# ACHIEVEMENTS AND CHALLENGES IN EXPLAINING DEEP LEARNING BASED COMPUTER-AIDED DIAGNOSIS SYSTEMS

**Adriano Lucieri**[*]
Department of Computer Science
Technical University of Kaiserslautern
67663 Kaiserslautern, Germany
Smart Data and Knowledge Services
German Research Center for AI GmbH (DFKI)
67663 Kaiserslautern, Germany
adriano.lucieri@dfki.de

**Muhammad Naseer Bajwa**[*]
Department of Computer Science
Technical University of Kaiserslautern
67663 Kaiserslautern, Germany
Smart Data and Knowledge Services
German Research Center for AI GmbH (DFKI)
67663 Kaiserslautern, Germany
naseer.bajwa@dfki.de

**Andreas Dengel**
Department of Computer Science
Technical University of Kaiserslautern
67663 Kaiserslautern, Germany
Smart Data and Knowledge Services
German Research Center for AI GmbH (DFKI)
67663 Kaiserslautern, Germany
andreas.dengel@dfki.de

**Sheraz Ahmed**
Smart Data and Knowledge Services
German Research Center for AI GmbH (DFKI)
67663 Kaiserslautern, Germany
sheraz.ahmed@dfki.de

November 30, 2020

## ABSTRACT

Remarkable success of modern image-based AI methods and the resulting interest in their applications in critical decision-making processes has led to a surge in efforts to make such intelligent systems transparent and explainable. The need for explainable AI does not stem only from ethical and moral grounds but also from stricter legislation around the world mandating clear and justifiable explanations of any decision taken or assisted by AI. Especially in the medical context where Computer-Aided Diagnosis can have a direct influence on the treatment and well-being of patients, transparency is of utmost importance for safe transition from lab research to real world clinical practice. This paper provides a comprehensive overview of current state-of-the-art in explaining and interpreting Deep Learning based algorithms in applications of medical research and diagnosis of diseases. We discuss early achievements in development of explainable AI for validation of known disease criteria, exploration of new potential biomarkers, as well as methods for the subsequent correction of AI models. Various explanation methods like visual, textual, post-hoc, ante-hoc, local and global have been thoroughly and critically analyzed. Subsequently, we also highlight some of the remaining challenges that stand in the way of practical applications of AI as a clinical decision support tool and provide recommendations for the direction of future research.

***Keywords*** Artificial Intelligence in Healthcare · Computer-Aided Diagnosis · Medical Image Analysis · Explainable Artificial Intelligence · Interpretability · Human-Centric Computing

---

[*]Authors contributed equally

[†]This is English translation of a German book chapter to appear in Springer Nature 'KI Im Gesundheitswesen'.

# 1 Introduction

Artificial Intelligence (AI) based methods to support medical professionals were developed as early as 1970s [1, 2]. However, hardly any of these methods ever found practical applications in practical clinical environment, which was partly due to the technological status of computer systems at that time and emerging ethical and legal concerns. The exponential growth in computing power and digital data volume in the past decades, coupled with maturing of data-driven algorithms like Deep Neural Networks (DNNs), led to a number of groundbreaking successes in areas such as image recognition [3] and speech recognition [4]. From the outset of the success story of image-based Deep Learning (DL) models, like Convolutional Neural Networks (CNNs), an interest in their application in medicine developed immediately. The supremacy of modern DL-based solutions in the field of automated diagnosis of medical images over conventional image processing methods quickly became apparent [5, 6]. Shortly afterwards, early comparative studies reported the superiority of such self-learning algorithms over human experts in fields such as dermatology [7, 8], ophthalmology [9, 10] and radiology [11]. In spite of enormous technical and political progress, most ethical concerns about the fairness and transparency of such algorithms still remain unresolved. Unlike earlier rule-based or linear AI models, increasingly complex DL methods often act as black-boxes which do not offer semantically traceable decision processes. Moreover, legal requirements such as European General Data Protection Regulation (GDPR) [12], which came into effect in 2018, stress compulsory obligations regarding the transparency of automated decision-making processes towards affected individuals.

These factors, among others, contributed to current trend in the field of AI-based diagnosis to move towards Computer-Aided Diagnosis (CAD) and so called "Augmented Doctor" [13]. Advantages such as speed, objectiveness, and thoroughness of AI methods are being used, for example, as an assistance system to point out relevant disease indicators to doctors, give diagnosis suggestions and to present similar past cases for comparison. This approach of clinical application benefits from being a Human-In-The-Loop (HITL) hybrid keeping the clinical experts in control of the process [14]. This HITL model is similar to driving aids like adaptive cruise control or lane keep assistance in automobiles where human driver retains control and bears responsibility for the final decisions but with a reduced workload and an added safety net. In healthcare, this implies that the indicators and explanations presented by AI move to the centre of attention and that the ultimate diagnosis will be of less relevance as the doctor will always make the final decision. The ability of a diagnostic assistance system to explain its decision path has been considered as most important characteristic by health professionals since long [15]. This and more general demand for justifications of automated high-stakes decisions led to an enormous increase in research in the field of explainable AI (xAI) [16]. The claim of this article is not to provide an extensive list of previous works and efforts to explain DNNs in medical problems nor to define the notion of explainability and interpretability. Interested readers will find broad overviews of the applications of xAI methods in general medical problems in [17, 18] as well as special works on the topics of medical image analysis [14] and digital pathology [19]. Also, there is no collectively agreed upon xAI vocabulary, but there have been numerous efforts aimed at defining key terms and taxonomies [20, 21, 16]. In this work, we follow the definitions of explainability and interpretability provided in [22].

In this article we take an application-oriented look at the current achievements in the field of the xAI methods to distil remaining challenges in the way of its practical application. We show that in many cases the technical bases for transparent AI applications in medicine already exist and that the missing piece of the puzzle is often extensive cooperation of AI developers and medical domain experts. We also caution readers to not overconfidently rely on any xAI method but be aware of the difficulties in objective evaluation of fidelity and quality of explanations. It is emphasised that AI developers should pay special attention to application-specific characteristics such as human-centricity, domain-orientation, diversity and completeness of explanations when developing xAI methods in future. The following section provides a short overview of some existing xAI methods. In section 3 we present current achievements in application of explainable DL methods to medical problems. This includes explanation-based improvement of algorithms, validation of DNNs decision processes, discovery of previously unknown disease criteria, as well as proposed DL-based xAI frameworks for CAD. Section 4 sheds a critical light some of the most strenuous challenges that bar the way to progress. Finally, the findings are concluded in section 5.

# 2 Overview of Common xAI Methods

Methods explaining the decision-making process of DNNs exist in a variety of forms. Not only the derivation of the explanations differs, but also the way it is communicated to the user. There are a number of taxonomies available in literature to differentiate these methods. An important distinction for AI users, for example, is made between post-hoc and ante-hoc methods. Methods which can explain the decision of a black-box model afterwards are called post-hoc (lat.: after-this event) methods. Ante-hoc (lat.: before-this event) methods, on the other hand, are already interpretable due to their architecture. Since this ante-hoc explanations are usually achieved by architectural or conceptual restrictions
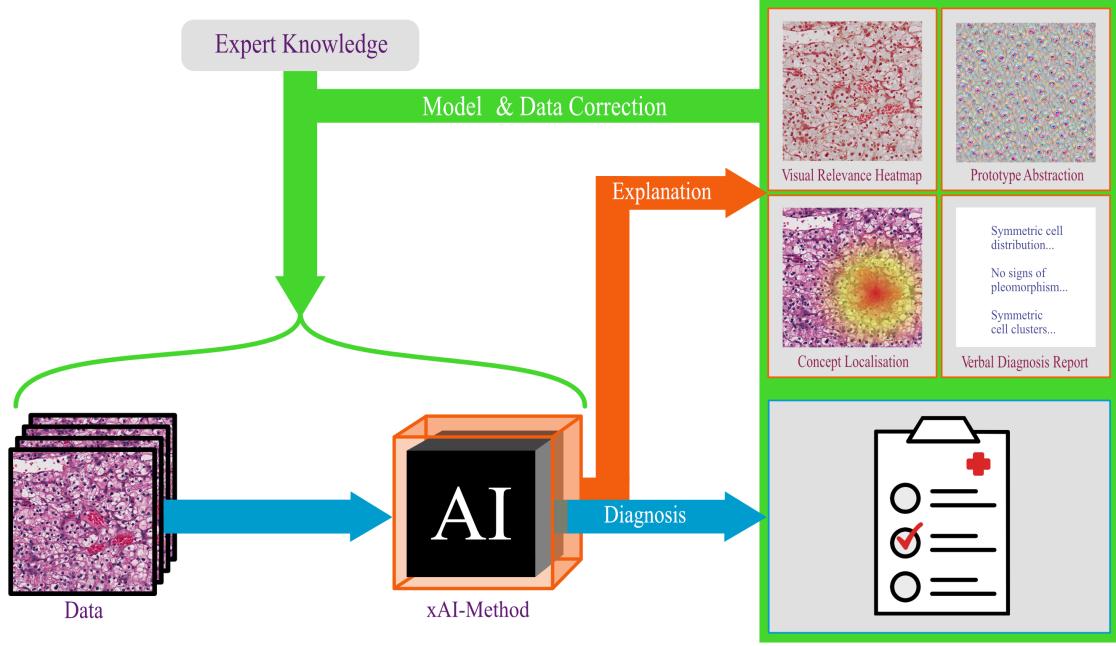
Figure 1: Topology of the xAI process with optional model and data correction as well as taxonomy of common and relevant xAI methods in medical image analysis.

in the learning process that limits modelling capacity, such inherently interpretable models are often thought to be inferior to their unrestricted conspecifics in terms of final model performance. However, this effect can sometimes be mitigated by pre-processing raw data with noisy features into meaningfully structured representations [23]. Another important distinction is made between explanation methods that attempt to explain a classifier's decision process on a global scale and those that focus on explaining single data sample at a time. Although local explanations might be initially sufficient for clinical applications as assistive diagnosis systems, global-level explanations will be crucial for understanding the model behaviour as a whole. This is specifically important for identifying decision biases and hence for the development towards autonomous decision systems. There exist various taxonomies for xAI methods in literature. We present a few types of methods that are specifically relevant to medical imaging. A visual overview of the grouping is provided in Figure 1.

## 2.1 Visual Relevance Heatmaps

Probably the most popular group of methods for explaining and interpreting image-based classification methods is the generation of visual heatmaps representing the influence of individual pixels on the result of the classification. Existing methods differ significantly in the computation of relevance values. The most obvious approach is the visualisation of the internal activations of a model [24]. Therefore, single or combinations of intermediate, two-dimensional activation values are scaled to input size and visualised. Other common methods rely on attribution of the classification results to the individual pixels. In practice, this is done using e.g. weighted activations in Class Activation Mapping (CAM) [25], gradient-based methods like Saliency [26], Gradient*Input [27], Grad-CAM [28], Integrated Gradient [29], DeepLift [30] or methods based on mathematical decomposition like Layerwise-Relevance Propagation (LRP) [31], Agglomerative Contextual Decomposition (ACD) [32] and SHapley Additive exPlanations (SHAP) [33]. All these methods require access to the model parameters and thus an understanding of the model architecture.

Perturbation-based methods, on the other hand, are completely model-agnostic and can thus be used for model-independent explanation without knowledge of their internals. In order to explain a given sample, it is modified several times and evaluated by the model again and again in order to systematically record the changes caused by the perturbations. Methods like Occlusion [24], RISE [34], and Extremal Perturbation [35] differ in the occlusion strategy (procedure and perturbation). LIME [36] goes one step further and trains local approximation models based on the results of the randomly modified images.

In addition to the post-hoc methods mentioned so far, there is also possibility to generate relevance heatmaps in an ante-hoc process. Here, model architectures can be extended by attention mechanisms that force the model to focus its

3

attention explicitly on certain parts of the input and to hide the remaining part. This distribution of attention can often be visualised in a heatmap [37], using pointers [38] or by explicitly cropping the input to the attended region [39] to gain insight into the network's decision-making process.

## 2.2    Class- and Prototype Abstraction

Visual Relevance Heatmaps (VRHs) usually help to explain the decision on individual samples. Another approach that aims towards both global and local explanations of DL models is the generalised representation of prototypes of individual classes or neurons as learned by the model. This includes, for instance, methods maximising the activation of particular outputs [26] or intermediate neurons [40] by optimising over an input image to determine their "prototypical" activation patterns. Many variations of this approach have already yielded interesting results and insights [41] for general image recognition tasks. However, only few works can be found applying abstraction methods to medical problems [42, 43, 44, 45]. This might be due to the complexity and entanglement of disease criteria and consequently complications in interpreting the prototypical results.

## 2.3    Conceptual Explainability and Biomarker Identification

The aim of concept-based explanation methods is to map human understandable semantic concepts to the concepts learned by DL models after training in order to make their decision-making processes more comprehensible. Such concepts can be very simple characteristics such as colours, shapes or textures. However, complex concepts can also be defined, consisting of combinations of simpler concepts. The TCAV method developed by Kim et al. [46] requires only a small number of sample images per concept to compute global concept influence scores. Further exploitation of this method allows explicit localisation of the concepts recognised by the network in the input domain [47], extend its application to regression tasks [42, 43] and introduce improved metrics [48]. Other concept-based approaches include the ones proposed by Bau et al. [49] and Zhou et al. [50].

Especially in the application of DL in medical problems, the detection and localisation of biomarkers by the model is popular in addition to the diagnosis of diseases. This approach allows intermediate steps of the models to be validated by experts. As has been shown in recent works [42, 43, 51], even post-hoc concept-based methods can be used to detect such biomarkers. However, more common approaches in literature are ante-hoc methods based on multi-task learning [52], where the models are trained for the combined classification or localisation of biomarkers [53, 54, 55]. Segmentation networks are often used for localization as in [54], however, such explicit approaches presuppose that correspondingly annotated data are available. An alternative approach by Zhang et al. [56] combines the optimisation of a CNN and a Generative Adversarial Network (GAN)[3] in a single end-to-end architecture for the localisation of biomarkers without the presence of explicit biomarker annotations.

## 2.4    Textual Explainability

There are different methods for generating verbal explanations of DL model decisions. One can distinguish between methods that use a template approach [57, 58, 59], rule-based methods [60, 61, 62, 63] and methods that use Natural Language Processing (NLP) models to generate an explanatory text [64, 65]. An early use case of NLP-based, textual explanation generation in the medical domain is MDNet framework developed by Zhang et al. [65]. This framework allows to generate a textual diagnostic report based on a medical image. In addition, a heatmap is generated for each word of the diagnostic report, which shows users the model's attention at that step.

# 3    Achievements of xAI in Medicine

The number of research papers on interpretability and explainability of AI has mushroomed in the last few years [16] and thereby the application and adaption of xAI methods to specific medical domains have also increased. In the following, we present research with the most practical significance towards clinical decision support systems.

## 3.1    Interventional Methods

The explanation of high performing AI algorithms that utilise spurious indicators for classification allows to reveal those biases. To make practical use of these explanations, methods that facilitate intervention and correction of current working of algorithms are required. Common methods for penalisation and correction of explanations in DL models

---

[3]Generative DNN that can be trained to learn underlying data generating process of a given training dataset. It can be used to interpolate between samples, generating unseen images.
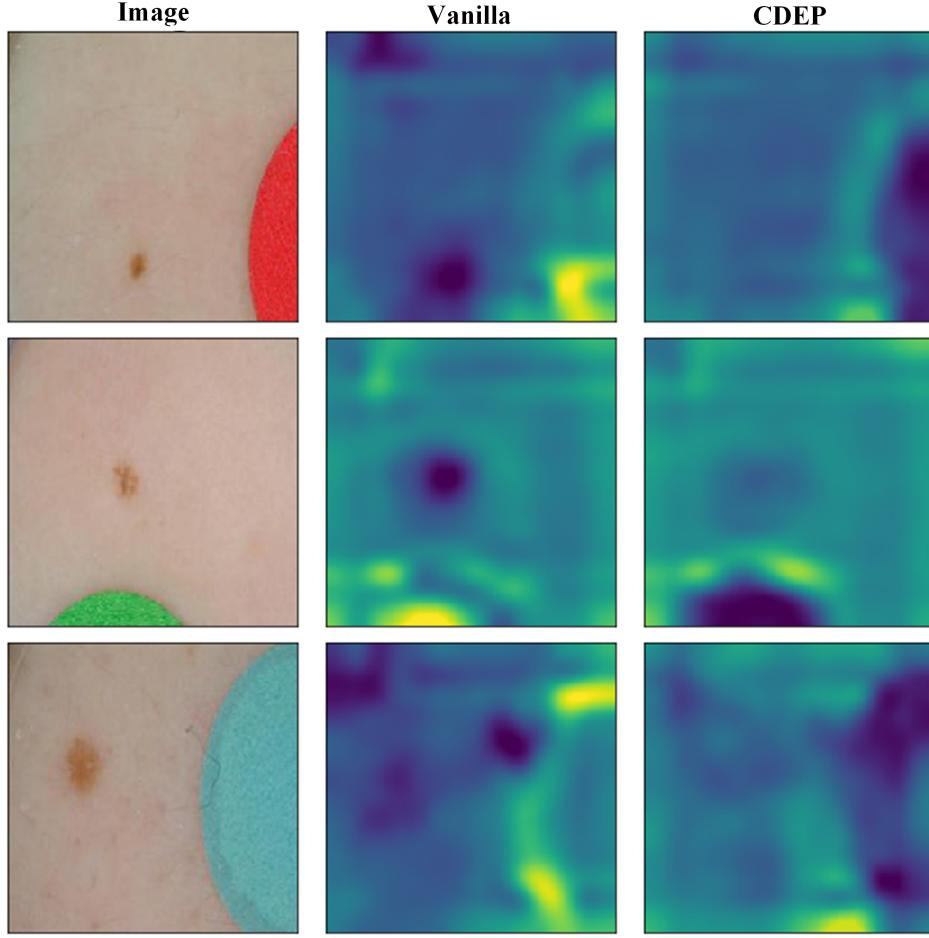
Figure 2: Results from Rieger et al. [71]. Left column shows original image samples from the dataset. The middle and right columns show grad-CAM heatmaps before and after correction using the developed model-correction method.

work by imposing a loss on explanation heatmaps (e.g. from VRH method) or conceptual predictions (e.g. TCAV) against ground truth explanations provided by human experts [66, 67]. This area is strongly related to the field of explicit expert knowledge incorporation. Examples of successful application of such methods in the medical domain are disease grading in diabetic retinopathy [68], lymph node histopathology [69] and dermoscopic skin lesion classification [70, 71]. Rieger et al. [71], for instance, were able to correct a classifier trained on the ISIC 2019 dataset, which is heavily biased towards benign predictions when coloured patches appear beside the lesion. A comparison between Grad-CAM maps generated before and after correction of the network can be seen in Figure 2. Inspired by the concept-based explanation method of TCAV, Graziani et al. [69] fine-tuned a deep classifier for histopathologic lymph node tumour detection. By penalising undesired control targets (concepts), they managed to increase Area Under the Curve (AUC) by 2%.

## 3.2 Revealing new criteria

Explainability methods are often employed in specific, sometimes medical application areas by expert computer scientist to prove their effectiveness. Lack of domain knowledge often hampers proper interpretation of presented results, rendering the provided explanations useless for assessing the correctness of the network. However, an increasing trend of collaboration between medical professionals and computer scientists is apparent in the application and tuning of DL models. The need for domain knowledge in order to understand and explain models has reflected in a growing number of publications on xAI including both computer scientists and domain experts.

A team of computer scientists and neurosurgeons succeeded in training a CNN for the localization of diagnostic features in confocal laser endomicroscopy images for glioma detection using only image-level annotations. Izadyyazdanabadi et al. [72] sequentially applied a visual relevance localization method (CAM) to a multi-head network, merging the

resulting maps by collateral integration as well as biologically inspired lateral inhibition principle. Their diagnostic localization maps correctly identified familiar diagnostic features and also revealed new diagnostic regions that were previously unknown to the neurosurgeons.

Using a complex model architecture consisting of two autoencoders and further processing steps, an interdisciplinary team of pathologists and computer scientists successfully predicted the recurrence of prostate cancer from digitised slides of histological sections in [73]. An especially developed method for calculating an impact score, which provides information about the direction of influence of an image section for diagnosis, provides further insights. It has been confirmed that the model independently learned the concept of the Gleason Score, an established prognostic value for prostate cancer among experts worldwide, and has identified the occurrence of stroma, an intermediate tissue running through the parenchymatous organs, in areas of the incision free of cancer cells as a prognostic factor for prostate cancer.

No clear physiological characteristics of insomnia are known yet. Researchers from Charité Berlin, HTW Berlin and University Medicine Göttingen have used machine learning models in [74] to detect insomnia in polysomnographic data with the aim of revealing such physiological features through AI. By applying DeepLift [30] method, some factors, such as increased and less synchronous eye movement, were highlighted as relevant for the prediction of insomnia. However, the authors themselves stress that the results should be interpreted with caution at first, as neither the bias of the results due to laboratory conditions can be excluded nor can the validity of the factors be definitively confirmed.

Lucieri et al. [51] used the concept-based TCAV method to investigate a high-performing CNN for the classification of skin lesions into malignant melanoma, benign nevi and seborrheic keratosis. Using only a few particularly detailed annotated images, the authors were able to show that the CNN uses proven concepts defined by the medical community to classify images. These results were also confirmed by a team of experienced dermatologists. Part of the results regarding seborrheic keratosis, however, are yet to be confirmed.

A team of computer scientists and biologists have used samples of microbiomes of human female skin to determine phenotypes such as age, skin moisture, menopause status and smoking status in [75]. The SHAP method was used to assess the relevance of each bacterial genus in the microbiome. As this method generates local explanations, SHAP values for all bacterial genera were averaged over the subset of samples with correct and good results for classification and regression. The most relevant bacterial genera and their influence on the respective task were reported. For the determination of all phenotypes, a number of relevant bacteria genera were identified. In the case of skin moisture determination, for example, the genera identified by the model as particularly important were already associated with skin moisture in previous studies.

Essemlali et al. [76] were able to determine whether patients suffer from mild cognitive impairment or even Alzheimer's dementia using their two-dimensional connectivity matrix of brain regions. They used a specially adapted CNN architecture for this purpose. To explain the disease prognosis, the gradients of all images of the respective classes were averaged to obtain a global explanation. These averaged heatmaps of different classes were subtracted to emphasize the crucial differences between two conditions. The results confirmed that the connectivity of the entorhinal cortex is crucial for the separation between healthy and Alzheimer's disease subjects and the hippocampus for the separation between healthy subjects and those with mild cognitive impairment. Their results have been discussed with an expert neuroanatomist.

### 3.3 xAI Frameworks

Since its rise in the early 2010s, DL and our understanding of its workings have greatly evolved. Thus, it is not surprising that the first commercial and non-commercial applications of explainable DL-based software solutions entered the market.

Data Language (UK) Ltd. [77] already use their commercial SCOPA Explainable AI Platform solution in classification of terrorist and extremist propaganda and claim it to be applicable to healthcare applications such as bone fracture detection from radiology images. The framework is described as using a layered ensemble of Deep Learning, Machine Learning and Algorithmic models for holistically explainable and scalable media classification.

The US-based start-up Decoded Health offers a commercial telehealth platform for radical automation of patient-doctor interactions and automated decision support for doctors [78]. Among other things, the company uses the Deep Adaptive Semantic Logic (DASL) [79] and ARSENAL [80] methodologies developed by its parent company SRI International to explain DL models using rule-based reasoning and providing natural language interaction.

Explainable AI for Dermatology (exAID) is a showcase framework for computer-aided diagnosis of diseases developed by German Research Center for Artificial Intelligence GmbH (DFKI) [81]. exAID is based on concept-based explanation methods [47, 51] to detect and localise biomarkers and generate verbal diagnostic explanations. In addition, the tool

offers an advanced mode that allows data scientists and physicians in residence training to explore collected data sets in an exploratory way, including the suggestions of the DL model. The tool can be applied to any DL model and is thus compatible with all state-of-the-art architectures.

Another DFKI project is specifically focused on the development of a Computer-Aided Diagnosis (CAD) system for the detection of skin diseases [82, 83]. The system developed in the Skincare project is capable of analysing images of skin diseases taken with a smartphone, generating a differential diagnosis, and segmenting the skin lesion and individual biomarkers. The explainability of the system is ensured through the calculation of expert scores and VRHs. A demo of the system can be tested on the project webpage[4].

# 4 Challenges for xAI Application

Since the initial applications of modern DL-based systems in medical domains, we have seen remarkable strides in the explanation of systems that in some cases already led to correction and verification of AI as well as disclosure of new potential diagnostic criteria. However, there are still a number of challenges pertinent to medical image diagnosis, which should be addressed by concerted efforts from AI researchers, medical practitioners and regulatory authorities.

## 4.1 Evaluation of Explanation Methods

Before xAI methods can be practically applied, it must be ensured that their explanations are reliable, trustworthy, and useful. In three steps, we distinguish between the evaluation of an explanation's truthfulness, usefulness and finally the interpretation by the user.

### 4.1.1 Evaluation of Truthfulness

One of the key challenges in explainable AI is difficulty in evaluating if the explanation of a model's behaviour is reliable, primarily because there is no ground truth available for evaluation [84]. Truthfulness or fidelity of an explanation refers to whether it is reliable and reflects the actual decision process of the AI. In order to practically deploy AI in clinical environments such that it increases the efficiency and accuracy of human doctors, it is of paramount importance to ensure the fidelity of xAI methods. However, due to lack of explanation ground truth, evaluation of such methods is largely subjective.

There have been attempts to quantify and measure the quality of explanations. Samek et al. [85] introduced Area Over the MoRF Perturbation Curve (AOPC) measure to quantitatively compare VRHs. The measure gradually perturbs input images starting from the regions that are marked as the most relevant according to a given explanation method. High AOPC values indicate that a model is sensitive to perturbations in those regions, thus confirming the validity. The RemOve And Retrain (ROAR) framework [84] is an advancement of AOPC approach. As image perturbations lead to a change in image distribution, they retrain the network on the perturbed images to avoid distribution gaps and evaluate the achieved accuracy. However, the evaluation of an altered model cannot give reliable insights into the sensitivity of the original model. In [17] a synthetic dataset with ground truth explanation has been generated for easier xAI method evaluation. Adebayo et al. [86] introduced randomisation tests in which model weights and data labels where systematically randomised to reveal if explanation methods where really model and data dependent. Although this method has not been used to quantify fidelity, its results are certainly meaningful for evaluation.

Truthfulness is the basis for robust and useful xAI. Results from works like [86] showed that some methods produce convincing explanations that are worth no more than simple edge detectors. Eitel et al. [87] performed a quantitative comparison of visual relevance methods for MRI-based Alzheimer's disease classification. They found that guided backpropagation attribution maps [88] averaged over all true positives for multiple training runs highlighted different regions in brain MRI. However, despite the variance, which makes it harder to compare and to replicate outcomes of individual experiments, some regions like hippocampus, cerebellum and edges of the brain were commonly identified as salient regions. Other visual relevance methods like Gradient*Input, Occlusion Sensitivity, and LRP also showed similar behaviour, which raises serious questions on the robustness and coherence of these explanation methods. However, this could also indicate abundance of biomarkers in the data that allows DNN's to perform the same task in a variety of ways.

### 4.1.2 Evaluation of Usefulness

Besides evaluating the fidelity and completeness of explanation methods, it is also crucial to quantify and qualify the usefulness of generated explanations. Doshi-Velez and Kim [89] proposed the distinction between application-grounded,

---

[4]http://www.dfki.de/skincare/classify.html

human-grounded and functional-grounded evaluation of explanations. In [90] the first functionally-grounded metrics where introduced, allowing to objectively judge the quality of an explanation. This quantification has the advantage of being independent from human subjectivity. On the other hand, human-grounded evaluation makes use of non-specialist human evaluators to subjectively compare or rate explanations. The evaluation approach that we find most important for xAI in medicine is the application-grounded evaluation. Depending on the domain or problem, medical practitioners have a very specific way of thinking of a problem, communicating or explaining a diagnosis. Hence, we argue that application-grounded evaluation is necessary to find and optimise the right explanation methods for a medical use-case.

### 4.1.3 Evaluation with respect to Evaluators

An equally decisive factor in the use of xAI methods is their interpretation by the end user. One and the same explanation can be interpreted differently by different individuals. A wrong or too naive interpretation of decision processes by developers or users can lead to serious consequences in the practical use of AI. The approach to the interpretation of explanations differs significantly for AI researchers and medical practitioners, but also overlaps to some extent.

For AI developers, explainability methods can help them design better models by understanding the interactions between the model and the data. However, AI developers and data scientists can sometimes over-trust or misuse these interpretability tools as noted by [91]. They conducted a small-scale study to learn how data scientists utilise publicly available interpretation tools and found that visual explanations are usually taken at their face values and used for rationalisation of suspicious observations instead of understanding how AI models worked. Experienced data scientists, on the other hand, were able to capitalise on these interpretability tools and effectively understand issues with models and data.

For medical practitioners, such tools can provide reasoning for model predictions and, therefore, develop trust and ease their acceptance into routine clinical workflow. Sayres et al. [92] evaluated the impact of DL-based diabetic retinopathy (DR) detection algorithm on the performance of human graders in computer-assisted setting. They found that the accuracy of human graders improved when assisted by the algorithm that provided only disease prediction without any explanation. However, when the graders were provided prediction plus visual explanation by the algorithm, their detection accuracy improved only for patients who had diabetic retinopathy (resulting in high sensitivity) and decreased for patients without DR (resulting in low specificity). Although the qualitative feedback of human graders on the explanations provided by the algorithm was generally positive, the participants were not able to harness this additional information to notably improve their performance. This could partly be because the pathologic features of DR are very tiny in size, inconspicuous and occupy only a fraction of the whole images space.

To meet the challenges in the evaluation of xAI, special focus should be placed on the evaluation of realistic applicability of methods in clinical environment. This includes truthfulness, robustness, quality and the actual usefulness of the methods. Through such detailed analyses, the agreement between medical expert knowledge and the knowledge gained from the model and data can be validated and evaluated and, possibly, new knowledge can be gained. A further dimension that should not be neglected when evaluating xAI applications in healthcare is the ethical assessment of the impact on individuals and society. As shown in section 3.3, there is an increasing commercial interest in explaining AI decisions. This requires development of regulatory measures that both take into account different needs of different individuals and user groups and are adaptable to the constantly evolving AI technology [93]. However, this also requires clearly defined evaluation and certification processes to assess the ethical conformity of the use of AI in a specific context. z-Inspection [94] is one of the first ethical evaluation and certification processes that integrates theoretical principles for the ethical evaluation of AI into a practically applicable framework.

### 4.2 Deployment in Clinical Workflow

Proof of concept studies and prototype methods are required to be tested rigorously to analyse their contextual fit in real-world clinical environment. However, many obstacles have been discovered and highlighted by researchers in implementing laboratory research in clinical settings. These challenges include lack of utility to clinicians' logistical hurdles that hamper clinical deployment and trials [95]. Ineffective use, or misuse, of these assistive systems can even lead to performance degradation of human graders [96, 97, 98]. Cai et al. [98] developed interactive user-centric techniques for pathologists to improve diagnostic utility and trust in algorithmic predictions, in lab settings. Previously, such Human-Computer Interface (HCI) techniques have been used only to improve the algorithm, however, these interactive tools have potential to enable users to test, understand and grapple with AI algorithms, leading to new ways for improving explainability of algorithms. Instead of waiting for algorithms to generate human understandable explanations [99, 36], interactive techniques can allow users to play an active role in the interpretation of algorithm predictions and hypothesis-test their intuitions. While these studies help understand the needs of clinicians as they interact with AI algorithms, they do not account for the highly situated nature of activities in clinical environments. In a study [100] designed for field assessment of a Decision Support System (DSS) for cardiologists, it was found that the

clinicians were more likely to embrace and use such systems if it was seamlessly and unobtrusively integrated into their existing workflow. However, the misuse of these systems can sometimes let the clinicians develop their own tolerance and workarounds in order to trust the algorithm results [101].

There are a few examples of such translation of AI into commercial applications for instance in detection of diabetic retinopathy [102], cancer, and analysing radiology images [103]. Deployment of CAD solutions in clinical settings can also help focus on the effects of workflow when new diagnostic and information systems are introduced into clinical environments. Arbabshirani et al. [104] integrated their AI based model for identification of Intracranial Haemorrhage (ICH) using head CT scans into clinical workflow for three months. During the trials, the model was able to reduce median time to diagnosis for routine studies from more than eight hours to only 19 minutes, while at the same time discovering some probable ICH cases which were overlooked by radiologists.

### 4.3 Diverse and Complete Explanations

Most applications of xAI in research focus on utilising single approaches and modalities for the explanation of AI models in given use-case. This can be seen in our analysis of achievements of xAI in section 3 as well as many reviews on xAI research [17, 18, 14, 19]. We believe that the integration of xAI in clinical workflows requires combination of multiple explanatory views to draw explanations that are diverse and as complete as possible. This is inspired by medical practitioners in routine clinical environments using textual descriptions alongside visualisations and temporal coherence to communicate decisions effectively and reliably. On one hand, this should motivate AI researchers to think of new, creative paths for xAI methods to complement existing methods and on the other to not only evaluate the effectiveness of approaches in isolation but in combination with diverse methods to leverage synergies. First efforts towards diverse explanations have been recently made in the Visual Question Answering community in works like [105] and [106]. Huk Park et al. [105] show the positive complementary effect of visual relevance and textual explanations which is backed up by human evaluation. Completeness of explanations can be considered from a model's and a user's point of view. Completeness from a model's point of view is directly related to fidelity. Yeh et al. [107] introduced a completeness measure that quantifies the completeness of a given concept-based explanation for a model prediction. Completeness from a user's point of view is subjective but equally relevant to usefulness.

### 4.4 Human-Centric Explanations

High-performing DL-based DNNs often utilise unintelligible notions of concepts to reach a prediction. Integration of AI assistants in clinical workflows requires human-centric explanation of decision that is able to not only explain a decision with high fidelity, but also conforms to human-understandable thought models. Compared to simpler use-cases like visual object classification or part segmentation, complex medical concepts used for diagnosis particularly necessitates to make explanations as human-understandable as possible.

#### 4.4.1 Human-Understandable Concepts

One way to explain the decisions of AI based CAD systems in a human-centric way is to investigate the role of human-understandable concepts, learned by DL-based algorithms. It is very important to analyse the learned features of an algorithm that makes right decisions but based on wrong reasons. It is a major issue that can affect performance when the system is deployed in the real world. Explaining the role of a model's concepts can reduce reliability concerns of medical practitioners and help develop their trust on CAD.

Application of concept-based xAI methods in medical image analysis has been problematic partly because these methods require concept datasets [50] or image patches corresponding to human-understandable concepts [46], which are not always available. An unsupervised approach, extending CAV method, is developed by Ghorbani et al. [108] to cluster object datasets by performing segmentation of single objects and clustering their relevant activations into semantically meaningful groups. This approach cannot be directly applied to, for example, skin lesion classification where there is a substantial overlap between various concepts that can not be segmented into distinct spatial patches. Also, this method does not guarantee discovery of human-understandable concepts and requires thorough human evaluation effort.

Sometimes general explanation methods cannot be readily used for certain medical image tasks due to technical requirements or inappropriateness to the domain. Besides the continuous development of advanced xAI methods, it is crucial that developers pay attention to the domain specific needs of particular medical applications and their users. There have been many studies extending existing methods to better suit the challenges of the medical image analysis. For example, Yang et al. [109] proposed Expressive gradients (EG), an extension of commonly used Integrated Gradients [29] to cover the retinal lesions better while [42] extended CAVs from [46] for continuous concepts like eccentricity and contrast. Lucieri et al. [47] extended the method for localising and highlighting image regions significant for network's concept recognition in a medical inspired dataset. This could allow doctors to verify the

9

network's concept learning and suggest precise image regions for concepts. Such studies lead to the advancement of the xAI domain and provide specialisation to application domains without designing new methods from scratch.

### 4.4.2 Challenges in Textual Explanations

Most disease classification algorithms using medical images attempt to answer Multiple Choice Questions in which the algorithm is expected to select one disease from a list of all possible diseases. In this type of experimental setting there is a fair chance that a correct prediction given by AI-based CAD is nothing more than a fluke – though the probability of fluke decreases with increase in total number of classes. Therefore, such classification algorithms require explicit interpretations of network predictions to validate their results.

In many medical domains like radiology and histopathology, doctors routinely write textual reports clearly noting salient findings before providing their impression (diagnosis). The nature of this type of detailed diagnosis substantiated by textual descriptions of the image is self-explanatory – at least for the domain experts. AI-based CAD can be enabled to process this multi-modal data (image and text) and generate textual reports to mimic the behaviour of radiologists and histopathologists. Such systems embed explanations of their decisions inside their predictions. These natural language explanations, using domain-specific terminology and mimicking structure of communication provide an intuitive and effective way of explaining decision processes to practitioners. However, providing textual explanations in the form of clinically accurate medical reports for medical images has some differences compared to other application areas where NLP is used to describe an image.

Generating long coherent reports (more than few tens of words) is one of the major challenges in textual xAI. Language generation models usually start with a few coherent sentences and after that its performance tapers off generating completely random words that have no association with the previously generated words or phrases. This happens generally due to very long temporal dependency among words which Long Short-Term Memory (LSTM) [110] models have difficulty handling. One way to address this problem is to use transformer networks [111] as language model decoder. These transformer models are able to capture relationship between words in a longer sentence better than Recurrent Neural Network (RNN) based models. Input text reports are tokenised and passed to the transformer network that consists of a single decoder layer and generates a query vector for another transformer model that generates reports by combining this query with information obtained from image processing model. The size of the generated reports and vocabulary can also be limited to ensure that the text is coherent and clinically meaningful.

Most of the reports written by doctors are free text report, which means that they don't always follow any defined template. Reports written by two radiologists, for example, for a given X-ray image can be vastly different. There can be superfluous information that does not contribute directly to the final diagnosis. This makes it very difficult to compare AI-generated reports with human-generated reports especially when some of the reports depend on the previous chest X-ray of the patients and provide a continuous diagnosis based on previous examination. This problem can be addressed by removing those parts of the input reports which bear no influence on the diagnosis such as at what time the doctor saw the patient, or who was the doctor on call.

### 4.4.3 Incorporation of Context

Traditional AI algorithms overwhelmingly rely on one input modality, for example images in medical image analysis. However, medical practitioners routinely incorporate context, in the form of, for instance, patient's clinical history, age, and sex etc., in their decision-making process. Compared to raw image pixels, this contextual information is much easier to understand for practitioners. Incorporation of this metadata into AI algorithms is tricky since context is difficult to represent in a form which is appropriate for processing by AI algorithms [112]. Not leveraging this useful context in AI algorithms can not only restrict their performance but also make explanations challenging. Therefore, another direction of research to make AI algorithms more transparent and explainable is to use multimodel data like medical images and patients' records together in the decision-making process and attribute the model decisions to each of them [14]. This approach simulates the diagnostic workflow of a clinician where both images and physical parameters of a patient are used to make the decision. It can not only improve diagnostic performance of these algorithms but also explain the phenomena more comprehensively.

An interesting example in which context mattered and lack of its inclusion resulted in a technically valid yet misleading ML prognostic model was the use of mortality risk prediction to make decisions about whether to provide treatment on an inpatient or outpatient basis for more than 14000 patients with pneumonia [113]. In this study the algorithm counter-intuitively suggested that patients with pneumonia and asthma were at lower risk of death compared to patients with only pneumonia, an indication that surprised the researchers who eventually ruled it out. A closer analysis of the data revealed that, at the hospital hosting this study, patients with history of asthma who presented with pneumonia were usually admitted directly to intensive care units to prevent complications. This led to a pattern in the data that

reflected better outcomes for such patients compared to patients with pneumonia and without history of asthma with approximately 50% less mortality rate. This example not only emphasises the importance of representative training data for such algorithms, but also that a contextually complete description of the data is of crucial importance.

## 5 Conclusion

Medical Image Analysis has benefited a lot from recent advances in deep learning. However, beyond academic research and proof of concept studies, there has been a healthy scepticism about to what extent, if at all, AI should make or support medical decisions in real clinical workflows [114]. Although the sequential computation of DL-based models is traceable, they often lack explicit declarative representation of knowledge. Justifying the decisions taken by AI using explanations can help bring such academic research one step closer to practical deployment in healthcare sector. We showed that the collaboration of AI developers and medical professionals already led to interesting advances in medical AI, including practical AI evaluation and discovery of new potential diagnostic criteria. It also appears that there is a growing interest in commercialising xAI solutions and developing use-case specific frameworks. However, we also caution to carefully gauge those achievements and continue investing efforts in standardised evaluation and investigation of xAI methods in close cooperation with domain experts.

In the medical domain, it is imperative to explain the output of algorithms in a human-understandable language as to support and not distract experts. Holzinger et al. [115] believe that the only way forward towards explainable CAD is to combine knowledge-driven and data-driven approaches, which could harness interpretability of former method and high accuracy of latter. We believe that the transition towards multimodal, diverse, and complete explanations that combine human-understandable modalities such as text, human-understandable concepts and context will substantially support the way of xAI in clinical assistive settings. In medical diagnosis, explanations can be different for different users. For instance, a doctor might use different language, modality or depth of explanation depending upon whether he/she is explaining to a patient, a regulator, or a fellow doctor. Similarly, explainable AI for healthcare serves different purpose for medical practitioners and AI developers. It is inevitable that AI engineers design solutions that provide diverse explanations fitting the need of specific use-cases. It has been observed that communication gap between AI developers and its users can lead to misuse of technology [96, 97]) and eventual performance degradation [92].

Qualitatively evaluating an explanation with regards to its interpretability and completeness can be substantially subjective. Recently, there have been many efforts to quantify and qualify xAI methods and their explanations in objective and subjective ways. However, there are no agreed-upon and standardised evaluation procedures for explanation methods that can guarantee fidelity and rate quality. Development of standardised and objective evaluation criteria can greatly help benchmark upcoming explainable CAD systems and is thus an extremely important requirement for application in routine clinical environments. Moreover, appropriate regulatory measures should provide an ethical framework for the application of AI in healthcare, which can ensure safety and transparency through standardised evaluation and certification procedures.

Since DL-based models are data driven, they suffer from limitations and biases inherent in the data. For example, some data might suggest that a cohort who took a certain drug recovered quickly compared to those who did not. DL models can detect this correlation easily. However, if the causality between drug and recovery is missing from the data, the models cannot propose an acceptable explanation of their decision either [116]. Therefore, explainable AI and particularly CAD solutions will hugely benefit from a carefully curated dataset which incorporates the context and does not leave out any confounding variables. Such dataset curation can be achieved by concerted and close collaboration between medical practitioners and AI developers right from the onset. Moreover, the advancement of interventional methods for the correction of algorithms and incorporation of explicit expert knowledge could account for small retrospective adjustments during trial phases.

We believe that modern AI technology has the potential to revolutionise healthcare in innumerable ways and that xAI plays a crucial role in creating a solid foundation of understanding and improving its functionality. Current advancements show that a close collaboration of medical domain experts and computer scientists paired with persistent efforts of AI experts to advance and develop new methods will eventually lead to many practical applications which are just an anticipation of what will be possible in the future.

## References

[1] BG BUCHANAN. Heuristic dendral: a program for generating explanatory hypotheses in organic chemistry. *Machine Intelligence*, 4:209–254, 1969.

[2] Edward Hance Shortliffe. MYCIN: a rule-based computer program for advising physicians regarding antimicrobial therapy selection. Technical report, STANFORD UNIV CALIF DEPT OF COMPUTER SCIENCE, 1974.

[3] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.

[4] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal processing magazine*, 29(6):82–97, 2012.

[5] Angel Alfonso Cruz-Roa, John Edison Arevalo Ovalle, Anant Madabhushi, and Fabio Augusto González 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*, pages 403–410. Springer, 2013.

[6] Varun Gulshan, Lily Peng, Marc Coram, Martin C Stumpe, Derek Wu, Arunachalam Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros, et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama*, 316(22):2402–2410, 2016.

[7] Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *nature*, 542(7639):115–118, 2017.

[8] Titus J Brinker, Achim Hekler, Alexander H Enk, Joachim Klode, Axel Hauschild, Carola Berking, Bastian Schilling, Sebastian Haferkamp, Dirk Schadendorf, Tim 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*, 113:47–54, 2019.

[9] Michael D Abràmoff, Philip T Lavin, Michele Birch, Nilay Shah, and James 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):1–8, 2018.

[10] James M Brown, J Peter Campbell, Andrew Beers, Ken Chang, Kyra Donohue, Susan Ostmo, RV Paul Chan, Jennifer Dy, Deniz Erdogmus, Stratis Ioannidis, et al. Fully automated disease severity assessment and treatment monitoring in retinopathy of prematurity using deep learning. In *Medical Imaging 2018: Imaging Informatics for Healthcare, Research, and Applications*, volume 10579, page 105790Q. International Society for Optics and Photonics, 2018.

[11] Pranav Rajpurkar, Jeremy Irvin, Robyn L Ball, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy Ding, Aarti Bagul, Curtis P Langlotz, et al. Deep learning for chest radiograph diagnosis: A retrospective comparison of the chexnext algorithm to practicing radiologists. *PLoS medicine*, 15(11):e1002686, 2018.

[12] Council of the European Union. Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation). available at http://data.europa.eu/eli/reg/2016/679/2016-05-04, April 2016.

[13] The American Medical Association (AMA, Ed.). AMA passes first policy recommendations on augmented intelligence.

[14] Amitojdeep Singh, Sourya Sengupta, and Vasudevan Lakshminarayanan. Explainable deep learning models in medical image analysis. *arXiv preprint arXiv:2005.13799*, 2020.

[15] Randy L Teach and Edward H Shortliffe. An analysis of physician attitudes regarding computer-based clinical consultation systems. *Computers and Biomedical Research*, 14(6):542–558, 1981.

[16] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115, 2020.

[17] Erico Tjoa and Cuntai Guan. Quantifying explainability of saliency methods in deep neural networks. *arXiv preprint arXiv:2009.02899*, 2020.

[18] Gregor Stiglic, Primoz Kocbek, Nino Fijacko, Marinka Zitnik, Katrien Verbert, and Leona Cilar. Interpretability of machine learning based prediction models in healthcare. *arXiv preprint arXiv:2002.08596*, 2020.

[19] Giulia Vilone and Luca Longo. Explainable artificial intelligence: a systematic review. *arXiv preprint arXiv:2006.00093*, 2020.

[20] Leilani H Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. Explaining explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th International Conference on data science and advanced analytics (DSAA)*, pages 80–89. IEEE, 2018.

[21] Zachary C Lipton. The mythos of model interpretability. *Queue*, 16(3):31–57, 2018.

[22] Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 73:1–15, 2018.

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

[24] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.

[25] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016.

[26] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.

[27] Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*, 2016.

[28] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.

[29] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. *arXiv preprint arXiv:1703.01365*, 2017.

[30] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. *arXiv preprint arXiv:1704.02685*, 2017.

[31] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015.

[32] Chandan Singh, W James Murdoch, and Bin Yu. Hierarchical interpretations for neural network predictions. *arXiv preprint arXiv:1806.05337*, 2018.

[33] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Advances in neural information processing systems*, pages 4765–4774, 2017.

[34] Vitali Petsiuk, Abir Das, and Kate Saenko. RISE: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*, 2018.

[35] Ruth Fong, Mandela Patrick, and Andrea Vedaldi. Understanding deep networks via extremal perturbations and smooth masks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2950–2958, 2019.

[36] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.

[37] Jimmy Ba, Volodymyr Mnih, and Koray Kavukcuoglu. Multiple object recognition with visual attention. *arXiv preprint arXiv:1412.7755*, 2014.

[38] Saumya Jetley, Nicholas A Lord, Namhoon Lee, and Philip HS Torr. Learn to pay attention. *arXiv preprint arXiv:1804.02391*, 2018.

[39] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in neural information processing systems*, pages 2017–2025, 2015.

[40] Aravindh Mahendran and Andrea Vedaldi. Visualizing deep convolutional neural networks using natural pre-images. *International Journal of Computer Vision*, 120(3):233–255, 2016.

[41] Chris Olah, Arvind Satyanarayan, Ian Johnson, Shan Carter, Ludwig Schubert, Katherine Ye, and Alexander Mordvintsev. The building blocks of interpretability. *Distill*, 3(3):e10, 2018.

[42] Mara Graziani, Vincent Andrearczyk, and Henning Müller. Regression concept vectors for bidirectional explanations in histopathology. In *Understanding and Interpreting Machine Learning in Medical Image Computing Applications*, pages 124–132. Springer, 2018.

[43] Mara Graziani, James M Brown, Vincent Andrearczyk, Veysi Yildiz, J Peter Campbell, Deniz Erdogmus, Stratis Ioannidis, Michael F Chiang, Jayashree Kalpathy-Cramer, and Henning Müller. Improved interpretability for computer-aided severity assessment of retinopathy of prematurity. In *Medical Imaging 2019: Computer-Aided Diagnosis*, volume 10950, page 109501R. International Society for Optics and Photonics, 2019.

[44] Shuhei Toba, Yoshihide Mitani, Hiroyuki Ohashi, Hirofumi Sawada, Noriko Yodoya, Hidetoshi Hayakawa, Masahiro Hirayama, Ayano Futsuki, Naoki Yamamoto, Hisato Ito, et al. Quantitative analysis of chest x-ray using deep learning to predict pulmonary to systemic flow ratio in patients with congenital heart disease. *Circulation*, 140(Suppl_1):A14250–A14250, 2019.

[45] Vincent Couteaux, Olivier Nempont, Guillaume Pizaine, and Isabelle Bloch. Towards interpretability of segmentation networks by analyzing deepdreams. In *Interpretability of Machine Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision Support*, pages 56–63. Springer, 2019.

[46] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). In *International conference on machine learning*, pages 2668–2677. PMLR, 2018.

[47] Adriano Lucieri, Muhammad Naseer Bajwa, Andreas Dengel, and Sheraz Ahmed. Explaining ai-based decision support systems using concept localization maps. *arXiv preprint arXiv:2005.01399*, 2020.

[48] Hugo Yeche, Justin Harrison, and Tess Berthier. UBS: A Dimension-Agnostic Metric for Concept Vector Interpretability Applied to Radiomics. In *Interpretability of Machine Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision Support*, pages 12–20. Springer, 2019.

[49] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6541–6549, 2017.

[50] Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Interpretable basis decomposition for visual explanation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 119–134, 2018.

[51] Adriano Lucieri, Muhammad Naseer Bajwa, Stephan Alexander Braun, Muhammad Imran Malik, Andreas Dengel, and Sheraz Ahmed. On interpretability of deep learning based skin lesion classifiers using concept activation vectors. *arXiv preprint arXiv:2005.02000*, 2020.

[52] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.

[53] Jeremy Kawahara, Sara Daneshvar, Giuseppe Argenziano, and Ghassan Hamarneh. Seven-point checklist and skin lesion classification using multitask multimodal neural nets. *IEEE journal of biomedical and health informatics*, 23(2):538–546, 2018.

[54] Jeremy Kawahara and Ghassan Hamarneh. Fully convolutional neural networks to detect clinical dermoscopic features. *IEEE journal of biomedical and health informatics*, 23(2):578–585, 2018.

[55] Davide Coppola, Hwee Kuan Lee, and Cuntai 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*, pages 734–735, 2020.

[56] Rong Zhang, Shuhan Tan, Ruixuan Wang, Siyamalan Manivannan, Jingjing Chen, Haotian Lin, and Wei-Shi Zheng. Biomarker localization by combining cnn classifier and generative adversarial network. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 209–217. Springer, 2019.

[57] Lisa Anne Hendricks, Ronghang Hu, Trevor Darrell, and Zeynep Akata. Grounding visual explanations. In *European Conference on Computer Vision*, pages 269–286. Springer, 2018.

[58] Pei Guo, Connor Anderson, Kolten Pearson, and Ryan Farrell. Neural network interpretation via fine grained textual summarization. *arXiv preprint arXiv:1805.08969*, 2018.

[59] Mohsin Munir, Shoaib Ahmed Siddiqui, Ferdinand Küsters, Dominique Mercier, Andreas Dengel, and Sheraz Ahmed. TSXplain: Demystification of DNN Decisions for Time-Series using Natural Language and Statistical Features. In *International Conference on Artificial Neural Networks*, pages 426–439. Springer, 2019.

[60] Shashank Srivastava, Igor Labutov, and Tom Mitchell. Joint concept learning and semantic parsing from natural language explanations. In *Proceedings of the 2017 conference on empirical methods in natural language processing*, pages 1527–1536, 2017.

[61] Braden Hancock, Martin Bringmann, Paroma Varma, Percy Liang, Stephanie Wang, and Christopher Ré. Training classifiers with natural language explanations. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2018, page 1884. NIH Public Access, 2018.

[62] Himabindu Lakkaraju, Ece Kamar, Rich Caruana, and Jure Leskovec. Faithful and customizable explanations of black box models. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 131–138, 2019.

[63] Johannes Rabold, Hannah Deininger, Michael Siebers, and Ute Schmid. Enriching visual with verbal explanations for relational concepts–combining lime with aleph. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 180–192. Springer, 2019.

[64] Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. Generating visual explanations. In *European Conference on Computer Vision*, pages 3–19. Springer, 2016.

[65] Zizhao Zhang, Yuanpu Xie, Fuyong Xing, Mason McGough, and Lin Yang. MDNet: A semantically and visually interpretable medical image diagnosis network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6428–6436, 2017.

[66] Andrew Slavin Ross, Michael C Hughes, and Finale Doshi-Velez. Right for the right reasons: Training differentiable models by constraining their explanations. *arXiv preprint arXiv:1703.03717*, 2017.

[67] Gabriel Erion, Joseph D Janizek, Pascal Sturmfels, Scott Lundberg, and Su-In Lee. Learning explainable models using attribution priors. *arXiv preprint arXiv:1906.10670*, 2019.

[68] Masahiro Mitsuhara, Hiroshi Fukui, Yusuke Sakashita, Takanori Ogata, Tsubasa Hirakawa, Takayoshi Yamashita, and Hironobu Fujiyoshi. Embedding human knowledge in deep neural network via attention map. *arXiv preprint arXiv:1905.03540*, 5, 2019.

[69] Mara Graziani, Sebastian Otálora, Henning Muller, and Vincent Andrearczyk. Guiding cnns towards relevant concepts by multi-task and adversarial learning. *arXiv preprint arXiv:2008.01478*, 2020.

[70] Yiqi Yan, Jeremy Kawahara, and Ghassan Hamarneh. Melanoma recognition via visual attention. In *International Conference on Information Processing in Medical Imaging*, pages 793–804. Springer, 2019.

[71] Laura Rieger, Chandan Singh, W James Murdoch, and Bin Yu. Interpretations are useful: penalizing explanations to align neural networks with prior knowledge. *arXiv preprint arXiv:1909.13584*, 2019.

[72] Mohammadhassan Izadyyazdanabadi, Evgenii Belykh, Claudio Cavallo, Xiaochun Zhao, Sirin Gandhi, Leandro Borba Moreira, Jennifer Eschbacher, Peter Nakaji, Mark C Preul, and Yezhou Yang. Weakly-supervised learning-based feature localization for confocal laser endomicroscopy glioma images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 300–308. Springer, 2018.

[73] Yoichiro Yamamoto, Toyonori Tsuzuki, Jun Akatsuka, Masao Ueki, Hiromu Morikawa, Yasushi Numata, Taishi Takahara, Takuji Tsuyuki, Kotaro Tsutsumi, Ryuto Nakazawa, et al. Automated acquisition of explainable knowledge from unannotated histopathology images. *Nature communications*, 10(1):1–9, 2019.

[74] Christoph Jansen, Thomas Penzel, Stephan Hodel, Stefanie Breuer, Martin Spott, and Dagmar Krefting. Network physiology in insomnia patients: Assessment of relevant changes in network topology with interpretable machine learning models. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 29(12):123129, 2019.

[75] Anna Paola Carrieri, Niina Haiminen, Sean Maudsley-Barton, Laura-Jayne Gardiner, Barry Murphy, Andrew Mayes, Sarah Paterson, Sally Grimshaw, Martyn Winn, Cameron Shand, et al. Explainable ai reveals key changes in skin microbiome associated with menopause, smoking, aging and skin hydration. *bioRxiv*, 2020.

[76] Achraf Essemlali, Etienne St-Onge, Maxime Descoteaux, and Pierre-Marc Jodoin. Understanding alzheimer disease's structural connectivity through explainable ai. In *Medical Imaging with Deep Learning*, pages 217–229. PMLR, 2020.

[77] Data Language (UK) Ltd. SCOPA - scalable, explainable AI.

[78] Decoded Health. The World's First Clinical Hyperautomation Platform - A force multiplier for physicians.

[79] Karan Sikka, Andrew Silberfarb, John Byrnes, Indranil Sur, Ed Chow, Ajay Divakaran, and Richard Rohwer. Deep Adaptive Semantic Logic (DASL): Compiling Declarative Knowledge into Deep Neural Networks. *arXiv preprint arXiv:2003.07344*, 2020.

[80] Shalini Ghosh, Daniel Elenius, Wenchao Li, Patrick Lincoln, Natarajan Shankar, and Wilfried Steiner. ARSENAL: automatic requirements specification extraction from natural language. In *NASA Formal Methods Symposium*, pages 41–46. Springer, 2016.

[81] Deutsches Forschungszentrum für Künstliche Intelligenz GmbH (DFKI). exAID - Bringing the power of Deep Learning to clinical practice.

[82] Fabrizio Nunnari and Daniel Sonntag. A CNN toolbox for skin cancer classification. *arXiv preprint arXiv:1908.08187*, 2019.

[83] Daniel Sonntag, Fabrizio Nunnari, and Hans-Jürgen Profitlich. The skincare project, an interactive deep learning system for differential diagnosis of malignant skin lesions. technical report. *arXiv preprint arXiv:2005.09448*, 2020.

[84] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability methods in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 9737–9748, 2019.

[85] Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE transactions on neural networks and learning systems*, 28(11):2660–2673, 2016.

[86] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages 9505–9515, 2018.

[87] Fabian Eitel, Kerstin Ritter, Alzheimer's Disease Neuroimaging Initiative (ADNI, et al. Testing the robustness of attribution methods for convolutional neural networks in mri-based alzheimer's disease classification. In *Interpretability of Machine Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision Support*, pages 3–11. Springer, 2019.

[88] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.

[89] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.

[90] An-phi Nguyen and María Rodríguez Martínez. On quantitative aspects of model interpretability. *arXiv preprint arXiv:2007.07584*, 2020.

[91] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan. Interpreting interpretability: Understanding data scientists' use of interpretability tools for machine learning. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2020.

[92] Rory Sayres, Ankur Taly, Ehsan Rahimy, Katy Blumer, David Coz, Naama Hammel, Jonathan Krause, Arunachalam Narayanaswamy, Zahra Rastegar, Derek Wu, et al. Using a deep learning algorithm and integrated gradients explanation to assist grading for diabetic retinopathy. *Ophthalmology*, 126(4):552–564, 2019.

[93] E Tjoa and C Guan. A survey on explainable artificial intelligence (xai): towards medical xai. 2019: 21. *arXiv preprint arXiv:1907.07374*.

[94] Roberto V. Zicari. Z-Inspection: A process to assess trustworthy AI.

[95] Glyn Elwyn, Isabelle Scholl, Caroline Tietbohl, Mala Mann, Adrian GK Edwards, Catharine Clay, France Légaré, Trudy van der Weijden, Carmen L Lewis, Richard M Wexler, et al. "many miles to go. . .": a systematic review of the implementation of patient decision support interventions into routine clinical practice. *BMC medical informatics and decision making*, 13(2):1–10, 2013.

[96] Elodia B Cole, Zheng Zhang, Helga S Marques, R Edward Hendrick, Martin J Yaffe, and Etta D Pisano. Impact of computer-aided detection systems on radiologist accuracy with digital mammography. *American Journal of Roentgenology*, 203(4):909–916, 2014.

[97] Ajay Kohli and Saurabh Jha. Why CAD failed in mammography. *Journal of the American College of Radiology*, 15(3):535–537, 2018.

[98] Carrie J Cai, Emily Reif, Narayan Hegde, Jason Hipp, Been Kim, Daniel Smilkov, Martin Wattenberg, Fernanda Viegas, Greg S Corrado, Martin C Stumpe, et al. Human-centered tools for coping with imperfect algorithms during medical decision-making. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–14, 2019.

[99] Carrie J Cai, Jonas Jongejan, and Jess Holbrook. The effects of example-based explanations in a machine learning interface. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, pages 258–262, 2019.

[100] Qian Yang, Aaron Steinfeld, and John Zimmerman. Unremarkable ai: Fitting intelligent decision support into critical, clinical decision-making processes. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–11, 2019.

[101] Marina Jirotka, Rob Procter, Mark Hartswood, Roger Slack, Andrew Simpson, Catelijne Coopmans, Chris Hinds, and Alex Voss. Collaboration and trust in healthcare innovation: The eDiaMoND case study. *Computer Supported Cooperative Work (CSCW)*, 14(4):369–398, 2005.

[102] Emma Beede, Elizabeth Baylor, Fred Hersch, Anna Iurchenko, Lauren Wilcox, Paisan Ruamviboonsuk, and Laura M Vardoulakis. A human-centered evaluation of a deep learning system deployed in clinics for the detection of diabetic retinopathy. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–12, 2020.

[103] Andrew L Beam and Isaac S Kohane. Translating artificial intelligence into clinical care. *Jama*, 316(22):2368–2369, 2016.

[104] Mohammad R Arbabshirani, Brandon K Fornwalt, Gino J Mongelluzzo, Jonathan D Suever, Brandon D Geise, Aalpen A Patel, and Gregory J Moore. Advanced machine learning in action: identification of intracranial hemorrhage on computed tomography scans of the head with clinical workflow integration. *NPJ digital medicine*, 1(1):1–7, 2018.

[105] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. Multimodal explanations: Justifying decisions and pointing to the evidence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8779–8788, 2018.

[106] Kamran Alipour, Jurgen P Schulze, Yi Yao, Avi Ziskind, and Giedrius Burachas. A study on multimodal and interactive explanations for visual question answering. *arXiv preprint arXiv:2003.00431*, 2020.

[107] Chih-Kuan Yeh, Been Kim, Sercan Arik, Chun-Liang Li, Tomas Pfister, and Pradeep Ravikumar. On completeness-aware concept-based explanations in deep neural networks. *Advances in Neural Information Processing Systems*, 33, 2020.

[108] Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based explanations. In *Advances in Neural Information Processing Systems*, pages 9277–9286, 2019.

[109] Hyun-Lim Yang, Jong Jin Kim, Jong Ho Kim, Yong Koo Kang, Dong Ho Park, Han Sang Park, Hong Kyun Kim, and Min-Soo Kim. Weakly supervised lesion localization for age-related macular degeneration detection using optical coherence tomography images. *PloS one*, 14(4):e0215076, 2019.

[110] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[111] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

[112] Federico Cabitza, Raffaele Rasoini, and Gian Franco Gensini. Unintended consequences of machine learning in medicine. *Jama*, 318(6):517–518, 2017.

[113] Rich Caruana, Yin Lou, Johannes Gehrke, Paul Koch, Marc Sturm, and Noemie Elhadad. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1721–1730, 2015.

[114] Andreas Holzinger, Georg Langs, Helmut Denk, Kurt Zatloukal, and Heimo Müller. Causability and explainability of artificial intelligence in medicine. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 9(4):e1312, 2019.

[115] Andreas Holzinger, Chris Biemann, Constantinos S Pattichis, and Douglas B Kell. What do we need to build explainable ai systems for the medical domain? *arXiv preprint arXiv:1712.09923*, 2017.

[116] Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic Books, 2018.