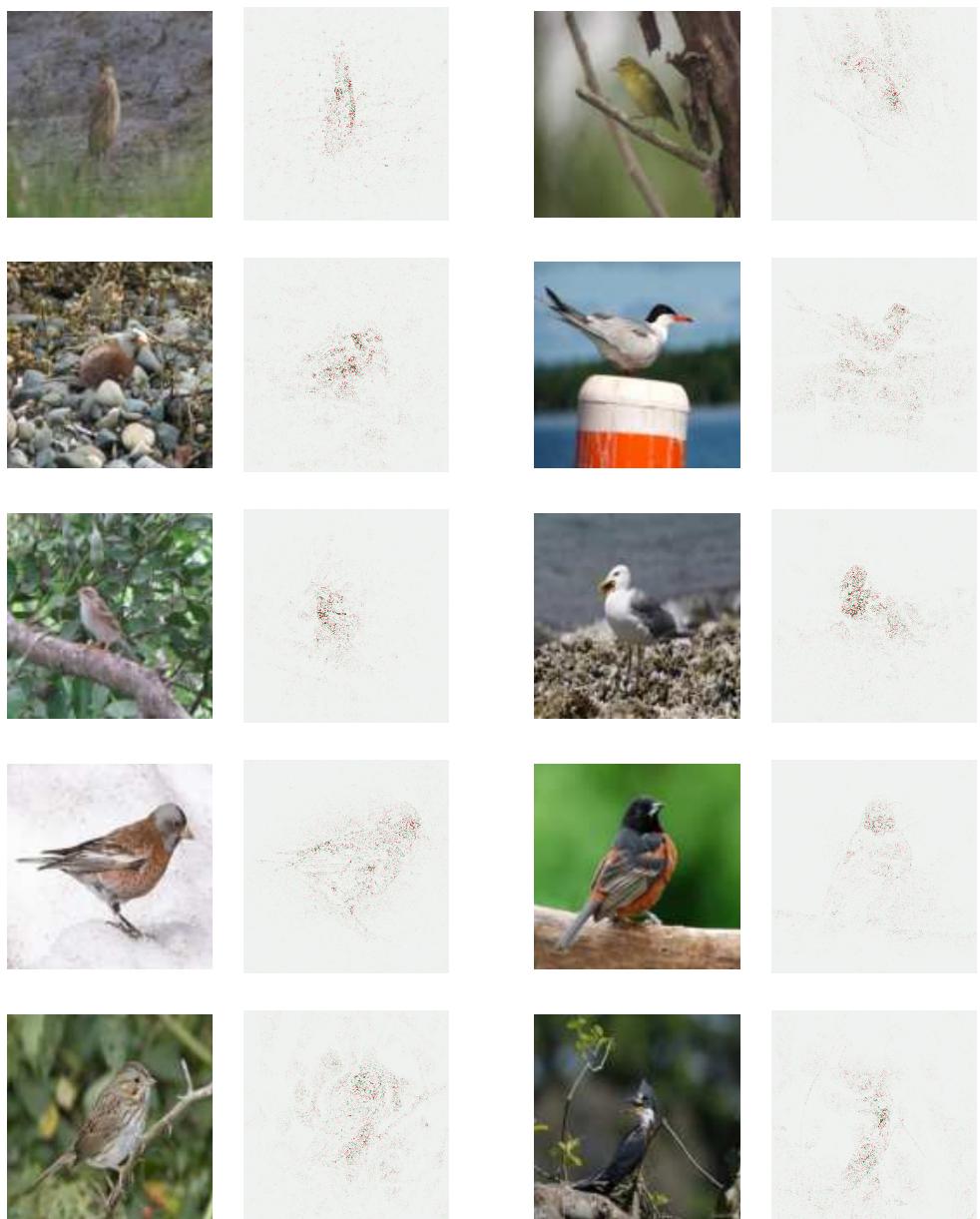
Figure 16: Concept-based saliency maps for ECG



(a) Brown belly

(b) Solid belly pattern

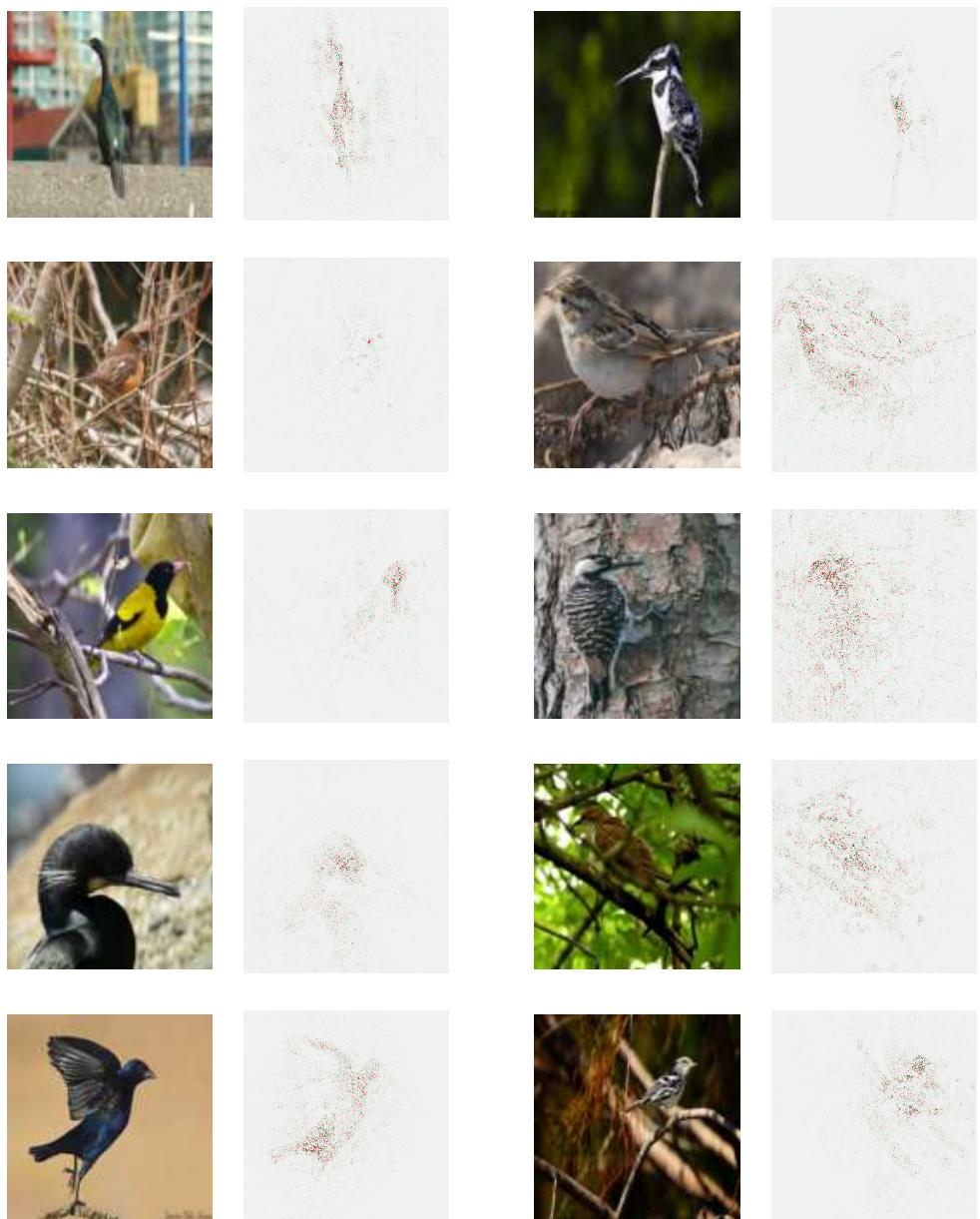
Figure 17: Concept-based saliency maps for CUB (1/3)



(a) Stripped back pattern

(b) Solid back pattern

Figure 18: Concept-based saliency maps for CUB (2/3)



(a) Black breast

(b) White breast

Figure 19: Concept-based saliency maps for CUB (3/3)