

# Generalizability of Concept Knowledge in Machine Learning Using TCAV Scores: A Case Study Using Different Skin-Lesion Datasets

Moritz C. Schwinghammer<sup>\*†</sup>, Laines Schmalwasser<sup>†</sup>, Sireesha Chamarthi<sup>†</sup>, Yuri A.W. Shardt<sup>\*</sup>

<sup>\*</sup>Technical University Ilmenau: Department of Automation Engineering, Ilmenau, 98684, Germany

<sup>†</sup>German Aerospace Center: Institute for Data Science, Jena, 07745, Germany

---

**Abstract:** The adoption of artificial intelligence in safety-critical fields such as skin-lesion classification remains rare, partly due to the lack of interpretability of the models' decisions. Concept-based interpretability methods like testing with concept activation vectors (TCAV) address this issue by quantifying how specific human-understandable concepts influence a model's decisions. For better comprehensibility and analyzability of the results, the paper also introduces concept detection scores (CDS). The CDS are aggregated TCAV scores which are directionally unified. Machine-learning models have issues with generalizability, meaning their predictive performance degrades when confronted with unseen data with a different distribution, *i.e.*, from a different domain. This paper examines the ability of concept knowledge, as measured by TCAV, to handle generalizability issues, specifically the robustness to domain shifts. The results show that interpretability provided by TCAV is robust against domain shifts (between technical domains).

**Keywords:** Interpretability, TCAV, Generalizability, Skin-Lesion Classification, Domain Shifts

---

## 1. INTRODUCTION

The adoption of artificial intelligence, particularly in safety-critical areas, has been slow. This can be partially attributed to two factors: the presence of uncertainties in AI systems and the prevalent distrust of the opaque decision-making processes in neural-network (NN) models [1]. These concerns are addressed by explainable artificial intelligence (xAI), which seeks to provide interpretability by opening the machine learning (ML) black box and explaining the models' decisions.

Various approaches to xAI and interpretability exist, varying in temporal type (*post-hoc* or *ante-hoc*), in the nature of the explanations provided (ranging from prototypes to decision-governing rules), and in the explained scope of explanations (from explaining individual samples to providing global interpretability for a NN) [2]. One area of interpretability is concept-based NN explanations. These explanations focus on explaining the factors that lead to a NN model's decisions in human-understandable terms. A prominent, global, *post-hoc* interpretability method that identifies representation vectors of human-understandable concepts in the activation (latent) space of a NN, and quantifies their influence on the model's predictions, is testing with concept activation vectors (TCAV) [3].

Skin lesion classification (SKL) is a safety-critical, high-risk field where limited-to-no adoption of AI for computer-aided diagnostics (CAID) has occurred, in part due to the lack of explanations for the classifications provided by the CAID NN models [4]. Dermatologists, as the medical professionals in this field, state that CAID systems could provide valuable insights as a "qualified second opinion," which dermatologists frequently provide to each other [5].

However, such a second opinion needs to be able to explain the reasoning behind its decisions, rather than simply stating a final classification [5].

Concept-based interpretability tools, such as TCAV, may prove beneficial in this context. This is because dermatologists themselves rely on specific criteria or concepts when learning to classify skin lesions, the collective knowledge is codified in several checklists, which in turn provide guidance for skin lesion classification depending on the presence, absence, or type of observable criteria on a lesion [6]. Thus, if done correctly, concept-based interpretability of the decisions of a CAID application could "speak the same language" as the experts it assists.

However, ML models in general have issues with generalizability. When a ML model is faced with unseen data from a distribution different from the one on which it was originally trained, then its predictive performance often degrades significantly [7, 8]. In this case, "different distribution" refers not merely to an unseen test-data split, but to a fundamentally different distribution of a dataset's inherent features [7]. A relevant example for a technical domain shift between different skin-lesion (SK) datasets would be image acquisition system settings like contrast and brightness, which differ between the datasets, since the comprising images were captured in different hospitals [8]. Addressing the generalizability issues is a major concern in ML-based SKL [8].

Therefore, before CAID tools based on concept-based explanations can be considered for real-world application, it is necessary to assess how these explanations are affected by domain shifts, *i.e.*, when the model must extrapolate to new, unseen distributions (domains). In this paper, TCAV and the newly introduced concept detection scores will be used to

provide an understanding of the degree to which different ML models learn and apply concepts on three different datasets. By comparing the concept knowledge detected on samples from a known distribution with concept knowledge on an unknown distribution, conclusions about generalizability can be drawn.

## 2. THEORY AND BACKGROUND

### 2.1. Dermatological concepts

Although studies show that early detection of skin lesions is challenging, dermatologists use heuristic approaches, such as the seven-point checklist, to diagnose and encode their expertise [6]. The diagnoses used in this paper are either *nevus* or *melanoma*. The concepts used in this paper are derived from the criteria of the seven-point checklist [6] and described in accordance with [9, 10] as:

- 1) Pigment Networks (PN): A pigment network is a grid-like pattern of interconnected lines surrounding lighter areas. A typical pigment network (PN\_T), which indicates a benign lesion, is symmetrical and consistent. An atypical pigment network (PN\_AT) is asymmetrical, with variable color, thickness, and spacing, and is suggesting *melanoma* presence.
- 2) Blue-Whitish Veil (BWV): This concept refers to an irregularly shaped, slightly blue lesion (-spot) covered by a whitish haze resembling ground glass.
- 3) Streaks (ST): Streaks can be either regular (ST\_R), indicating a benign lesion, or irregular (ST\_IR), suggesting *melanoma*. In general, they appear as straight extensions, bulbous projections, or a widened network along the lesion edge.
- 4) Dots and Globules (DG): Regular dots and globules (DG\_R) are centered within the lesion middle or aligned on the network lines and are uniform, indicating a benign lesion. Irregular dots and globules (DG\_IR) exhibit higher variability, suggesting *melanoma*.
- 5) Regression Structures (RS): The presence of fine grey-bluish dots, light areas without blood vessels, or shiny-white structures indicates regression structures, suggesting *melanoma*.

### 2.2. Concept detection and testing with concept activation vectors

Testing with concept activation vectors (TCAV), introduced by Been Kim *et al.* [10], is a method for interpreting neural networks by analyzing concept relevance. TCAV involves splitting a NN  $m$  at a specified layer  $l$  with  $e$  neurons. The model  $m$  described by  $g_m(x): \mathbb{R}^a \rightarrow \mathbb{R}^k$ , maps inputs  $x \in \mathbb{R}^a$  to logit space  $\mathbb{R}^k$  where  $x$  represents the pixel-based input images. The activations at layer  $l$ , denoted by  $o_{l,m}(x): \mathbb{R}^a \rightarrow \mathbb{R}^e$ , are used to train a binary classifier for each labeled individual concept  $c$  with  $c \in C$  and  $C$  being the set of all observed concepts. The normal to the hyperplane of the linear classifier is the concept activation vector (CAV),  $v_{c,m}^l$ , which points toward the concept encoded in the  $e$ -dimensional activation space  $\mathbb{R}^e$ .

To compute TCAV scores, input data must have classification class labels  $k$ , corresponding to the class predicted by the NN's final (logit) layer. The "classifying"

part of the NN complements the "feature extracting" NN part  $o_{l,m}(x)$  and is described with  $h_{l,m}: \mathbb{R}^n \rightarrow \mathbb{R}^k$ . The sensitivity  $S_{c,k,l,m}^{Sens}(x)$  measures how sensitive a model  $m$  is to a concept  $c$  for classification class  $k$  at layer  $l$ . It is described as the directional derivative:

$$S_{c,k,l,m}^{Sens}(x) = \lim_{\varepsilon \rightarrow 0} \frac{h_{l,k,m}(o_{l,m}(x) + \varepsilon v_{c,m}^l) - h_{l,k,m}(o_{l,m}(x))}{\varepsilon}, \quad (2.1)$$

$$S_{c,k,l,m}^{Sens}(x) = \nabla h_{l,k,m}(o_{l,m}(x)) \cdot v_{c,m}^l,$$

with  $S_{c,k,l,m}^{Sens}(x)$  providing quantifiable information about the sensitivity of the classification of a single input image  $x$ . For global *post-hoc* interpretability, TCAV scores aggregate  $S^{Sens}$  globally, thus quantifying the relevance of a concept  $c$  across all samples belonging to a classification class  $k$ . The TCAV scores are defined as:

$$S_{c,k,m,w}^{TCAV} = \frac{|\{x \in X_k: S_{c,k,l,m}^{Sens}(x) > 0\}|}{|X_k|}, \quad (2.2)$$

which is the fraction of all predictions of  $X_k$  where the concept was classification-relevant, over all images of  $X_k$ , the set containing all samples of a specific class  $k$  [3]. Repeated calculations, with  $w > 1$ , of  $S^{TCAV}$  for each combination of variables account for differences in data pre-processing and the initialization of the classifiers [3].

### 2.3. Concept detection scores

While  $S^{TCAV}$  are averaged  $S^{Sens}$ , there are still  $|C||K|w$  individual TCAV scores per model  $m$  with  $c \in C$  and  $k \in K$ , which presents a challenge for analysis. Here, set  $C$  is the set containing all concepts  $c$  and  $K$  is the set containing all classes  $k$ . The multitude of TCAV scores, as well as the issue of directionality, motivated the introduction of concept detection scores (CDS) as is done in this paper.

The issue of directionality refers to the fact that depending on concept  $c$  and classification class  $k$ , "perfect" concept detection results in either a high  $S^{TCAV}$  of 1 or a low  $S^{TCAV}$  with a value of 0. This is due to the fact that the concepts  $c$  are all indicative of precisely one of the classification classes  $k$ . Thus, if  $S^{TCAV}$  is calculated for a concept  $c$  which is indicative of another classification class  $k$  than the one examined during calculation of the  $S^{TCAV}$ , a perfect concept detection would result in a score of 0. This inversion, which depends on the combination of concept  $c$  and classification class  $k$ , poses a significant hindrance for later analysis of concept knowledge.

To address this, CDS serve two purposes, first, the reduction of the dimensionality of the  $S^{TCAV}$  results through aggregation, second, a partial inversion of the calculated  $S^{TCAV}$  to norm the direction of the detected concept knowledge, with 1 being always indicative of a "perfect" concept detection.

The reduction in dimensionality is feasible because all observed concepts  $c$  are indicative of one element of the binary class  $k$ . This allows aggregation of all concepts for each element of class  $k$ , which is controlled by the variable concept type group  $B$  defined as:

$$B = \begin{cases} b_1 & c \in C_{negative} \\ b_2 & c \in C_{positive} \\ b_3 & c \in C \end{cases}, \quad (2.3)$$

with  $C_{negative}$  and  $C_{positive}$  being the sets of all concepts  $c$  indicative of either classification class  $k$ . While  $B$  accounts for the indicative  $k$  of the individual  $c$ , CDS are also dependent on the basic  $k$ , since the underlying  $S^{TCAV}$  depend on it. For the purpose of CDS,  $k$  is grouped into the target class group  $D$ , defined as:

$$D = \begin{cases} d_1 & S_{c,k_{negative},m,w}^{TCAV} \\ d_2 & S_{c,k_{positive},m,w}^{TCAV} \\ d_3 & d_1 \cup d_2 \end{cases} \quad (2.4)$$

The individual CDS are described:

$$S_{b_1,d_1}^{CD} = \frac{\sum_{c_{neg}}^{|C_{negative}|} \sum_{i=0}^w S_{w,k_{negative},c_{neg}}^{TCAV}}{|w| \cdot |C_{negative}|} \quad (2.5)$$

$$S_{b_1,d_2}^{CD} = \frac{\sum_{c_{neg}}^{|C_{negative}|} \sum_{i=0}^w 1 - S_{w,k_{negative},c_{neg}}^{TCAV}}{|w| \cdot |C_{negative}|} \quad (2.6)$$

$$S_{b_2,d_1}^{CD} = \frac{\sum_{c_{pos}}^{|C_{positive}|} \sum_{i=0}^w 1 - S_{w,k_{negative},c_{pos}}^{TCAV}}{|w| \cdot |C_{positive}|} \quad (2.7)$$

$$S_{b_2,d_2}^{CD} = \frac{\sum_{c_{pos}}^{|C_{positive}|} \sum_{i=0}^w S_{w,k_{positive},c_{pos}}^{TCAV}}{|w| \cdot |C_{positive}|} \quad (2.8)$$

where  $c_{neg}$  and  $c_{pos}$  are individual concepts indicative of either  $k_{neg}$  or  $k_{pos}$ .

### 3. METHODOLOGY

#### 3.1. Datasets

Three types of datasets were used in the paper. The first type, the *concept dataset*, consists of a single dataset containing the concept information used to train the CAVs. This dataset includes SK images with labels for both elements of  $k$ , as well as the presence or absence of concepts. The second dataset type, referred to as the *random dataset*, also consists of a single dataset from which images were randomly selected and paired with random concept labels. The final dataset type comprises three *domain-shifted datasets*. Domain shifts mean that they contain different distributions (domains) which have observable shifts in their features and/or characteristics [8].

The *concept dataset* used was the 7-point Criteria Evaluation Database defined as  $f_{d7pt}^{cnc}$  [10]. The concepts, the classification class  $k$  they indicate, and the number of images containing each concept are shown in Table 1.

Table 1: Overview of all applied concepts

Concept of seven-point checklist	Pheno-type	Abbreviation	Indicates $k$	# of imgs
<b>Pigment network</b>	typical	PN_T	<del>mel</del>	335
	atypical	PN_AT	<del>mel</del>	216
<b>Blue whitish veil</b>		BWV	<del>mel</del>	183
<b>Streaks</b>	regular	ST_R	<del>mel</del>	96
	irregular	ST_IR	<del>mel</del>	237
<b>Dots and globules</b>	regular	DG_R	<del>mel</del>	301
	irregular	DG_IR	<del>mel</del>	392
<b>Reg. structures</b>		RS	<del>mel</del>	183

“Indicates  $k$ ” with a value of ~~mel~~ refers to the fact that the presence of said concept indicates the absence of *melanoma*

The *random dataset* was the train split of the ISIC2018 dataset [11]. Finally, the *domain-shifted datasets* were datasets BCN20000, HAM10000 and MSK. These datasets were separated by [8] into their underlying domains, which could be both technical, as well as biological. This paper focuses on the technical domains.

#### 3.2. Procedure

In accordance with [9], concept training was repeated twenty times on each NN model for each concept to guard against statistical influence on classifier capability, thus setting  $w = 20$ . Consequently, this required twenty individual stratified splits of  $f_{d7pt}^{cnc}$  for each combination of  $S_{c,k,m,w}^{TCAV}$  variables. These splits were further split into cluster-based undersampled train splits, and as-is evaluation splits for each individual concept  $c$ . The per-concept-balanced training splits were then used to train twenty individual classifiers per concept for each NN model. The training splits were used as input and passed through the NN models until reaching the split layer  $l$ , where activations were extracted. Across all NN models, regardless of the model architecture, layer  $l$  was consistently set as the first regularization or flattening layer immediately following the final imported convolutional layer of the primary model architecture (VGG16, VGG19, InceptionNetV3, or ResNet50).

Training of  $w$  separate classifiers was one safeguard against erroneous results to which CAV/TCAV is susceptible. An additional safeguard was the calculation of a random baseline. Using randomly selected images from the *random dataset* as input verified whether the calculated results were significantly different from those produced by classifiers trained on data without useful information. To achieve this, the *random dataset* ISIC2018 was used as the source for the randomly selected images, and thusly defined as  $f_{ISIC2018}^{rnd}$ . The images of  $f_{ISIC2018}^{rnd}$  were then assigned a partially random label, either  $k_{nevus}$  or  $k_{melanoma}$ . In accordance with the literature,  $k_{melanoma}$  is the positive class or  $k_{positive}$ , since in SKL it is imperative to detect all *melanoma* cases. The class distribution of the real per-concept-balanced train splits which were not equal because of the cluster-based undersampling was mirrored with the randomly assigned labels. The total number of classifiers trained on random images from  $f_{ISIC2018}^{rnd}$  were 160 as determined by  $w|C| = 160$ , since eight individual concepts  $c$  comprise set  $C$ . All 160 random baseline classifiers, along with the corresponding TCAV scores, formed the population for a single random concept.

The classifier training process begins with extraction of activations from activation space  $\mathbb{R}^e$  of the NN model for each per-concept-balanced training split. The activations, together with their corresponding labels, were then shuffled and passed to a linear classifier. The type of classifier used for the results presented in the paper was a linear classifier with  $L_2$  regularization. For evaluation of the classifiers, the per-concept-balanced test splits were fed as input into the primary NN model, where the activations were extracted at layer  $l$ , before being fed into the trained classifiers. The performance of each classifier was then tracked using the metrics accuracy, precision, recall, and F1 score.

The process of model training, classifier training and then concept knowledge detection on the concept target datasets  $f^{tcn}$  is shown in Figure 1. Concept target datasets  $f^{tcn}$  are the *domain-shifted datasets* upon which the concept knowledge of the models is evaluated. In turn the source dataset of a model, which was also one of the three *domain-shifted datasets* was defined as  $f^{src}$ . For concept detection using TCAV scores, the CAV first had to be determined, *i.e.*, the coefficient vector of each classifier was extracted as the concept activation vector  $v_{c,m}^l$ . Multiplication of the dot product of  $v_{c,m}^l$  and the gradient of the classification part of the model  $\nabla h_{l,k}(o_l(x))$  resulted in the sensitivity  $S_{c,k,l}^{Sensitivity}(x)$  per image  $x$ . Aggregation of all sensitivities using TCAV scores resulted in  $w$ -times  $S^{TCAV}$  for each individual concept  $c$  and for both classification classes  $k_{nevus}$  and  $k_{melanoma}$ .

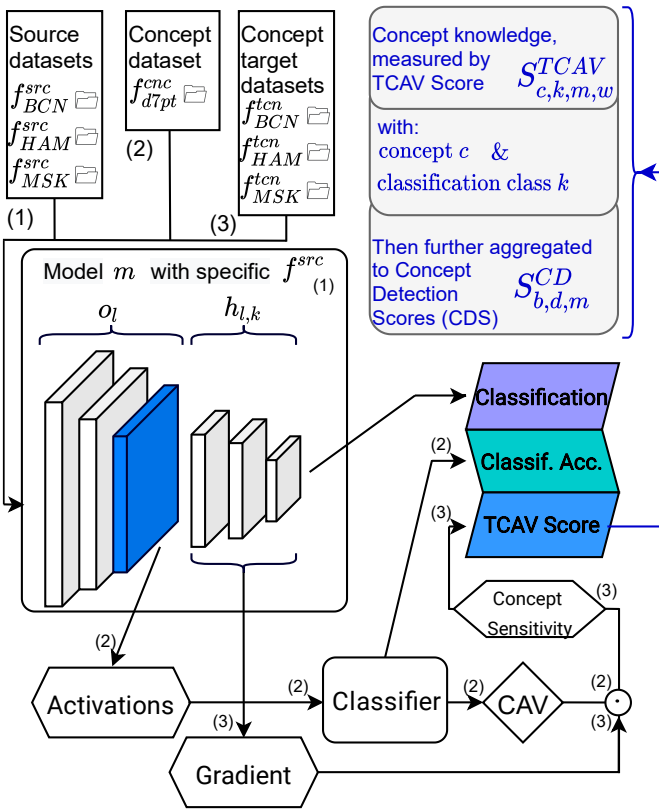


Figure 1: Overview of the process for a single model  $m$

#### 4. RESULTS AND ANALYSIS

Before the exploration of concept detection on domain-shifted datasets, it is necessary to examine the validity of the detected or missed concepts. The validation baseline is calculated on randomly selected images from dataset  $f_{ISIC2018}^{rnd}$  and is expected at the center of the result space with no bias towards concept presence or absence [9]. Thus, for the results of this paper which all are between 0 and 1, the baseline is expected at 0.5. In this paper, the results of the calculation of “random classifiers” on a single model  $m$  corroborate this observation, with the mean F1 score being 0.510 and the mean accuracy being 0.505. Thus, the classifier baseline can be set to 0.5, in accordance with expectations

and observed results. A larger difference would indicate a hidden bias in the model.

Since none of the concept target datasets  $f_{BCN}^{tcn}$ ,  $f_{HAM}^{tcn}$ , and  $f_{MSK}^{tcn}$  provide concept labels, the test splits of the *concept dataset*  $f_{d7pt}^{cnc}$  are used for validation. A general overview of the fidelity of concept detection through the classifiers on  $f_{d7pt}^{cnc}$  is presented in Figure 2. The figure confirms that all concepts, regardless of which class  $k$  they indicate, are detected by the classifiers with an average total F1 score of 0.705. Moreover, almost all classifiers correctly recognize their respective concept  $c$ . Out of the 960 trained classifiers per  $c$  across all NN models, fifteen classifiers for concept BWV, six classifiers for concept PN\_AT, and four classifiers for concept ST\_R have F1 scores below the baseline of 0.5. Thus, since most classifiers perform better than the baseline, the general detectability of concept knowledge is validated.

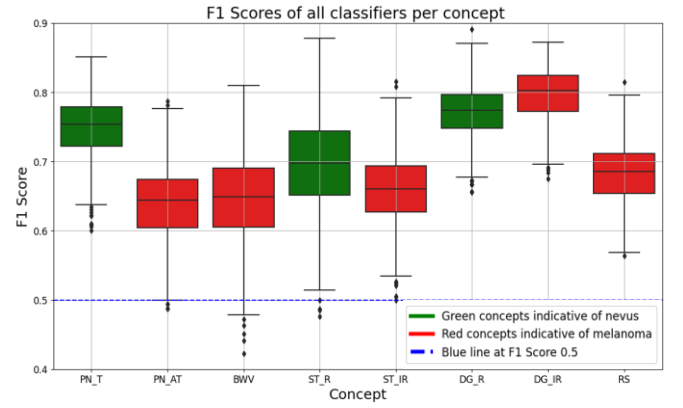


Figure 2: F1 scores of all twenty classifiers per concept  $c$ , calculated on the unseen test split of  $f_{d7pt}^{cnc}$

With the validation of the validity of the concept detection on  $f_{d7pt}^{cnc}$  concluded, the concept detection on the unlabelled *domain-shifted datasets* can be examined. Here, too, a random baseline for TCAV scores is expected at around 0.5 [9]. This is confirmed by the calculation of 160 TCAV scores for each of the classes  $k$ , which resulted in a mean of 0.51 for  $k_{nevus}$  and 0.48 for  $k_{melanoma}$ . Thus, for TCAV as well, the baseline can be set to 0.5. Compared to this baseline, 78.1% of TCAV results are significantly different for a  $p$ -value  $< 0.05$ .

Individual TCAV scores for a single  $m$ , selected for its concept detection capability, with  $f_{BCN}^{src}$  and  $f_{BCN}^{tcn}$  as source and concept target dataset are shown in Figure 3. Since  $f_{BCN}^{src} = f_{BCN}^{tcn}$ , the results are within domain. In the left subplot of the figure, it can be observed that  $S_{c_{negative}, k_{nevus}, m, w}^{TCAV}$  are mostly close to 1, while  $S_{c_{positive}, k_{nevus}, m, w}^{TCAV}$  are near zero. This matches the expectations, since it is expected that concepts  $c_{negative}$  which are indicative of  $k_{nevus}$  are detected and influential on  $x$  belonging to  $X_{k_{nevus}}$ . The signs of most of the TCAV scores per concept, switch for TCAV scores calculated on images belonging to  $k_{melanoma}$ , as is also expected.

While TCAV scores visualized with mean and violin plots, as in Figure 3, provide a viable gauge of the detected concepts for a single  $m$  for the case where  $f_{BCN}^{src} = f_{BCN}^{tcn}$ , the

comprehensibility decreases greatly when 16 models are trained on all three  $f^{src}$  and evaluated on all three  $f^{tcn}$ . This is where CDS provide an advantage.

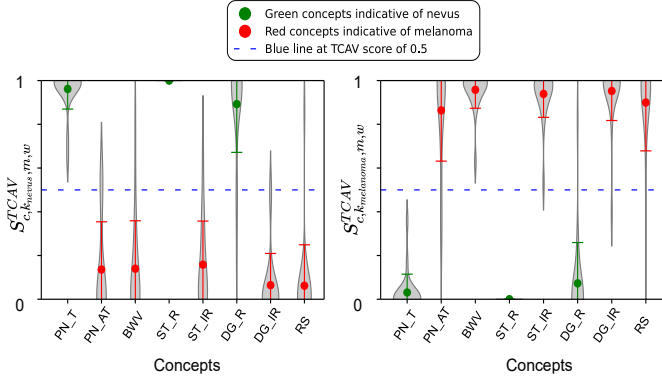


Figure 3: Violin plots of the distribution of TCAV scores  $S^{TCAV}$  per concept  $c$  of a single  $m$  with  $f^{src}_{BCN} = f^{tcn}_{BCN}$

This advantage is shown in Figure 4, which shows the concept detection of all models for all  $S^{CD}$  and within- and across domain. From top to bottom, the first subplot shows  $S^{CD}_{b_1,d_1}$  and thus concepts  $c_{negative}$  detected on images of  $k_{nevus}$ . The following subplots show in order  $S^{CD}_{b_2,d_1}$ ,  $S^{CD}_{b_1,d_2}$ , and  $S^{CD}_{b_2,d_2}$ . Within each subplot, the boxplot color signifies  $f^{tcn}$ , while the background color distinguishes  $f^{src}$ . The color codes are orange for BCN, purple for HAM, olive for MSK, and neon green where d7pt is used as  $f^{tcn}_{d7pt}$ .

The primary observation in Figure 4 is the apparent irrelevance of  $f^{tcn}$ . Models trained on the same  $f^{src}$  seem to report highly similar  $S^{CD}$  for all  $f^{tcn}$ . This pattern is even more pronounced when only the median is considered, and it consistently applies across all  $b$  and  $d$ . The highest difference between the medians of the average concept knowledge measured across all  $f^{tcn}$  is 0.054 for  $f^{src}_{BCN}$  at  $S^{CD}_{b_2,d_1}$ . The average difference for all  $f^{src}$  and  $S^{CD}$  is 0.022.

The distribution from Figure 4 and the miniscule differences between CDS of different  $f^{tcn}$  both suggest that the specific concept target dataset  $f^{tcn}$  is largely irrelevant for the detected concept knowledge and its influence. However, Figure 4 only considers the distribution and median of the CDS. Thus, to observe on the level of individual NN models, Table 2 shows the maximum difference ( $\Delta^{max}$ ) between all instances of  $f^{tcn}$  for each individual CDS,  $m$ , and  $f^{src}$ . The table shows that across all  $f^{src}$ , the mean and median  $\Delta^{max}$  between the instances of  $f^{tcn}$  is below 0.056 for the mean and 0.021 for the median, thus further underscoring the observation that  $f^{tcn}$  has a negligible impact on the concept knowledge. However, some outliers are also present in Table 2. The total maximum of all  $\Delta^{max}$  across all  $f^{tcn}$  shows the highest outliers (value  $> 0.1$ ) at  $S^{CD}_{b_1,d_2,m,f^{src}_{HAM}}$ ,  $S^{CD}_{b_2,d_1,m,f^{src}_{HAM}}$ , and  $S^{CD}_{b_2,d_2,m,f^{src}_{HAM}}$ . Albeit six of nine possible  $\Delta^{max} S^{CD}_{b,d,m}$ , which itself already are the highest differences between CDS of different  $f^{tcn}$  across all types of CDS,  $m$ , and  $f^{src}$ , are below 0.1. Thus, concept knowledge appears to be largely independent of the dataset on which it is measured, in contrast to predictive capability, which generalizes poorly and is highly domain-dependent.

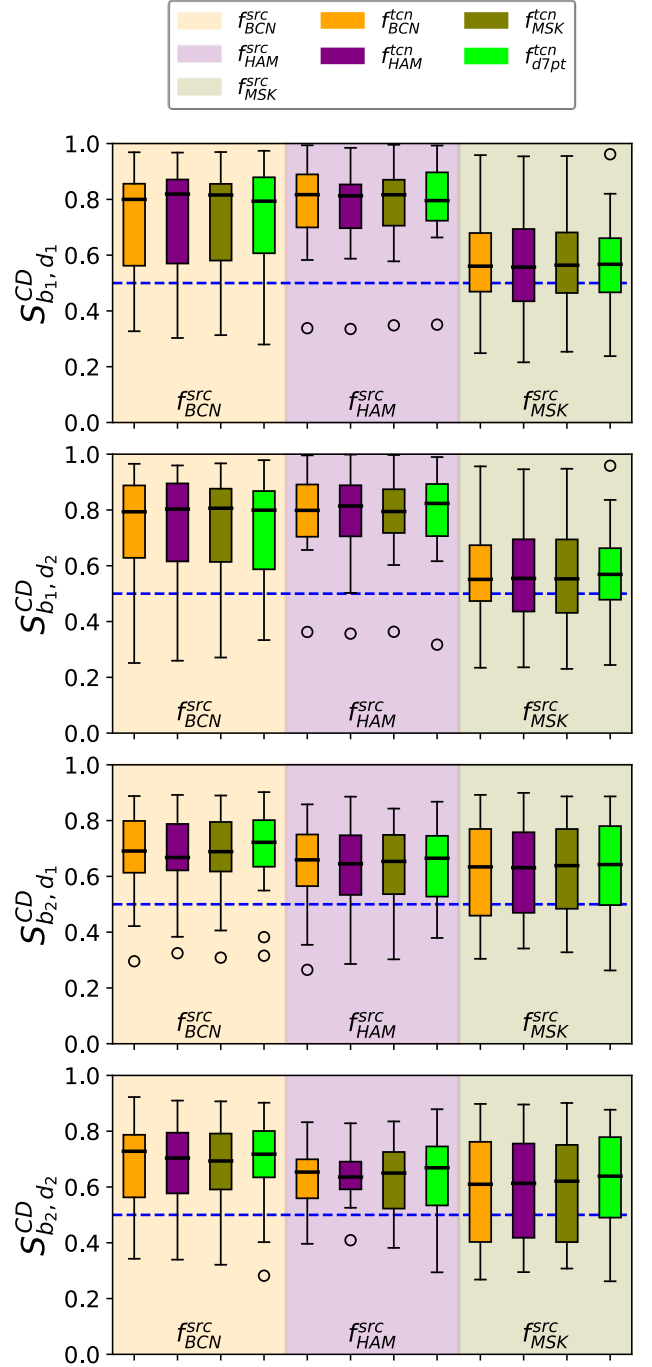


Figure 4: Boxplots of the distribution of concept detection scores  $S^{CD}_{b,d}$  for all combinations of  $b$  and  $d$  separated by source datasets  $f^{src}$  and concept target datasets  $f^{tcn}$

The three highest outliers are all trained on  $f^{src}_{HAM}$ . However, while  $f^{src}_{BCN}$  is overall the best source dataset with respect to (w.r.t.) concept knowledge, as shown in Figure 4,  $f^{src}_{HAM}$  is the dataset with the smallest spread between its CDS. Lastly, NN models trained on  $f^{src}_{MSK}$  command the least concept knowledge, they consistently have the lowest median CDS across all  $b$  and  $d$ , particularly for  $b_1$ , where the concepts indicate  $k_{nevus}$ . This suggests that the concepts in general are harder to learn on  $f^{src}_{MSK}$ , especially for  $c_{nevus}$ , although learning is still possible.

Table 2: Maximum difference between lowest and highest  $S_{b,d,m,f^{tcn}}^{CD}$  between  $f^{tcn}$  and per  $m$  for all  $b$ ,  $d$ , and  $f^{src}$ .

$f^{src}$	Statistic	$\Delta^{max} \text{ of } S_{b_1,d_1,m}^{CD}$	$\Delta^{max} \text{ of } S_{b_1,d_2,m}^{CD}$	$\Delta^{max} \text{ of } S_{b_2,d_1,m}^{CD}$	$\Delta^{max} \text{ of } S_{b_2,d_2,m}^{CD}$
BCN	Max.	0.063	0.047	0.052	0.096
BCN	Median	0.018	0.020	0.019	<b>0.026</b>
BCN	Mean	0.022	0.021	0.021	0.031
HAM	Max.	0.063	<b>0.238</b>	<b>0.141</b>	<b>0.226</b>
HAM	Median	0.013	<b>0.029</b>	0.023	0.021
HAM	Mean	0.019	<b>0.050</b>	<b>0.033</b>	<b>0.056</b>
MSK	Max.	<b>0.066</b>	0.106	0.076	0.068
MSK	Median	<b>0.024</b>	0.026	<b>0.030</b>	0.021
MSK	Mean	<b>0.028</b>	0.035	0.031	0.028

Highest values for each  $S_{b,d,m,f^{tcn}}^{CD}$  highlighted in **bold**, lowest values in *italics*.

This difference between concepts  $c_{nevus}$  and  $c_{melanoma}$  is furthermore observable in the other  $f^{tcn}$ . Models trained on  $f_{BCN}^{tcn}$  and  $f_{HAM}^{tcn}$  seem to learn and use  $c_{nevus}$  concepts significantly more than  $c_{melanoma}$  concepts. Lastly there exists at least one  $m$  with  $S^{CD}$  close to 1 for  $S_{b_1,d_1}^{CD}$  and  $S_{b_1,d_2}^{CD}$ . Thus, a selection of  $m$  by high concept knowledge and usage would have been possible.

## 5. CONCLUSIONS

This paper has examined the by TCAV quantified concept knowledge learned by different NN models. To enhance comprehensibility and to unify the direction of the detected concept knowledge, concept detection scores CDS were introduced and used. Before evaluating the concept knowledge detectability on different datasets, the validity of the results was established. The primary focus of the paper was the examination of detected concept knowledge on three different domain-shifted datasets and the assessment of the influence of said domain shifts. The results showed that concept knowledge is largely domain-shift agnostic, meaning the models are able to apply their learned concept knowledge without suffering from generalizability issues that typically affect their predictive capabilities. Since concept knowledge and usage seems to transcend domains, further research should be done towards the exploration of the relationship between concept knowledge and predictive capabilities within and across domains. If a positive relationship could be established, then NN model selection for high concept knowledge would be a way to address generalizability issues w.r.t. predictive capabilities.

## REFERENCES

- [1] A. J. London, "Artificial Intelligence and Black-Box Medical Decisions: Accuracy versus Explainability," *The Hastings Center report*, no. 1, p. 15–21, 2019.
- [2] Y. Zhang, P. Tiño, A. Leonardis und K. Tang, „A Survey on Neural Network Interpretability,” *IEEE Trans. Emerg. Top. Comput. Intell.*, Bd. 5, Nr. 5, p. 726–742, 2021.
- [3] B. Kim, M. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viegas und R. Sayres, „Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV),“ in *Conference: 35th International Conference on Machine Learning (ICML)*, Stockholm, Sweden, 2018.
- [4] N. Natasha, U. Muhammad, K. S. Muhammad, I. Shahid und A. Douhadji, „A Deep Learning Approach Based on Explainable Artificial Intelligence for Skin Lesion Classification,” *IEEE Access*, Bd. 10, p. 113715–113725, 2022.
- [5] R. V. Zicari, S. Ahmed, J. Amann, S. A. Braun, J. Brodersen, F. Bruneault, J. Brusseau, E. Campano, M. Coffee, A. Dengel, B. Döder, A. Gallucci, T. K. Gilbert und P. Gottfrois, „Co-Design of a Trustworthy AI System in Healthcare: Deep Learning Based Skin Lesion Classifier,” *Front. Hum. Dyn.*, Bd. 3, 2021.
- [6] G. Argenziano, G. Fabbrocini, P. Carli, V. d. Giorgi, E. Sammarco und M. Delfino, „Epiluminescence microscopy for the diagnosis of doubtful melanocytic skin lesions. Comparison of the ABCD rule of dermoscopy and a new 7-point checklist based on pattern analysis,” *Archives of dermatology*, Bd. 134, Nr. 12, p. 1563–1570, 1998.
- [7] I. Goodfellow, Y. Bengio und A. Courville, *Deep learning (Adaptive computation and machine learning)*, Cambridge: MIT Press, 2016.
- [8] K. Fogelberg, S. Chamarthi, R. C. Maron, J. Niebling und T. J. Brinker, „Domain shifts in dermoscopic skin cancer datasets: Evaluation of essential limitations for clinical translation,” *New biotechnology*, Bd. 76, p. 106–117, 2023.
- [9] A. Lucieri, M. N. Bajwa, S. A. Braun, M. I. Malik, A. Dengel und S. Ahmed, „On Interpretability of Deep Learning based Skin Lesion Classifiers using Concept Activation Vectors,” in *2020 International Joint Conference on Neural Networks (IJCNN)*, Glasgow, UK, 2020.
- [10] J. Kawahara, S. Daneshvar, G. Argenziano und G. Hamarneh, „7-Point Checklist and Skin Lesion Classification using Multi-Task Multi-Modal Neural Nets,” *IEEE journal of biomedical and health informatics*, Bd. 23, Nr. 2, pp. 538–546, 2018.
- [11] N. Codella, V. Rotemberg, P. Tschandl, M. E. Celebi, S. Dusza, D. Gutman, B. Helba, A. Kalloo, K. Liopyris, M. Marchetti, H. Kittler und A. Halpern, „Skin Lesion Analysis Toward Melanoma Detection 2018: A Challenge Hosted by the International Skin Imaging Collaboration (ISIC),“ in *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, Washington, D.C., USA, 2018.