

## Reviewer f9CQ

### 1. Analysis with SOTA Architecture

**The paper considers mostly convolutional networks. It lacks analysis of more state-of-the-art architectures like residual networks or transformers.**

Redo CUB with a ResNet. Put results here when available

Copy the NLP experiment from reviewer qiiw.

### 2. Increasing Explainability at Training Time

**Would it be possible to get inspiration from the insights in the paper to improve neural network training, with the purpose of increasing explainability?**

Propose a regularization/pretraining scheme? If time permits, do a small experiment.

### 3. Sensitivity to Adversarial Attacks

**How are the proposed explanations sensitive to adversarial attacks?**

As suggested by the reviewer, 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 from our paper. 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, this 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  $x_{\text{adv}} \in \mathcal{D}_{\text{adv}}$  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  $x \in \mathcal{X}$  that maximizes the prediction shift for the black-box  $f : \mathcal{X} \rightarrow \mathcal{Y}$ :

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

The adversarial image is then defined as  $x_{\text{adv}} \equiv x + \epsilon$ . 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 below.

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

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 to adversarial perturbations in this setting.

For completeness, we have also adapted the background shift robustness experiment in Section 7 from *Koh, P. et al. (2020). Concept Bottleneck Models*. As in our paper, 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 Koh et al., 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. 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.

## Reviewer znZB

### 1. Hyperparameter Choice

**The introduction of kernels may lead to much more hyper-parameter choice to use in practice for very complicated concepts and networks, an additional RBF kernel may be insufficient to make the data linearly separable in practice.**

Since our CAR classifiers are kernel-based, they indeed come with extra hyperparameters (e.g. the kernel width). We would like to emphasize that, to ensure fair comparisons with the CAV classifiers, none of these hyperparameter has been optimized in the experiments from Section 3.1.1 and 3.1.2. We have used the default hyperparameters in the scikit-learn implementation of support vector

classifiers. In all our experiments, CAR classifiers substantially outperform CAV hyperparameters without having to tune the 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 hyperparameters  $\theta_h$  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. That said, we believe that the above approach to tune the hyperparameters of CAR classifiers might be useful in other cases so we will add it to the manuscript.

We agree that not all concepts can be captured by CAR classifiers, even after hyperparameter optimization (e.g. Figure 4.c from the paper, we see that CAR classifiers don't generalize well on the layer Mixed-5d). As we argue in the paper, this is a strong indication that our concept smoothness assumption (Assumption 2.1) is violated. This implies that concepts are not smoothly encoded in the geometry of the model's representation space  $\mathcal{H}$ . The user can therefore deduce that the concept is unlikely to be salient to interpret  $\mathcal{H}$ . In that sense, the inability to fit CAR classifiers that generalize well is as informative as the ability to fit CAR classifiers that generalize well. Note that this whole reasoning is made more quantitative through statistical hypothesis testing in Section 3.1.1 of our paper.

## 2. Layer Selection

The layer selection problem is not dealt with, as TCAV suggested that sometimes earlier layers may be more fruitful even if the accuracy seems lower (for more low-level concepts). In my personal experience, using the mixed-7b layer inception-V3 for TCAV usually produces a

much better result than the penultimate layer. Moreover, how does one choose the layer to apply TCAR is not addressed. **note: TCAV does not work well with the penultimate layer since  $d c/d \text{ activation} = W\_c$  (which is independent to the instance), and thus for all instances in the same class the directional derivative to the same concept would be fixed.**

With our CAR formalism, the user is free to choose the layer they want to interpret. In creating TCAR explanations for various layers, we decided to select the layer for which the concept classifiers (for both CAR and CAV) generalize better to unseen examples, as measured in our experiment from Section 3.1.1 from our paper. Following the reviewer’s recommendation, we decided to repeat the comparison between TCAV and TCAR from Section 3.1.2 with the layer Mixed-7b of our Inception-V3 classifier. In doing so, we measured the following correlation between the scores and the groundtruth proportion of examples within a class that exhibit the concept:

$$r(\text{TCAV}, \text{TrueProp}) = .46 \quad r(\text{TCAR}, \text{TrueProp}) = .71$$

For both TCAV and TCAR, these correlations are lower than the ones obtained in the model’s penultimate layer. In this case, it appears that the association between classes and concepts are more meaningfully encoded in the deeper layers of the neural network. In a future work, it would be interesting to extend this study to other architectures.

### 3. Discrepancy between CAV and TCAV Accuracy

**In the inception experiment, CAV has a pretty close accuracy to CAR (showing that the linear separability is not a huge issue), but the evaluation of TCAV score in table 1 seems completely wrong. The failure of TCAV in MNIST and ECG is understandable since the model seems under-represented, and an additional kernel classifier would help a lot. However, the result in Inception-v3 is not convincing.**

Double check our results. CAV accuracy describes how accurate a linear classifier is to separate concept positives and negatives in the model’s representation space. TCAV scores, to be accurate, need an extra dimension: the concept need to appropriately be associated to predicted labels. It is perfectly possible to have accurate separation of the concept positive/negative in a model’s representation space without having this concept-to-class association. Point at Inception examples in the supplementary material (complete by giving the accuracy of the selected concepts).

### 4. Feature Importance Evaluation

The evaluation of concept-based feature importance is only a sanity check, how is this useful?

The experiment is in the same spirit as the one to validate TCAR scores: we would like to assess whether TCAR leads to explanations that are consistent with what humans would expect. This does not only validate TCAR as an explanation method but the model itself. The difficulty in the experiment is the following: since the ground-truth concept-specific concept importance is unknown, we cannot assess the feature importance directly (per image concept-level segmentations are unavailable in the datasets). Even if the ground-truth feature importance is unknown, it is legitimate to expect that concepts that are identified with the same features should lead to similar feature importance scores. Quantitatively, this translates into a high correlation between the respective concept feature importance scores. Write this a bit more mathematically.

## Reviewer qiiw

### 1. Generalizing CAV Sensitivity Interpretations

One of the critical weaknesses of CAR is that since explanations are generated through a Kernel-based technique, it loses the nice interpretations CAV offers, like if one increases the presence of a concept, how does it affect the model predictions. Would it be possible to replicate the interpretations CAV offers like concept sensitivity? An alternate would be to perform test time interventions on the Kernel-based Concept classifier.

Explain that the concept density can be seen as our equivalent of the CAV concept sensitivity. If this density is high, this implies that the example lies within a cluster of concept positives. Proposition to generalize sensitivity:

$$S_k^c(x) = (\nabla_h \rho^c[g(x)])^\top (\nabla_h l_k[g(x)])$$

### 2. Using CAR with Unsupervised Concepts

Is it possible to run CAR with concepts (or latent variables) generated by Autoencoders or VAE that are inherently noisy or abstract? Could CAR be used to analyze the concepts generated by techniques like SENN?

Read SENN and see if we can realistically run an experiment with it. Otherwise, focus on the latent variables discovered by any unsupervised model really.

### 3. Robustness of CAR Explanations

How robust are explanations generated by CAR? Is it robust to change in backgrounds of the images?

Copy and adapt the robustness point from reviewer f9CQ.

## 4. CAR for NLP

### Would CAR work for other domains like NLP?

CAR is a general framework and can be used in a wide variety of domains that involve neural networks. In our paper, we show that CAR provides explanations for models trained on large image datasets, medical time series and medical tabular data. As suggested by the reviewer, we perform a small 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 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 CAR to analyze large-scale language model would be an interesting study that we leave for future work.

## 5. Using CAR without Human Annotations

**Is it possible to relax the assumption that such a class of techniques requires additional annotation of concepts? As of now, CAR requires a user to specify the positive and negative examples for each concept.**

This could probably be merged with Point 2 above. Note that human concepts need to be defined by humans.

## 6. Minor Points

**Given that CAR assumes to have access to the feature extractor of the model, it isn't truly a black-box setup, unlike paper portrays.**

**I would encourage the authors to list the limitations of the proposed approach.**

We thank the reviewer for these additional remarks. We will make sure that to implement those changes in the final manuscript.

## Experiments TODO

- DONE *Small* Kernel width optimization
- RUNNING *Medium* CUB with ResNet
- DONE *Small* CUB TCAR/TCAV and Acc with mixed-7b layer
- DONE *Large* NLP Experiment
- DONE *Small* Robustness to adversarial attacks
- DONE *Medium* CUB with background shift
- ? CAR and unsupervised concepts?
- ? Training regularization?