# Deep Learning Based Decision Support for Medicine - A Case Study on Skin Cancer Diagnosis

Adriano Lucieri[a,b,*], Andreas Dengel[a,b], Sheraz Ahmed[a]

[a]*Deutsches Forschungszentrum für Künstliche Intelligenz GmbH (DFKI), Smart Data & Knowledge Services (SDS), Trippstadter Straße 122, 67663 Kaiserslautern, Germany*
[b]*TU Kaiserslautern, Erwin-Schrödinger-Straße 52, 67663 Kaiserslautern, Germany*

**Abstract**

Early detection of skin cancers like melanoma is crucial to ensure high chances of survival for patients. Clinical application of Deep Learning (DL)-based Decision Support Systems (DSS) for skin cancer screening has the potential to improve the quality of patient care. The majority of work in the medical AI community focuses on a diagnosis setting that is mainly relevant for autonomous operation. Practical decision support should, however, go beyond plain diagnosis and provide explanations. This paper provides an overview of works towards explainable, DL-based decision support in medical applications with the example of skin cancer diagnosis from clinical, dermoscopic and histopathologic images. Analysis reveals that comparably little attention is payed to the explanation of histopathologic skin images and that current work is dominated by visual relevance maps as well as dermoscopic feature identification. We conclude that future work should focus on meeting the stakeholder's cognitive concepts, providing exhaustive explanations that combine global and local approaches and leverage diverse modalities. Moreover, the possibility to intervene and guide models in case of misbehaviour is identified as a major step towards successful deployment of AI as DL-based DSS and beyond.

*Keywords:* Decision Support, Explainability, Deep Learning, Medical Image

*Corresponding author

*Email addresses:* adriano.lucieri@dfki.de (Adriano Lucieri), andreas.dengel@dfki.de (Andreas Dengel), sheraz.ahmed@dfki.de (Sheraz Ahmed)

## 1. Introduction

In the last decade, Deep Learning (DL) and particularly the Convolutional Neural Network (CNN) has demonstrated its exceptional ability to efficiently solve a variety of image-based tasks in medical diagnosis such as the detection of skin cancer [1], Alzheimer's [2] and Diabetic Retinopathy (DR) [3]. With IDx-DR [3] and 3DermSpot [4], Digital Diagnostics Inc. have launched the first ever autonomous Artificial Intelligence (AI) diagnostic systems, receiving breakthrough designation by the FDA[1]. Although having been approved in the United States (US) under close cooperation with the FDA [5], the inner workings of DL-based black-box algorithms still raise numerous questions regarding their decision-making processes [6] and unexpected malfunctions [7]. Bissoto et al. [8] explicitly demonstrate that *Clever Hans* behaviour might often be mistaken with high diagnostic performance in the classification of skin cancer from dermoscopic images. In contrast to the US, the strategy pursued by the European Union (EU) is heavily focusing on human oversight, transparency and building an ecosystem of trust [9]. Specially the application of AI in the medical domain bears complex ethical and legal implications which demand the use of explainable and traceable solutions within such an ecosystem [10]. Attempting to explain DL-based systems might also be the best way to validate a system's performance beyond quantitative metrics, preventing issues arising from effects as demonstrated in [8]. The development of explainable AI (xAI) and methods that explain existing AI algorithms is therefore a blooming research area.

Skin cancer is one of the most common types of cancer worldwide [11]. Melanoma is a particularly dangerous type of skin cancer. In the US, it accounts for only 1% of skin cancers diagnosed, but resulting in the largest share of skin cancer related deaths [12]. The manual diagnosis of skin lesions involves

---

[1]Food and Drug Administration

the regular examination of a patient's skin lesions by dermatologists that went through years of specialized training. The cancer is diagnosed by thoroughly analyzing skin lesions, applying rules of dermoscopic pattern recognition and algorithms like the ABCD-rule [13] or 7-point checklist [14] that have been developed and described in medical literature throughout years of research. This routine process requires constant attention of the expert. Therefore, diagnostic performance is highly dependent on fatigue and emotional state of the physician. This and the subjectivity of dermoscopic criteria recognition might be contributing to the fact that doctor's decisions are hard to reproduce [15]. The augmentation of human doctors through AI-based Decision Support Systems (DSS) is highly desirable as it potentially increases reproducibility of results, speeding up tedious examinations, and therefore allowing more thorough treatment of the broader population. Fortunately, initiatives for public large-scale repositories of digital skin images such as the International Skin Imaging Collaboration (ISIC)[2] drove an increasing interest of AI application in this domain.

This work provides an overview of current explainability methods for DLS applied to the problem of skin lesion classification. The critical assessment of current approaches ought to present the bigger picture of AI in dermatology to identify missing traits that impede the deployment of AI assistants in the medical domain. Section 2 provides basic knowledge of skin lesion classification and common explainability methods for DL-based classifiers. Then, the reviewed work is summarized and classified into four non-exclusive groups describing the main explanation techniques in Section 3. Section 4 critically assesses the current state of explainable DL-based skin lesion classification. Moreover, four major objectives to be followed for clinically practical DLS in medicine are proposed, followed by concluding remarks in Section 5.

---

[2]https://www.isic-archive.com/

## 2. Background

### 2.1. Skin Cancer Diagnosis

The classification of skin lesions is usually attempted by analysing one of the three most popular domains, i.e. clinical, dermoscopic and histopathologic imaging. **Clinical images** are digitized naked-eye observations of the skin and exhibit the lowest information value. **Dermoscopic images** are taken using a dermatoscope, which uses light and high magnification rates to visualize patterns present in upper layers of skin in more detail. **Histopathology** is the microscopic examination of cell tissue. Whole Slide Images (WSIs) are detailed images of enormous size (up to Gigapixels) that are generated by scanning glass slides of skin tissue. Histopathologic examination is regarded as the gold standard in skin lesion diagnosis, whereas the information content declines for dermoscopic and clinical images accordingly.

### 2.2. XAI Taxonomy for Skin Cancer Diagnosis

A variety of taxonomies for xAI methods exist in literature [16, 17]. In contrast to existing categorizations for the broad application of xAI methods, we group the reviewed work in four main categories, i.e. *Visual Relevance Localization*, *Dermoscopic Feature Prediction & Localization*, *Similarity Retrieval* and *Intervention*. This distinction is specifically selected to fit the dermatological use-case as well as being user-centric, focusing on the utility of the output rather than the producing mechanisms. An overview of the reviewed works in the four non-exclusive categories is presented in Fig. 1.

*Visual Relevance Localization.* A popular approach towards explaining CNNs is the generation of visual maps that provide spatial information about the input's relevance to a prediction. We further distinguish post-hoc, relevance-based methods that visualize network's activations (*Activation* e.g. [48, 49]) or trace the input's relevance to prediction (*Attribution*, e.g. [50, 51, 52, 53, 54, 55, 56, 57]) from attention-based methods (*Attention*, e.g. [58, 59, 60]) that

4

Attempts on Explainable
AI in Dermatology

| Visual Relevance Localization | | Similarity Retrieval |
|---|---|---|

**Visual Relevance Localization**

**Activation**
Van Molle et al. [18]
Barata et al. [19]

**Attention**
Jia et al. [25]
Yan et al. [38]
Zhang et al. [39]
Barata et al. [19]

**Attribution**
Cruz-Roa et al. [20]          Yang et al. [29]
Cruz-Roa et al. [21]          Mishra et al. [30]
Radhakrishnan et al. [22]     Xiang et al. [31]
Esteva et al. [23]            Rieger et al. [32]
Ge et al. [24]                Young et al. [33]
Jia et al. [25]               Xie et al. [34]
Li et al. [26]                Mikołajczyk [35]
Codella et al. [27]           Sonntag et al. [36]
Thandiackal et al. [28]       Hägele et al. [37]

**Similarity Retrieval**
Kawahara et al. [40]
Codella et al. [27]

**Intervention**
Yan et al. [38]
Rieger et al. [32]
Mikołajczyk [35]
Chen et al. [46]

**Dermoscopic Feature Prediction & Localization**

**Prediction**
Kawahara et al. [41]      Lucieri et al. [44]
Codella et al. [27]       Coppola et al. [45]
Veltmeijer et al. [42]    Chen et al. [46]
Murabayashi et al. [43]

**Localization**
Kawahara et al. [40]      Zhang et al. [47]
Kawahara et al. [41]      Veltmeijer et al. [42]
Thandiackal et al. [28]   Sonntag et al. [36]

Figure 1: Overview of reviewed publications, grouped by the proposed taxonomy. Single papers can belong to multiple groups.

are usually incorporated into the network architecture to enforce a network's attention during training.

*Dermoscopic Feature Prediction & Localization.* These methods aim at predicting whether dermoscopic features, relevant to the classification of skin diseases, are present in the input image or their localization. This type of explainability is often sometimes approached by means of multitask learning [61], which is problematic. A more sophisticated method that identifies abstract concepts in a pre-trained classifier was proposed by Kim et al. [62].

*Similarity Retrieval.* Another way of making predictions more intelligible is by allowing the user to retrieve cases, that share relevant similarities according to the classifier. Content-based Image Retrieval (CBIR) systems are usually trained for the very purpose of finding a semantically meaningful representation to retrieve similar images [63]. However, the notion of CBIR can be applied to models trained on the primary objective of classification as well [64].

*Intervention.* The last group of methods is characterized by actively intervening and improving a model's explanations. These *Intervention* methods include

e.g. the penalization of wrong explanations during the training process [32] or identification of biases in the data as well as methods to remove them [65]. The continuous improvement of AI-based DSS in the health sector during their development and life cycle is crucial for their successful application [66] and is therefore considered as a complement to explanation methods.

## 3. Overview of Methods

Tables 1 and 2 provide a comprehensive summary of all reviewed methods, including informations on the dataset, additional annotations and methods used.

### 3.1. Visual Relevance Localization

### 3.1.1. Visual Activation

Various methods have been proposed for the visualization of a CNNs intermediate feature representations (e.g. [48], [49]). Van Molle et al. [18] visualized the activation maps of a custom CNN trained on the ISIC archive by rescaling and mapping them onto the input image. They discovered separate activation maps excited by surrounding skin, hair-like features, variations in lesion colour and lesion borders. However, no activation maps have been found, focusing on common dermoscopic criteria as defined by dermatologists.

### 3.1.2. Visual Attribution

As early as 2013, Cruz-Roa et al. [20] worked on automated detection of BCC on the HistologyDS[3] dataset using a patch-based DL approach. Their system processes WSIs patch-wise at a resolution of $8 \times 8$ through an encoder, a CNN and a final softmax classifier to get a prediction. A *Digital Staining* procedure generates heatmaps on top of the input image by weighting each image patch with its resulting prediction score. In 2014, the authors proposed a revised method for detection of Invasive Ductal Carcinoma (IDC) on a private dataset in [21]. Again, they applied a supervised, patch-base approach with a

---

[3]Available at: http://www.informed.unal.edu.co/histologyDS

Table 1: Clinical images and Dermoscopy: Comparison of explanatory approaches for clincial and dermoscopic image classifiers. Any non-public dataset used has been marked as private. Additional manual labeling and the use of quantitative measures of explainability are marked in separate columns.

| Author | Year | Type | | Datasets | | | | | Added Annotations | Method | | | | | | | Quantification | Evaluation | Comment |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | Visual | | | Feature | | | | | | |
| | | Clinical | Dermoscopic | ISIC Archive | Dermofit | Derm7pt | PH² | Private | | Activation | Attribution | Attention | Similarity | Prediction | Localization | Intervention | | | |
| Radhakrishnan et al. [22] | 17 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | – | – | – | – | – | – | ✓ | Pseudo-Importance; Coarse patch-wise heatmaps. |
| Esteva et al. [23] | 17 | ✓ | ✓ | ✓ | ✓ | – | – | ✓ | – | – | ✓ | – | – | – | – | – | – | – | Basic attribution method (Saliency). |
| Ge et al. [24] | 17 | ✓ | ✓ | – | – | – | – | ✓ | – | – | ✓ | – | – | – | – | – | – | – | Coarse heatmaps (CAM-BP). |
| Jia et al. [25] | 17 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | ✓ | – | – | – | – | – | – | Coarse heatmaps (CAM). |
| Kawahara et al. [41] | 18 | – | ✓ | ✓ | – | – | – | – | – | – | – | – | – | – | ✓ | – | – | ✓ | Dermoscopic feature segmentation. |
| Kawahara et al. [40] | 18 | ✓ | ✓ | – | – | ✓ | – | – | – | – | – | – | ✓ | ✓ | ✓ | – | – | – | Use of metadata. Multi-task. |
| Li et al. [26] | 18 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | – | – | – | – | – | – | – | Computationally intensive attribution. |
| Codella et al. [27] | 18 | – | ✓ | – | – | – | – | – | – | – | ✓ | – | ✓ | ✓ | – | – | – | – | Coarse heatmaps (QAM & RAM). |
| Thandiackal et al. [28] | 18 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | – | – | – | ✓ | – | ✓ | ✓ | Multi-task. |
| Van Molle et al. [18] | 18 | – | ✓ | ✓ | – | – | – | – | – | ✓ | – | – | – | – | – | – | – | – | |
| Yang et al. [29] | 18 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | – | – | – | – | – | – | – | Coarse heatmap (CAM); Multi-task. |
| Yan et al. [38] | 19 | – | ✓ | ✓ | – | – | – | – | – | – | – | ✓ | – | – | – | ✓ | – | ✓ | |
| Zhang et al. [47] | 19 | – | ✓ | ✓ | – | – | – | – | – | – | – | – | – | – | ✓ | – | – | – | Unsupervised biomarker localization; No distinction between biomarkers. |
| Mishra et al. [30] | 19 | – | ✓ | – | – | – | – | ✓ | – | – | ✓ | – | – | – | – | – | – | – | Coarse heatmaps (Grad-CAM, GBP). |
| Veltmeijer et al. [42] | 19 | – | ✓ | ✓ | – | – | – | – | ✓ | – | – | – | – | ✓ | ✓ | – | – | – | Multi-task. |
| Murabayashi et al. [43] | 19 | – | ✓ | ✓ | – | – | – | ✓ | – | – | – | – | – | ✓ | – | – | – | – | Multi-task; Semi-supervised VAT. |
| Xiang et al. [31] | 19 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | – | – | – | – | – | – | – | LIME attribution. |
| Rieger et al. [32] | 19 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | – | – | – | – | ✓ | – | – | Coarse heatmap (Grad-CAM). |
| Zhang et al. [67] | 19 | – | ✓ | ✓ | – | – | – | – | – | – | – | ✓ | – | – | – | – | – | – | Multi-stage attention. |
| Young et al. [33] | 19 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | – | – | – | – | – | – | – | Grad-CAM; Kernel SHAP. |
| Mikołajczyk et al. [35] | 20 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | – | – | – | – | ✓ | – | – | Global explanation. |
| Sonntag et al. [36] | 20 | – | ✓ | ✓ | – | – | – | – | – | – | ✓ | – | – | – | ✓ | – | – | – | Coarse heatmaps (Grad-CAM, RISE); Multi-task. |
| Barata et al. [19] | 20 | – | ✓ | ✓ | – | – | – | – | – | ✓ | – | ✓ | – | – | – | – | – | – | Hierarchical approach. |
| Lucieri et al. [44] | 20 | – | ✓ | ✓ | – | ✓ | ✓ | ✓ | – | – | – | – | – | ✓ | – | – | ✓ | – | Feature influence analysis on unconstrained DLS; Global explanation. |
| Coppola et al. [45] | 20 | – | ✓ | – | – | ✓ | – | – | – | – | – | – | – | ✓ | – | – | – | – | Multi-task. |
| Chen et al. [46] | 20 | – | ✓ | ✓ | – | – | – | – | – | – | – | – | – | ✓ | – | ✓ | ✓ | – | Allows for concept discovery. |
| Σ | | 3 | 26 | 22 | 1 | 3 | 1 | 5 | 1 | 2 | 14 | 4 | 2 | 7 | 6 | 4 | 3 | 4 | |

Table 2: Histopathology: Comparison of explanatory approaches for skin tissue classifiers. Any non-public dataset used has been marked as private. Additional manual labeling on public data and the use of quantitative measures of explainability are noted in separate columns.

| Author | | Datasets | | | Added Annotations | Method | | | | | | | Quantification | Comments |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | Visual | | | | Feature | | | | |
| | Year | TCGA | histologyDS | Private | Added Annotations | Activation | Attribution | Attention | Similarity | Prediction | Localization | Intervention | Quantification | |
| **Cruz-Roa et al. [20]** | 13 | – | ✓ | ✓ | – | – | ✓ | – | – | – | – | – | – | Pseudo-Importance; Coarse, patch-wise heatmap. |
| **Cruz-Roa et al. [21]** | 14 | – | – | ✓ | – | – | ✓ | – | – | – | – | – | – | Pseudo-Importance; Coarse, patch-wise heatmap. |
| **Xie et al. [34]** | 19 | ✓ | – | ✓ | – | – | ✓ | – | – | – | – | – | – | Coarse heatmaps (CAM, Grad-CAM). |
| **Hägele et al. [37]** | 20 | ✓ | – | – | ✓ | – | ✓ | – | – | – | – | – | ✓ | Fine-grained heatmaps (LRP). |
| Σ | | **2** | **1** | **3** | **1** | **0** | **4** | **0** | **0** | **0** | **0** | **0** | **1** | |

CNN to generate an IDC pseudo-probability map for visual analysis. The 2017 work from Esteva et al. [23] is well cited for being the first DL-based classifier outperforming 21 board-certified dermatologists in the classification of skin cancer. The authors used the GoogleNet InceptionV3 [68] architecture, pre-trained on ImageNet [69], and fine-tuned on clinical and dermoscopic data from ISIC Archive, Dermofit and data from Stanford Hospital. Explanations are provided through saliency maps generated as gradients of the loss function w.r.t. the input. In the same year, Radhakrishnan et al. [22] proposed PatchNet, a patch-based architecture that processes each patch of a dermoscopic input image using a CNN. The patch-based paradigm is leveraged to compute attribution maps by combining the scores of every patch in the binary task of benign vs. malignant. Resulting attribution maps express the network's pseudo-probability of a patch being malignant. The patch-wise heatmaps are compared to Class Activation Maps (CAM) [51] and Grad-CAM [52] and quantitatively evaluated on segmentations of dermoscopic features provided in the ISIC 2017 challenge dataset. Ge et al. [24] explored different multi-modal network structures, classifying pairs of clinical and dermoscopic images in one of 15 skin conditions. The private dataset used was provided by MoleMap[4]. Moreover, they propose CAM-BP

---

[4]https://www.molemap.co.nz/

which generates attribution maps by weighting Bilinear Pooling (BP) [70] using CAM. In [26], the authors applied the perturbation-based *Prediction Difference Analysis* method proposed in [55] to produce visual attribution maps. In 2018, Yang et al. [29] proposed *Region Average Pooling* to improve dermoscopy image classification by jointly segmenting the lesion, using this segmentation map as ROI for pooling. CAM heatmaps were provided as well to prove the classifier's focus on the lesion region. In 2019, Xie et al. [34] used VGG19 [71] and ResNet50 [72] networks to classify WSIs from the TCGA[5] dataset into Melanoma and Nevi. CAM attribution has been computed to explain the classifiers. Mishra et al. [30] trained several ResNet architectures and applied Grad-CAM [52] as well as Guided Backpropagation [53] to analyze their models. The authors found that the models were having trouble with bad lighting, image blur, image noise and a wide field of view. Xiang et al. [31] trained an ensemble of different DL-based architectures. Explainability of their classifiers has been attempted by applying LIME [57] to obtain visual attribution maps. It was confirmed that the classifiers made use of the lesion region for classification. Young et al. [33] trained 30 models on Melanoma and Nevi from a publicly available dermoscopic dataset and applied GradCAM and KernelSHAP [73] for explanation. Despite high accuracies, models sometimes assigned relevance to features that were non-related to the actual classification task. Such artifacts can arise through process-specific conditions during clinical examination, bearing the risk of causing spurious correlations in the data. Well-known examples from public skin lesion datasets are the appearance of vignetting effects, coloured patches, scales or markings on skin lesion images as shown in Fig. 2. Recently, Hägele et al. [37] emphasized on the importance of high-resolution attribution maps for the explanation and disclosure of biases in DL-based classifiers for histopathology. They compare high-resolution Layer-Wise Relevance Propagation (LRP) [56] maps to probability maps and GradCAM, demonstrating how fine-grained solutions allow explanations on cell-level instead of patch-level and how this is

---
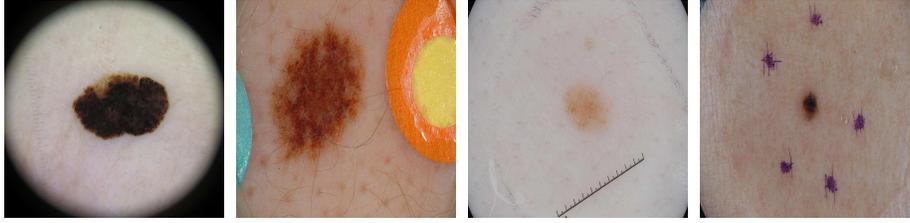
[5]Available at: https://portal.gdc.cancer.gov/

Figure 2: Examples of skin images suffering from different artifacts like (from left to right) vignetting effects, coloured patches, scale, markings.

useful to discover biases in classifiers and data. To this end, they even provide a quantitative comparison of LRP and Grad-CAM by means of the area under the receiver operating characteristic curve (AUROC) for the detection sensitivity of cancer cells.

### 3.1.3. Visual Attention

In 2017, Jia et al. [25] developed a two-stage framework using a single CNN architecture that is first trained on ISIC 2017 dataset images. Then, Class Activation Maps are computed to crop images to the respective region of interest (ROI). The same CNN is then again trained with cropped images. The resulting classifier combines attribution and explicit attention for its explanation. In the following year, Zhang et al. [67] proposed an *Attention Residual Learning* (ARL) CNN that provides attention maps at every ARL block stage. Yan et al. [38] incorporated an attention mechanism into a VGG16 network to obtain attention maps. In 2020, Barata et al. [19] treated the problem of skin lesion classification as a hierarchical problem of sequential word prediction. Using a CNN for feature encoding and a softmax attention mechanism for channel-wise and spatial attention masks, a masked embedding of the original image is fed into a recurrent long short-term memory (LSTM) module that sequentially predicts increasingly specific classes (Melanocytic vs. Non-Melanocytic; Benign vs. Malignant; Diagnosis). The attention maps were visualized, aiding the interpretability of the model. Gu et al. [74] propose a comprehensive attention CNN (CA-Net) for medical image segmentation that combines the advantages of

spatial, channel and scale attention to understand and interpret how pixel-level predictions are obtained.

## 3.2. Dermoscopic Feature Prediction & Localization

In 2018, Kawahara et al. [41] presented a fully convolutional architecture for the segmentation of four dermoscopic features annotated in the 2017 ISIC challenge dataset. Results are evaluated using a fuzzy Jaccard Index. In further work, Kawahara et al. [40] developed a multi-modal DLS trained in a multi-task setting for the simultaneous detection and localization of diagnosis class and dermoscopic criteria. The CNN architecture is trained on sets of clinical and dermoscopic images along with meta data. It provides localization maps for all dermoscopic criteria, differianted in sub-types (e.g. regular / irregular streaks), can handle missing data (e.g. dermoscopic image without meta data) during inference due to the use of multiple loss combinations at train time and serves as a multi-modal image retrieval system. Thandiackal et al. [28] trained U-Nets [75] to segment dermoscopic features. Those segmentation maps were stacked channel-wise with the original image and fed into a final CNN classifier . Segmentation results are evaluated against the results from [41]. Additionally, they allow to quantify the amount of network attribution (based on Input × Gradient) channel-wise and location-wise using partial derivatives. In 2019, Zhang et al. [47] attempted at localizing biomarkers without explicit labels. They trained a generative adversarial network (GAN) [76] on removing potential biomarkers from diseased images, to make them appear like benign images. Those fake benign images were subtracted from the GAN's input to localize potential dermoscopic criteria. The resulting representation was fed to a subsequent classifier to predict the initial image's class label. Veltmeijer et al. [42] manually annotated three dermoscopic criteria (Asymmetry, Irregular Borders and Multiple Colours) in the 2017 ISIC dataset. Using those annotations they trained two multi-task networks of which one predicted the criteria and the other segmented it, along with the diagnosis task. In [43], Murabayashi et al. made use of semi-supervised *Virtual Adversarial Training* (VAT) to leverage

unlabeled data, training a multi-task CNN on the prediction of diagnosis labels and dermoscopic criteria. In 2020, Sonntag et al. [36] proposed a DSS for the explainable classification of skin lesions. The system allows lesion segmentation, dermoscopic criteria segmentation, binary or 8-class classification as well as the generation of visual attribution maps using GradCAM or RISE [54]. Lucieri et al. [44] used TCAV [62], a concept-based explanation method, to find directions of dermoscopic features in the latent space of a trained CNN. Unlike all previously described methods, [44] analyse a trained CNN without constraining its training procedure by explicitly enforcing dermoscopic criteria learning. The authors found that a publicly available skin lesion classifier module, scoring third in the 2017 ISIC lesion classification challenge, is naturally able to distinguish between absence and presence of some common dermoscopic criteria defined by dermatologists. Furthermore, it was found that those criteria influence the CNNs decision in accordance with the medical consensus. Coppola et al. [45] recently proposed a novel multi-task learning network with learnable gates that allow to share features between subtasks. They discovered correlations between the features used for the detection of specific dermoscopic criteria and concluded that the diagnosis tasks uses features from most of the other subtasks, which confirms their intuition.

### 3.3. Similarity Retrieval

Codella et al. [27] were the first to explicitly expand the explanation of dermoscopic classifiers by allowing content-based image retrieval (CBIR). By manually labeling images in non-expert similarity groups (e.g. visual appearance) and employing a triplet loss, their network's latent space is structured to provide meaningful nearest neighbours as evidence for its classification. In addition, they generate *Result Activation Maps* (RAM) based on CAM, highlighting regions that contribute most to the proximity of the nearest neighbour of a sample. The previously mentioned multi-modal DLS by Kawahara et al. [40] also allowed image retrieval. A special property is that a single feature vector can be used to retrieve multi-modal training inputs.

12

### 3.4. Intervention

In 2019, Yan et al. [38] improved the classification performance of their CNN by explicitly guiding the attention of the network through ground truth segmentation maps during training. Rieger et al. [32] introduced *Contextual Decomposition Explanation Penalization* (CDEP). CDEP is a model-inherent method that allows to constrain the explanation of a deep network during training using ground-truth segmentation maps. As a result, bias can be removed iteratively during optimization, resulting in an explainable, accurate network. GradCAM is used to generate attribution maps. In 2020, Mikołajczyk et al. [35] attempted global explanations by proposing *Global Explanations for Discovering Bias in Data* (GEBI). The idea is based on *SpRAy* [65] which applies spectral clustering on a CNN's attribution maps to explain its behaviour. In [35], this method is extended by concatenating the input image to the corresponding attribution map before spectral clustering. Moreover, a counterfactual method is proposed in which a hypothesized bias is artificially induced in all images of the dataset to observe the change in predictive behaviour. Variations of LRP are used as attribution method. Building on the idea of TCAV, Chen et al. [46] propose Concept Whitening (CW) as a normalization technique which is able to align a network's latent space with a set of desired concepts. Training a ResNet18 with CW on ISIC skin lesion images slightly improved the classification performance and resulted in a significantly more disentangled latent space. The authors concluded that the network did not focus on the age of a patient but instead on the size of the lesion. CW moreover reveals other concept directions that do not correspond to predefined concepts, simplifying the discovery of new, unknown biomarkers. The most influential concept axis in their experiments indicated the utilization of irregular borders for the detection of malignity.

## 4. Discussion

Tables 1 and 2 summarize all reviewed publications including additional information like datasets used, manual annotations, whether the approach allowed

13

quantitative explanation and whether the explanations have been evaluated. It appears that most works aimed at explaining dermoscopic image classifiers. In spite of numerous publications on the explanation of DL-based classifiers for histopathology slides [39, 77, 78, 79], only little attention has been drawn to skin tissue in particular. This might also be related to the limited availability of datasets and annotations, as is also reflected in the limited number of publicly available datasets used by the reviewed studies.

It appears that two explanation modalities dominate previous work, namely visual attribution and dermoscopic criteria prediction & localization. Visualization has always been a particularly popular way of explaining CNNs due to the importance of spatial arrangement in feature maps and visual input. However, the utility of different visualization methods varies a lot with their implementation. Numerous gradient-based attribution methods, for instance, have been shown to be independent of the model at hand [80] (e.g. Guided Backpropagation). Also, some methods produce coarse heatmaps owing to their implementation (e.g. CAM, GradCAM). Perturbation-based methods are known to suffer from out-of-distribution sampling which can significantly influence explanations [81]. Attention methods often explicitly enforce focus on specific image regions, thus limiting the information used by a DLS. Other approaches providing non-constraining attention cannot guarantee influence of their provided explanations on the prediction. The same holds for prediction and localization of dermoscopic criteria through multi-task optimization as well as simple activation visualization. Some of the reviewed methods used the predicted scores of sub-regions (e.g. patches) to explain the network decision spatially. However, an opaque and uninterpretable high-level sub-decision can hardly be considered an explanation.

The implicit prediction and localization of dermoscopic criteria as practiced in [46] and [44] is closely related to the utilization of auxiliary classification objectives and hierarchization of classifier structures. These methods implicitly allow to monitor the DLS' understanding of predictions by verifying that a set of simpler sub-tasks are correctly performed and allow to quantify the influence of a

sub-task to the final classifcation. The localization of single criteria is specifically useful to increase the efficiency of doctors by leading their attention and allowing them to validate network's decisions more easily. The Concept Localization Map (CLM) method proposed in [82] allows to extend frameworks based on CAVs for the implicit localization of biomarkers. Except for the method proposed in [47] and [46], all methods require the manual definition and annotation of criteria. This constrains the model performance in favour of explainability and constitutes a complex and laborious task, requiring hours of work from experts.

Attempts beyond visual and conceptual interpretability like content-based image retrieval approaches and methods that allow to explicitly adjust DLS's behaviour represent a significant step towards clinical usefulness of such systems. Some of the reviewed methods, for instance, allow to explicitly guide a network's focus and explanation through the incorporation of domain-knowledge or by detecting bias in the training data that can consequently be removed.

Analysis of existing work revealed that there is an ongoing trend of user-centered and adaptable explainable AI not only in the dermatology domain, but in AI in general. We suggest that future attempts to successfully deploy XAI in medicine should focus on following key points:

1. User-centered explanations
2. Diverse explanations
3. Global explanations
4. Interventions

*User-Centered Explanations.* The process of interpreting explanations is sometimes described as a translation of abstract concepts into a human-understandable domain [83]. In order to create an understanding, it is essential to transfer knowledge in a way that is accessible to the human interpreter. Thus, a user-centric orientation of explanations is a major requirement for effective explanations. The easiest way to fullfil this requirement is the communication by means of commonly used mental constructs (concepts) in the explainee's cultural milieu or the domain of expertise. This can, for instance, be implemented by

15

explicitly handling intermediate tasks through modularization of architectures or the disclosure of such concepts in unconstrained DL-based algorithms.

*Diverse Explanations.* Another crucial point that can be observed in the explanation of human beings and that will aid the interpretation of algorithmic decisions is the utilization of diverse modalities for explanation. Spatial localization of important input features using attribution maps is often helpful. However, only a fraction of the information necessary to explain a decision in its entirety is conveyed this way. The augmentation of visual relevance maps by more abstract, concept-based explanations and more fine-grained concept localization can result in an enriched and more complete explanation. Additionally, existing explanation approaches can be combined with different input modalities of the same data to reveal new relationships and increase diversity. Novel ideas considering input data in different data spaces such as frequency domain (e.g. [84]), uncommon colour domains or other transformations (e.g. [85]) could potentially result in the alternative expression of relevance beyond spatial location. The combination of different explanation methods with diverse modalities on one hand allows for more complete explanations and on the other hand for the validation of the different explanations' congruence, hence the meaningfulness of the whole explanation.

*Global Explanations.* Most of today's xAI approaches consider only local explanations. However, a proper validation of network's working that will gain the trust of users and authorities can only be achieved by analyzing the bigger picture. In a DL setting, this means that all of the data space needs to be taken into account. So far, the consideration of the complete input space is - in most cases - impossible. Therefore, methods for the approximation of data spaces are required which in the best case provide some theoretical guarantees. As has been shown, current global explanations such as concept-based methods used in [44] or the data-centric method proposed in [35] approach global explanations by approximating the input space with as much data available. The combination of local and approximated global measures bears the potential to increase

the scope of the explanation and therefore the user's understanding. Without any guarantees of causality, assuring that a statistical model learnt only relevant coherences is only possible by identifying and removing systemic bias from the representation space.

*Interventions.* Ultimately, an explainable classifier is not worth much if it does not base its decision on legitimate grounds. Therefore, intervention is required to correct issues identified through explanations. Moreover, this is a way of explicitly incorporating expert knowledge into a learning system. The requirement for extensive annotation collection that is often seen as a drawback can be alleviated by applying classifiers with good base performance in clinical support settings, making use of interventional active learning. Here, the DLS represents the learner and an expert can act as a tutor that communicates potential misconceptions to the network on a case-by-case basis. The tutor benefits from the good and consistent performance of the learner, who in turn benefits form the correction of misconceptions. If proper channels of communication from tutor to student are provided, such scenarios might be comparable to a tutor teaching a real doctor in training.

Looking towards the future, another path should be followed in addition to the supervised incorporation of expert knowledge. As already mentioned, the definition and incorporation of human concepts and ideas into computer-understandable language is not trivial and requires exceptional efforts. Moreover, it is also known that many human decisions, processes and conventions emerged through time and are not always based on optimal solutions. Hence, it is not unlikely that there exist alternative representations of problems that, so far, seem abstract to humans but are well-grounded in theory. Decrypting the representations learned by high performing DL-based classifiers is therefore still of great potential value. First attempts towards unsupervised concept discovery have been proposed in [86, 87, 88]. The most critical aspect is however the assignment of meaning exploitable by humans.

## 5. Conclusions

The present work reviewed and discussed the current state-of-the-art in explaining DL-based skin lesion classifiers. Four non-exclusive main groups of explanation approaches have been identified, namely, *Visual Relevance Localization*, *Dermoscopic Feature Prediction & Localization*, *Similarity Retrieval* as well as *Intervention*. It has been shown that the explanation of histopathologic skin lesion classifiers has so far received little attention and that the main focus in the field lies on the explanation of dermoscopic classifiers by means of visual relevance maps and prediction & localization of dermoscopic criteria. In order to meaningfully advance the field of medical AI towards practical, clinical deployment of DLS, we suggest to focus future efforts on user-centered explanation attempts that combine diverse modalities with local and global evidence. The continuous shaping of a classifier's behaviour through interventional active learning has the potential to revolutionize the way knowledge is transferred from human experts to DLS. Moreover, first attempts towards the unsupervised discovery of knowledge learned by data-driven algorithms indicate promising prospects for both, DL and the medical communities.

## References

[1] T. J. Brinker, A. Hekler, A. H. Enk, J. Klode, A. Hauschild, C. Berking, B. Schilling, S. Haferkamp, D. Schadendorf, T. Holland-Letz, et al., Deep learning outperformed 136 of 157 dermatologists in a head-to-head dermoscopic melanoma image classification task, European Journal of Cancer 113 (2019) 47–54.

[2] S. Qiu, P. S. Joshi, M. I. Miller, C. Xue, X. Zhou, C. Karjadi, G. H. Chang, A. S. Joshi, B. Dwyer, S. Zhu, et al., Development and validation of an interpretable deep learning framework for alzheimer's disease classification, Brain.

[3] M. D. Abràmoff, P. T. Lavin, M. Birch, N. Shah, J. C. Folk, Pivotal trial of an autonomous ai-based diagnostic system for detection of diabetic retinopathy in primary care offices, NPJ digital medicine 1 (1) (2018) 1–8.

[4] P. Newswire, Fda permits marketing of artificial intelligence-based device to detect certain diabetes-related e www.prnewswire.com.
URL https://www.prnewswire.com/news-releases/3derm-announces-two-fda-breakthrough-device

[5] A. Stark, Fda permits marketing of artificial intelligence-based device to detect certain diabetes-related e
U.S. Food and Drug Administration (FDA).
URL https://www.fda.gov/news-events/press-announcements/fda-permits-marketing-artificial

[6] C. Rudin, Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead, Nature Machine Intelligence 1 (5) (2019) 206–215.

[7] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, R. Fergus, Intriguing properties of neural networks, in: International Conference on Learning Representations, 2014.
URL http://arxiv.org/abs/1312.6199

[8] A. Bissoto, M. Fornaciali, E. Valle, S. Avila, (de) constructing bias on skin lesion datasets, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2019, pp. 0–0.

[9] E. Commission, White paper on artificial intelligence - a european approach to excellence and trust.

[10] D. Schneeberger, K. Stöger, A. Holzinger, The european legal framework for medical ai, in: International Cross-Domain Conference for Machine Learning and Knowledge Extraction, Springer, 2020, pp. 209–226.

[11] Z. Khazaei, F. Ghorat, A. Jarrahi, H. Adineh, M. Sohrabivafa, E. Goodarzi, Global incidence and mortality of skin cancer by histological subtype and its relationship with the human development index (hdi); an ecology study in 2018, World Cancer Research Journal (WCRJ) 6 (2) (2019) e13.

[12] A. C. Society, Cancer facts & figures 2020, Am. Cancer Soc.

[13] F. Nachbar, W. Stolz, T. Merkle, A. B. Cognetta, T. Vogt, M. Landthaler, P. Bilek, O. Braun-Falco, G. Plewig, The abcd rule of dermatoscopy: high prospective value in the diagnosis of doubtful melanocytic skin lesions, Journal of the American Academy of Dermatology 30 (4) (1994) 551–559.

[14] G. Argenziano, G. Fabbrocini, P. Carli, V. De Giorgi, E. Sammarco, M. Delfino, Epiluminescence microscopy for the diagnosis of doubtful melanocytic skin lesions: comparison of the abcd rule of dermatoscopy and a new 7-point checklist based on pattern analysis, Archives of dermatology 134 (12) (1998) 1563–1570.

[15] J. G. Elmore, R. L. Barnhill, D. E. Elder, G. M. Longton, M. S. Pepe, L. M. Reisch, P. A. Carney, L. J. Titus, H. D. Nelson, T. Onega, et al., Pathologists' diagnosis of invasive melanoma and melanocytic prolifera- tions: observer accuracy and reproducibility study, Bmj 357 (2017) j2813.

[16] V. Arya, R. K. Bellamy, P.-Y. Chen, A. Dhurandhar, M. Hind, S. C. Hoff- man, S. Houde, Q. V. Liao, R. Luss, A. Mojsilović, et al., One explanation does not fit all: A toolkit and taxonomy of ai explainability techniques, arXiv preprint arXiv:1909.03012.

[17] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Bar- bado, S. García, S. Gil-López, D. Molina, R. Benjamins, et al., Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and chal- lenges toward responsible ai, Information Fusion 58 (2020) 82–115.

[18] P. Van Molle, M. De Strooper, T. Verbelen, B. Vankeirsbilck, P. Simoens, B. Dhoedt, Visualizing convolutional neural networks to improve decision support for skin lesion classification, in: Understanding and Interpreting Machine Learning in Medical Image Computing Applications, Springer, 2018, pp. 115–123.

[19] C. Barata, M. E. Celebi, J. S. Marques, Explainable skin lesion diagnosis using taxonomies, Pattern Recognition (2020) 107413.

[20] A. A. Cruz-Roa, J. E. A. Ovalle, A. Madabhushi, F. A. G. Osorio, A deep learning architecture for image representation, visual interpretability and automated basal-cell carcinoma cancer detection, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2013, pp. 403–410.

[21] A. Cruz-Roa, A. Basavanhally, F. González, H. Gilmore, M. Feldman, S. Ganesan, N. Shih, J. Tomaszewski, A. Madabhushi, Automatic detection of invasive ductal carcinoma in whole slide images with convolutional neural networks, in: Medical Imaging 2014: Digital Pathology, Vol. 9041, International Society for Optics and Photonics, 2014, p. 904103.

[22] A. Radhakrishnan, C. Durham, A. Soylemezoglu, C. Uhler, Patchnet: Interpretable neural networks for image classification, arXiv preprint arXiv:1705.08078.

[23] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, S. Thrun, Dermatologist-level classification of skin cancer with deep neural networks, nature 542 (7639) (2017) 115–118.

[24] Z. Ge, S. Demyanov, R. Chakravorty, A. Bowling, R. Garnavi, Skin disease recognition using deep saliency features and multimodal learning of dermoscopy and clinical images, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2017, pp. 250–258.

[25] X. Jia, L. Shen, Skin lesion classification using class activation map, arXiv preprint arXiv:1703.01053.

[26] X. Li, J. Wu, E. Z. Chen, H. Jiang, What evidence does deep learning model use to classify skin lesions?, arXiv preprint arXiv:1811.01051.

[27] N. C. Codella, C.-C. Lin, A. Halpern, M. Hind, R. Feris, J. R. Smith, Collaborative human-ai (chai): Evidence-based interpretable melanoma classification in dermoscopic images, in: Understanding and Interpreting Machine Learning in Medical Image Computing Applications, Springer, 2018, pp. 97–105.

[28] K. Thandiackal, O. Goksel, A structure-aware convolutional neural network for skin lesion classification, in: OR 2.0 Context-Aware Operating Theaters, Computer Assisted Robotic Endoscopy, Clinical Image-Based Procedures, and Skin Image Analysis, Springer, 2018, pp. 312–319.

[29] J. Yang, F. Xie, H. Fan, Z. Jiang, J. Liu, Classification for dermoscopy images using convolutional neural networks based on region average pooling, IEEE Access 6 (2018) 65130–65138.

[30] S. Mishra, H. Imaizumi, T. Yamasaki, Interpreting fine-grained dermatological classification by deep learning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2019, pp. 0–0.

[31] A. Xiang, F. Wang, Towards interpretable skin lesion classification with deep learning models, in: AMIA Annual Symposium Proceedings, Vol. 2019, American Medical Informatics Association, 2019, p. 1246.

[32] L. Rieger, C. Singh, W. J. Murdoch, B. Yu, Interpretations are useful: penalizing explanations to align neural networks with prior knowledge, arXiv preprint arXiv:1909.13584.

[33] K. Young, G. Booth, B. Simpson, R. Dutton, S. Shrapnel, Deep neural network or dermatologist?, in: Interpretability of Machine Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision Support, Springer, 2019, pp. 48–55.

[34] P. Xie, K. Zuo, Y. Zhang, F. Li, M. Yin, K. Lu, Interpretable classification from skin cancer histology slides using deep learning: A retrospective multicenter study, arXiv preprint arXiv:1904.06156.

[35] A. Mikołajczyk, M. Grochowski, A. Kwasigroch, Global explanations for discovering bias in data, arXiv preprint arXiv:2005.02269.

[36] D. Sonntag, F. Nunnari, H.-J. Profitlich, The skincare project, an interactive deep learning system for differential diagnosis of malignant skin lesions. technical report, arXiv preprint arXiv:2005.09448.

[37] M. Hägele, P. Seegerer, S. Lapuschkin, M. Bockmayr, W. Samek, F. Klauschen, K.-R. Müller, A. Binder, Resolving challenges in deep learning-based analyses of histopathological images using explanation methods, Scientific reports 10 (1) (2020) 1–12.

[38] Y. Yan, J. Kawahara, G. Hamarneh, Melanoma recognition via visual attention, in: International Conference on Information Processing in Medical Imaging, Springer, 2019, pp. 793–804.

[39] Z. Zhang, Y. Xie, F. Xing, M. McGough, L. Yang, Mdnet: A semantically and visually interpretable medical image diagnosis network, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 6428–6436.

[40] J. Kawahara, S. Daneshvar, G. Argenziano, G. Hamarneh, Seven-point checklist and skin lesion classification using multitask multimodal neural nets, IEEE journal of biomedical and health informatics 23 (2) (2018) 538–546.

[41] J. Kawahara, G. Hamarneh, Fully convolutional neural networks to detect clinical dermoscopic features, IEEE journal of biomedical and health informatics 23 (2) (2018) 578–585.

[42] E. Veltmeijer, S. Karaoglu, T. Gevers, et al., Integrating clinically-relevant features into skin lesion classification., in: BNAIC/BENELEARN, 2019.

[43] S. Murabayashi, H. Iyatomi, Towards explainable melanoma diagnosis: Prediction of clinical indicators using semi-supervised and multi-task learn-

ing, in: 2019 IEEE International Conference on Big Data (Big Data), IEEE, 2019, pp. 4853–4857.

[44] A. Lucieri, M. N. Bajwa, S. A. Braun, M. I. Malik, A. Dengel, S. Ahmed, On interpretability of deep learning based skin lesion classifiers using concept activation vectors, arXiv preprint arXiv:2005.02000.

[45] D. Coppola, H. Kuan Lee, C. Guan, Interpreting mechanisms of prediction for skin cancer diagnosis using multi-task learning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020, pp. 734–735.

[46] Z. Chen, Y. Bei, C. Rudin, Concept whitening for interpretable image recognition, Nature Machine Intelligence 2 (12) (2020) 772–782.

[47] R. Zhang, S. Tan, R. Wang, S. Manivannan, J. Chen, H. Lin, W.-S. Zheng, Biomarker localization by combining cnn classifier and generative adversarial network, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2019, pp. 209–217.

[48] M. D. Zeiler, R. Fergus, Visualizing and understanding convolutional networks, in: European conference on computer vision, Springer, 2014, pp. 818–833.

[49] K. Simonyan, A. Vedaldi, A. Zisserman, Deep inside convolutional networks: Visualising image classification models and saliency maps, arXiv preprint arXiv:1312.6034.

[50] A. Shrikumar, P. Greenside, A. Shcherbina, A. Kundaje, Not just a black box: Learning important features through propagating activation differences, arXiv preprint arXiv:1605.01713.

[51] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, A. Torralba, Learning deep features for discriminative localization, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 2921–2929.

[52] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra, Grad-cam: Visual explanations from deep networks via gradient-based localization, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 618–626.

[53] J. T. Springenberg, A. Dosovitskiy, T. Brox, M. Riedmiller, Striving for simplicity: The all convolutional net, arXiv preprint arXiv:1412.6806.

[54] V. Petsiuk, A. Das, K. Saenko, Rise: Randomized input sampling for explanation of black-box models, arXiv preprint arXiv:1806.07421.

[55] L. M. Zintgraf, T. S. Cohen, T. Adel, M. Welling, Visualizing deep neural network decisions: Prediction difference analysis, arXiv preprint arXiv:1702.04595.

[56] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, PloS one 10 (7) (2015) e0130140.

[57] M. T. Ribeiro, S. Singh, C. Guestrin, " why should i trust you?" explaining the predictions of any classifier, in: Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining, 2016, pp. 1135–1144.

[58] J. Ba, V. Mnih, K. Kavukcuoglu, Multiple object recognition with visual attention, arXiv preprint arXiv:1412.7755.

[59] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, Y. Bengio, Show, attend and tell: Neural image caption generation with visual attention, in: International conference on machine learning, 2015, pp. 2048–2057.

[60] J. Fu, H. Zheng, T. Mei, Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.

[61] R. Caruana, Multitask learning, Machine learning 28 (1) (1997) 41–75.

[62] B. Kim, M. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viegas, et al., Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav), in: International conference on machine learning, 2018, pp. 2668–2677.

[63] J. Wan, D. Wang, S. C. H. Hoi, P. Wu, J. Zhu, Y. Zhang, J. Li, Deep learning for content-based image retrieval: A comprehensive study, in: Proceedings of the 22nd ACM international conference on Multimedia, 2014, pp. 157–166.

[64] A. Qayyum, S. M. Anwar, M. Awais, M. Majid, Medical image retrieval using deep convolutional neural network, Neurocomputing 266 (2017) 8–20.

[65] S. Lapuschkin, S. Wäldchen, A. Binder, G. Montavon, W. Samek, K.-R. Müller, Unmasking clever hans predictors and assessing what machines really learn, Nature communications 10 (1) (2019) 1–8.

[66] D. B. Larson, H. Harvey, D. L. Rubin, N. Irani, R. T. Justin, C. P. Langlotz, Regulatory frameworks for development and evaluation of artificial intelligence–based diagnostic imaging algorithms: Summary and recommendations, Journal of the American College of Radiology.

[67] J. Zhang, Y. Xie, Y. Xia, C. Shen, Attention residual learning for skin lesion classification, IEEE transactions on medical imaging 38 (9) (2019) 2092–2103.

[68] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, Rethinking the inception architecture for computer vision, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 2818–2826.

[69] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., Imagenet large scale vi-

sual recognition challenge, International journal of computer vision 115 (3) (2015) 211–252.

[70] T.-Y. Lin, A. RoyChowdhury, S. Maji, Bilinear cnn models for fine-grained visual recognition, in: Proceedings of the IEEE international conference on computer vision, 2015, pp. 1449–1457.

[71] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.

[72] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.

[73] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: Advances in neural information processing systems, 2017, pp. 4765–4774.

[74] R. Gu, G. Wang, T. Song, R. Huang, M. Aertsen, J. Deprest, S. Ourselin, T. Vercauteren, S. Zhang, Ca-net: Comprehensive attention convolutional neural networks for explainable medical image segmentation, IEEE Transactions on Medical Imaging.

[75] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical image computing and computer-assisted intervention, Springer, 2015, pp. 234–241.

[76] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: Advances in neural information processing systems, 2014, pp. 2672–2680.

[77] Y. Yamamoto, T. Tsuzuki, J. Akatsuka, M. Ueki, H. Morikawa, Y. Numata, T. Takahara, T. Tsuyuki, K. Tsutsumi, R. Nakazawa, et al., Automated acquisition of explainable knowledge from unannotated histopathology images, Nature communications 10 (1) (2019) 1–9.

[78] P. Sabol, P. Sinčák, P. Hartono, P. Kočan, Z. Benetinová, A. Blichárová, L. Verbóová, E. Štammová, A. Sabolová-Fabianová, A. Jašková, Explainable classifier for improving the accountability in decision-making for colorectal cancer diagnosis from histopathological images, Journal of Biomedical Informatics 109 (2020) 103523.

[79] A. B. Tosun, F. Pullara, M. J. Becich, D. L. Taylor, S. C. Chennubhotla, J. L. Fine, Histomapr™: An explainable ai (xai) platform for computational pathology solutions, in: Artificial Intelligence and Machine Learning for Digital Pathology, Springer, 2020, pp. 204–227.

[80] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, B. Kim, Sanity checks for saliency maps, in: Advances in Neural Information Processing Systems, 2018, pp. 9505–9515.

[81] B. Mittelstadt, C. Russell, S. Wachter, Explaining explanations in ai, in: Proceedings of the conference on fairness, accountability, and transparency, 2019, pp. 279–288.

[82] A. Lucieri, M. N. Bajwa, A. Dengel, S. Ahmed, Explaining ai-based decision support systems using concept localization maps, in: International Conference on Neural Information Processing, Springer, 2020, pp. 185–193.

[83] G. Montavon, W. Samek, K.-R. Müller, Methods for interpreting and understanding deep neural networks, Digital Signal Processing 73 (2018) 1–15.

[84] D. de Mijolla, C. Frye, M. Kunesch, J. Mansir, I. Feige, Human-interpretable model explainability on high-dimensional data, arXiv preprint arXiv:2010.07384.

[85] K. Abhishek, G. Hamarneh, M. S. Drew, Illumination-based transformations improve skin lesion segmentation in dermoscopic images, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2020, pp. 728–729.

[86] A. Ghorbani, J. Wexler, J. Y. Zou, B. Kim, Towards automatic concept-based explanations, in: Advances in Neural Information Processing Systems, 2019, pp. 9277–9286.

[87] J. Hu, R. Ji, Q. Ye, T. Tong, S. Zhang, K. Li, F. Huang, L. Shao, Architecture disentanglement for deep neural networks, arXiv preprint arXiv:2003.13268.

[88] P. Esser, R. Rombach, B. Ommer, A disentangling invertible interpretation network for explaining latent representations, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 9223–9232.