

Article

# Analysis of Explainers of Black Box Deep Neural Networks for Computer Vision: A Survey

Vanessa Buhrmester , David Münch , and Michael Arens 

Fraunhofer IOSB, Gutleuthausstraße 1, 76275 Ettlingen, Germany

\* Correspondence: vanessa.buhrmester@iosb.fraunhofer.de; Tel.: +49 7243-992-167

**Abstract:** Deep Learning is a state-of-the-art technique to make inference on extensive or complex data. As a black box model due to their multilayer nonlinear structure, Deep Neural Networks are often criticized to be non-transparent and their predictions not traceable by humans. Furthermore, the models learn from artificial datasets, often with bias or contaminated discriminating content. Through their increased distribution, decision-making algorithms can contribute promoting prejudice and unfairness which is not easy to notice due to lack of transparency. Hence, scientists developed several so-called explanators or explainers which try to point out the connection between input and output to represent in a simplified way the inner structure of machine learning black boxes. In this survey we differ the mechanisms and properties of explaining systems for Deep Neural Networks for Computer Vision tasks. We give a comprehensive overview about taxonomy of related studies and compare several survey papers that deal with explainability in general. We work out the drawbacks and gaps and summarize further research ideas.

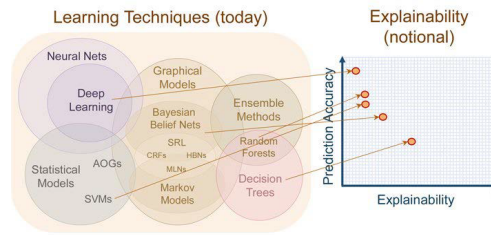
**Keywords:** Interpretability, explainer, explainer, explainable AI, Trust, Ethics, Black Box, Deep Neural Network.

## 1. INTRODUCTION

Artificial Intelligence (AI)-based technologies are increasingly being used to make inference on classification or regression problems: Automated image and text interpretation in medicine, insurance, advertisement, public video surveillance, job applications, or credit scoring save staff and time and are moreover practical successful. The severe drawback is that many of these technologies are black boxes and referenced results can hardly be understood by the user.

### 1.1. Motivation: More complex algorithms are hardly comprehensible.

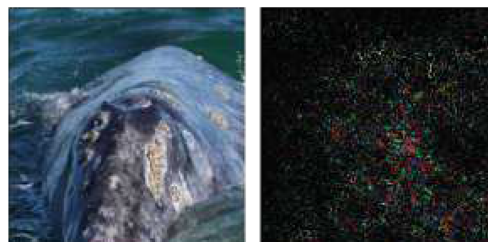
Latest models are more complex and Deep Learning (DL) architectures are getting deeper and deeper and millions of parameters are calculated and optimized by a machine. For example, the common network VGG-19 incorporates about 144 millions parameters that were optimized over millions or hundred thousands of images [1]. ResNet has about  $5 \cdot 10^7$  trainable parameters and for classifying one image it needs to execute about  $10^{10}$  floating point operations, see [2]. It is hardly traceable and not recalculate-able by humans. Metrics like accuracy or the mean average precision are depending on the quality of manually hand annotated data. However, these metrics are often the only values that evaluate the learning algorithm itself. The explainability of a Machine Learning (ML) technique is decreasing with an increasing prediction accuracy, and the prediction accuracy is growing with more complex models like Deep Neural Networks (DNNs), see Figure 1. But a trade-off between explainability and accuracy is not satisfying, especially in critical tasks.



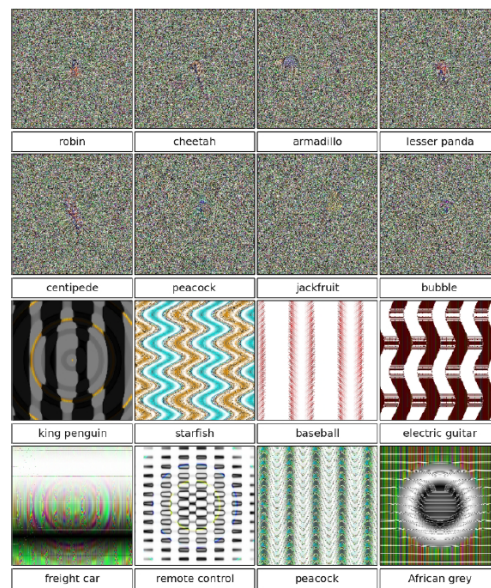
**Figure 1.** The explainability is decreasing with a growing prediction accuracy while the prediction accuracy is growing with more complex models like Deep Learning, Figure from [3] .

In the past, many scientists found weaknesses of Deep Learning models – even high performing models are affected. They showed, for instance, how easy one could fool object detectors with small changes in the input image or created adversarial examples to make them collapse, see [4], [5]. Furthermore, they draw attention on supposedly good models which focus on totally wrong features and just made good decisions by chance. The problem is that the neural network learns only from the training data, which should characterize the task. But suitable training data is tedious to create and annotate, hence it is not always perfect. If in the training data is a bias, the algorithm will learn it as well. The following examples present attacks on neural networks:

A classifier of enemy tanks and friendly tanks with high accuracy did not deliver good results in the application and was discovered to be a just a good classifier of sunny or overcast days [6]. The reason was that most of the photos for the training set of enemy tanks were taken on days with clouds on the sky, while the friendly ones were shot during sunny weather. The same problem happened as a dog-or-wolf-classifier turned out to be just a good snow detector [7] because of a bias in the background of the training images. There are several cases which underline the negative characteristics of a DNN. Changing only one pixel in an input image or one letter in a sentence could change the prediction, or even adding small-magnitude perturbations, [8], [9], [10], see also Figure 2. Adversarial examples with serious impact exist; fixed stickers on road signs [11] or extended adversarial patch attacks in optical flow networks [12] could lead to dangerous misinterpretation. Worn prepared glasses confuse face detectors by imitating special persons, [13], [14]. Further cases regard on Figure 3, 4, and 5. All this examples show how harmful it could be to rely on a black box with supposedly well performing results. However, currently applied DNN-methods and models are such vulnerable black boxes.



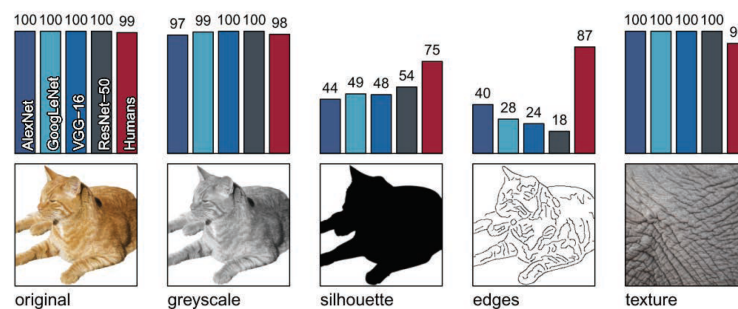
**Figure 2.** DeepFool [8] examines the robustness of neural networks. Very small noisy images were added on right classified images, humans cannot see the difference, but the algorithms changes its prediction:  $x$  (left) is correctly classified as whale, but  $x + r$  as turtle,  $r$  (right) is very small.



**Figure 3.** [15] created artificial images, that where unrecognizable by humans, but the state-of-the-art classifier was very confident that they were known objects. Images are either directly (top) or indirectly (bottom) encoded.



**Figure 4.** Texture-shape cue conflict [16]: texture (left) is classified as elephant, content (middle) is classified as cat, texture-shape (right) is classified as elephant because of a texture bias.



**Figure 5.** [16] show that Convolutional Neural Networks focus more on texture than on shapes and that it is a good idea to improve shape bias to get more reliable results in the way like humans would interpret the content of an image.

### 1.2. Also an ethical question: Why is explainability so important?

Machine Learning-models are tested by the engineer on validation data, this is the basis on which he evaluates how correct they work according to the chosen metrics. Although, the applicant applies the

model to the real world, he should believe that the results still are reliable. But how could he? However, fairness and reliability should be important properties of such models. To win the users trust, the link between the features and the predictions should be understandable. If a human is even judged by a machine, he or she should have the right of explanation. Since the European Union's new General Data Protection Regulation (in Germany DSGVO) was passed in the EU in May 2018, it will restrict the use of machine learning and automated individual decision-making and focus on protection of sensitive data of persons like their age, sex, ancestry, name or place of residence for instance. If a result affects users, they should be able to demand for explanations of the algorithmic decision that was made about them, see [17]. For example, if a doctor makes a mistake the patient wants to know why. Was the mistake excusable? Or does the mistake reflect negligence, even intentionally or occurred due to other factors? Similarly, if an algorithm fails and contributes to an adverse clinical event or malpractice, doctors need to be able to understand why it produced the result and how it reached a serious decision.

Although or just because a human decision is not free from prejudice, it must be ensured that an algorithm for selection procedure – e.g. for job proposals or suspended sentences – may not discriminate humans because of sex or origin etc. Disadvantages can arise here for individuals, groups of persons or a whole society. The prevailing conditions can thereby further deteriorate more and more caused by a gender-word bias [18]. An interesting example are word embeddings [19] that a Natural Language Processing (NLP) algorithm creates from training datasets, which are just any available texts of a special origin [20]. The procurement of women, disabled, black people etc. seem to be deep anchored in these texts, with the serious consequence that the model learns that to be currently. This reinforces discrimination and unfairness. Implicit Association Tests [21] have uncovered stereotype biases which people are not aware of. If the model supposes – and studies [19], [22], [23] show that – that doctors are male and nurses are female, furthermore, women are sensitive and men are successful, it will sort out all women who apply as a chief doctor – only because of their sex – no need to check their qualification. If in the train data foreigners have predominantly less income and increased unemployment, an automatic credit scoring model will suggest a higher interest rate or even refuse the request only because of the origin and without considering the individual financial situation. Unfairness will progress through these algorithms.

Hence, getting understanding and insights in the mechanisms should uncover these problems. Properties of a model like transparency and interpretability are basics to build patient, provider trust, and fairness. If this succeeds, the causes of the discrimination and serious mistakes can be remedied, additionally. There is the opportunity to improve the society through making automated decisions free of prejudice. Our contribution on the way achieving this goal is giving an overview about state-of-the-art-explainers with regard to their taxonomy through differing their mechanisms and technical properties. We do not just limit our work to explaining methods, but also look at the meaning of understanding a machine in general. To our knowledge this is the first survey paper that focuses on ML-black box DNNs for Computer Vision tasks.

## 2. Overview about explaining systems of DNNs

We just give a short introduction in early approaches to explain inner Machine Learning operations. After that we will focus only on understanding DNNs.

### 2.1. Early Machine Learning explaining systems

Early explaining systems for ML-black boxes go back to 1986 with generalized additive models (GAM) [24]. This is a global statistic model which uses smooth functions as a diagnostic tool. Later **Decision Trees** were successful classification tools that provide individual explanations, see [25], [26]. A Decision Tree is a tree-like graph of decisions and their possible consequences that visualizes an algorithm



that only contains condition control statements. Another approach [27] shows the marginal effect of one or two features on the prediction of learning techniques using Partial Dependence Plots (**PDP**). The method gives a statement about the global relationship of a feature and whether its relation to the outcome is linear, monotonous or more complex. PDP is the average of Individual Conditional Expectation (**ICE**) over all features. ICE [28] points to how the prediction changes if a feature changes. PDP is limited to two features. Another example is the early use of explainable AI [29] which was developed as a simulator game for the commercial platform training aid Full Spectrum Command motivated by previous work such as [30]. The proposed procedure of [31] based on a set of assumptions, which allows to explain the decision for a particular label of a single data instance for several classification models. The framework provides local explanation vectors as class probability gradients which yield the relevant features of every point of interest in the data space.

We do not respond further to early studies of explaining systems of **Random Forests** [32], **Naïve Bayesian** classifiers [33], [34], [35], Support Vector Machines (**SVMs**) [36], [37], or other early Machine Learning prediction methods [38].

## 2.2. Mechanism and properties of DNN-explainers



**Figure 6.** [7] presented Local Interpretable Model-agnostic Explanations (LIME) which can explain the predictions of any agnostic black box classifier and any data. Here the superpixels – areas of an input image – are highlighted that are most responsible for top three image classification predictions. (1) original image (2) explaining electric guitar, (3) explaining acoustic guitar, and (4) explaining labrador.

In the last years, the importance of DNNs in inference tasks grew rapidly, and with the increasing complexity of these models, also the need of better explanations did. Most used DNNs for image or video processing [39], [40], [41] are Convolutional Neural Networks (CNNs) [42], for videos [41] or sequences of text Recurrent Neural Networks (RNNs) [43], and especially for language modeling [44] Long-Short Term Memories (LSTMs) [45]. There exist some general surveys of methods for explaining several machine learning black boxes, so called eXplanatory Artificial Intelligence (XAI), that cover a wide spectrum of AI-based black boxes, for instance see [46], [47]. However, we just want to focus on the black box DNN and deepen the insights, especially in Computer Vision. In the following we describe state-of-the-art explainers in these tasks. There are explainers that create visualizations which give an explanation by digesting the mechanisms of a model down to images which themselves have to be interpreted. Saliency maps of important features are calculated, for example, they show super-pixels or given keywords that have influenced the prediction most, see [7] and Figure 6. The problem is that these explainers are in turn black boxes. White box explainers use methods that gain insights and show all causal effects, for instance linear regression or Decision Trees. Black box explainers do not require access to the internals and do not disclose all feature interaction. There are mainly two kinds of explaining models:

- **Ante-hoc or intrinsically interpretable** models, see [48]. Ante-hoc systems give explanations starting from the beginning of the model. For instance enabling one to gauge how certain a neural network is about its predictions.

- **Post-hoc** techniques entail baking explainability into a model from its outcome, like marking what part of the input data is responsible to the final decision, for example in LIME. These methods can be applied more easily to different models, but say less about the whole model in general.

They can be also split in

- **local**: the algorithm can be explained only for each single prediction and
- **global**: the whole system can be explained and the logic can be followed from the input to every possible outcome,

and even in

- **model specific**, tied to a particular type of black box or data and
- **model agnostic**, indifferently usable.

[49] defined key terms like “explanation”, “interpretability”, and “explainability” philosophically. Important is that interpretability and explainability is not the same, although it is often used interchangeably. An explanation is much more concrete, a coherence of facts can be describe with words while an interpretation is just a substantial formation that arises in the head. They also describe the difficulties of both interpretability and completeness, so a compromise is needed. For more details see the following. We compile the most important definitions of properties of explainers and their obvious connections in a more technical way:

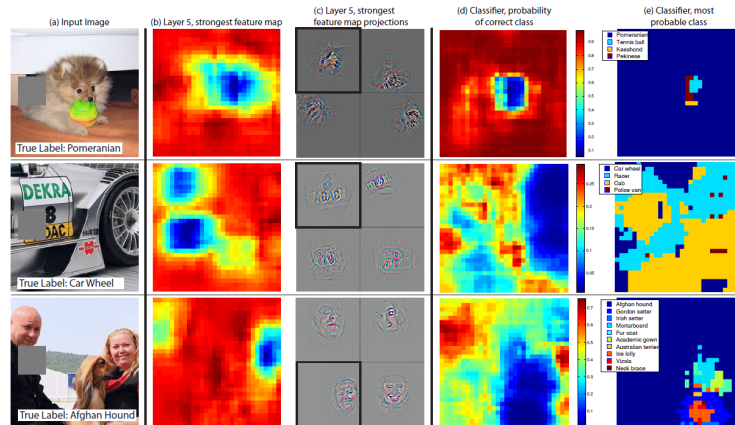
- An **explainer** or **explanator** is a synonym for an explaining system, that gives an answer to a questions. For instance if the question is how a machine is working, the explainer makes the internal structure of a machine more transparent to humans. A further question could be why a prediction was made instead of another, so the explainer should point to where the decision boundaries between classes are and why particular labels are predicted for different data points [50].
- **Interpretability** is a substantial first step to reach a comprehension of a complex coherence in some level of detail but is insufficient alone, see [49].
- **Explainability** includes interpretability but not always reverse. It provides relevant responses to questions and subdivides their meaning in understandable terms to a human, see [49], [51].
- **Comprehensibility** or **understandable explanation**. An understandable explanation must be created by a machine in a given time and can be comprehend by a user, who need not to be an expert, but has an educational background, in a given time (e.g. one hour or one day).
- **Completeness**. A **complete explanation** records all possible factors from input to output of a model. A DNN with its millions of parameters is too complex, hence a complete explanation would not be understandable. That makes it necessary to focus on the most important reasons and not all of them.
- **Compactness**. A **compact explanation** has a finite number of aspects. Because the parameters and operations of a DNN are finite, one can get a complete explanation of a DNN after a finite number of steps. That is why compactness follows from completeness if the regarded connections are finite. A DNN can be explained for instance completely and compactly or compactly and understandably.

Other, in our eyes less gentle aspects that are mentioned in the literature are fidelity, trust, intelligibility, privacy, usability, monotonicity, causality, scalability, and generality, see [52], [51], [46]. Now we investigate the employed mechanisms of explainers:

- **Visualizations**. To visualize an explanation there are many options [53]. One tool is to look at the **activations** produced on each layer of a trained CNN as it processes an image or video. Another one enables visualizing features at each layer of a DNN via regularized optimization in image space [54]. Visualizations of particular neurons or neuron layers show responsible features that

lead to a maximum activation or highest possible probability of a prediction, see [55] and can be split in **generative models** or **saliency maps**. To create a map of import pixels one can repeatedly feed an architecture with several portions of inputs and compare the respective output. Or one visualizes them directly by going rearwards through the inverted network from an output of interest. Also grouped in this category is exploiting neural networks with **activation atlases** through feature inversion. This method can reveal how the network typically represents some concepts [56].

- **Gradients** or variants of (guided) **backpropagation** can emphasize important unit changes and thereby draw attention to sensitive features or input data areas, see [57], [58], [59]. With these techniques it is also able to produce artificial prototype class member images that maximize a neuron activation or class confidence value, see [60], [61].
- Regarding image or text portions that **maximize the activation of interesting neurons or whole layers** can lead to interpretation of the responsible area of individual parts of the architecture.
- **Deconvolution** or **inverting DNNs** is applied to create typical inputs or parts of an input, that fits to a desired output of the network, a special layer or single unit, see [62], Figure 7, [63].
- Another method is **decomposition**, isolation, transfer or limitation of portions of networks, e.g. layers to get further insights in which way single parts of the architecture influences the results, see [64], or Deep Taylor Decomposition (DTD), [65]. **Automatic Rule Extraction** and **Decision Trees** are anchored in this area, too.



**Figure 7. Deconvnet [62]:** They plot three examples of input image (a), strongest feature map (2) and feature map projections (3) of layer 5, and the classifier with the probability of correct class (4) and most probable class (5), respectively.

We give some examples of explaining approaches that can be placed in the last point: A global ante-hoc method for tabulator data is Bayesian Rule List (**BRL**) [66], [67]. BRL is a generative model that yields a posterior distribution over possible decision lists which consist of a series of if-then-statements. The “if”-statements define a partition of a set of features and the “then”-statements correspond to the predicted outcome of interest. The work was developed from preliminary versions that used a different prior, called Bayesian List Machine [68]. Similar to **DeepRED** [69] the rule generation “KnowledgeTron” (**KT**) [70] applied an if-then-rule for each neuron, layer by layer. Another option in this field is to decompose a DNN in **Decision Trees**, e.g. DeepRED or [26]. Decision trees were used since the 1990s to explain machine learning tasks, but applied to DNNs their generation is quite expensive and the comprehensibility suffers from the necessary size and number of trees. Decision Trees are able to explain an algorithm completely but with DNNs there is a conflict with comprehensibility.

### 2.3. Selected DNN-explainers presented

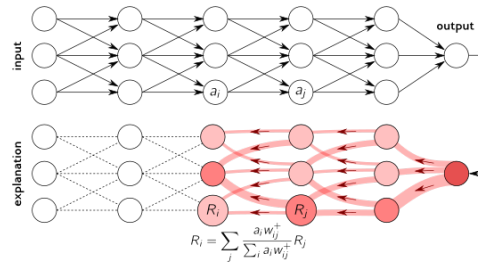
Let us give a technical overview about selected explainer models, we will put a focus on the category Computer Vision:

The Counterfactual Impact Evaluation (CIE) method, [71], [72], is a local method of comparison for different predictions. Counterfactuals are contrastive. They explain why a decision was made instead of another. A counterfactual explanation of a prediction may be defined as the smallest change to the feature values that changes the prediction to a predefined output. They could be employed to DNNs with any data type.

Famous work to visualize and understand Convolutional Neural Networks is **Deconvnet** [62]: Deconvnet is a calculation of a backward convolutional network that reuses the weights at each layer from the output layer back to the input image. The employed mechanisms were deconvolution and unpooling which are especially designed for CNNs with convolutions, maxpooling, and REctified Linear Units (ReLUs). The method makes it possible to create feature maps of an input image that activates certain hidden units most, linked to a particular prediction, see Figure 7. With their propagation technique they identify the most responsible patterns for this output. The patterns are visualized in the input space. Deconvnet is limited to max-pooling layers and in the absence of a particular theoretical criterion which could directly connect the prediction to the created input patterns. To close that gap [57] proposes a new and efficient way following the initial attempt of Zeiler and Fergus. They replaced max-pooling by a convolutional layer with increased stride and called the method **The All Convolutional Net**. The performance on image recognition benchmarks was similar well. With this approach they were able to analyze the neural network by introducing a novel variant of the Deconvnet to visualize the concepts learned by higher network layers of the CNN. The problem of max-pooling layers is that they are not invertible in general. That is the reason why Zeiler and Fergus computed positions of maxima within each pooling region and used these “switches” in the Deconvnet for a discriminative reconstruction. Not using max-pooling Springenberg et al. could directly display learned features and was not conditioned on an image. Furthermore, for higher layers, they produced sharper, more recognizable visualizations of descriptive image regions than previous methods. This is in agreement with the fact that higher layers learn more invariant representations.

[54] introduced two tools to aid interpreting DNNs in a global way. First they have displayed the neurons activations produced on each layer of a trained CNN processing an image or sequence of images. They found that looking at live activations that change in response to input images helps to build valuable intuitions about the inner mechanisms of these neural networks. The second tool was built on previous versions which calculated less recognizable images. Some novel regularization methods that were combined produce qualitatively clearer and more interpretable visualizations and enable plotting features at each layer via regularized optimization in image space.

Gradient-based is layer-wise Relevance Propagation (LRP) [64] which is suffering from shattered gradient problems, see Figure 8. It relies on a conservation principle to propagate the outcome decision back without using gradients. The idea behind is a decomposition of prediction function as a sum of layer-wise relevance values. When LRP is applied to deep ReLU networks, LRP can be understood as a deep Taylor decomposition of the prediction. This principle ensures that the prediction activity is fully redistributed through all the layers onto the input variables. More about how to explain nonlinear classification decisions with Deep Taylor Decomposition see [65]. They decompose the network classification decision into contributions of its input elements and assess the importance of single pixels in image classification tasks. Their method efficiently utilizes the structure of the network by backpropagating the explanations from the output to the input layer and display the connections in heat maps.



**Figure 8.** Layer-wise Relevance Propagation (LRP) [64] is a gradient method suffering from shattered gradient problems. The idea behind is a decomposition of the prediction function as a sum of layer-wise relevance values. When LRP is applied to deep ReLU networks, LRP can be understood as a deep Taylor decomposition of the prediction.

Also a global and ante-hoc model is a joint framework for description and prediction, presented by [73]. The model creates Black box Explanations through Transparent Approximations (**BETA**). It learns a compact two-level decision set in which each rule explains parts of the model behavior unambiguously and is a combined objective function to optimize these aspects: high agreement between explanation and the model, little overlaps between decision rules in the explanation, and the explanation decision set is lightweight and small.

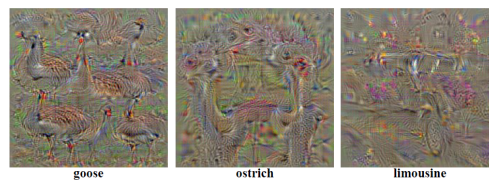
An interpretable end-to-end explainer for healthcare is the REverse Time AttentIoN mechanism **RETAIN** [74] for application to Electronic Health Records (EHR) data. The approach mimics physician practice by attending the EHR data, two RNNs are trained in a reverse time order with the goal of efficiently generating the appropriate attention variables. It is based on a two-level neural attention generation process that detects influential past visits and significant clinical variables to improve accuracy and interpretability.

Another technique was realized by [61]. To find **prototype class members**, they created input images that have the highest probability to be predicted as certain classes of a trained CNN. Their tools are Taylor series, based on partial derivatives to display input sensitivities in images. A few years later [60] developed this idea further by synthesizing the preferred inputs for neurons in neural networks via deep generator networks for activation maximizing. The first algorithm is the generator and creates synthetic prototype class members that look real. The second algorithm is the black box classifier of the artificial image whose classification probability should be maximized. To view the prototype images, see Figure 9. Another related derivative-based method is **DeepLift** [75]. It propagates activation differences instead of gradients through the network. Partial derivatives do not explain a single decision but point to what change in the image could make a change in the prediction.

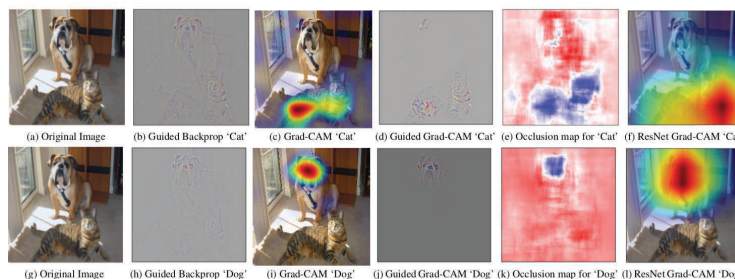




**Figure 9.** To explain what a black box classifier network comprehends as a class member being [60] made synthetic prototype images that look real. They were created by a deep generator network and classified by the black box neural network.



**Figure 10.** Deep inside convolutional networks [61]: created input images that have the highest probability to be predicted as certain classes of a trained CNN. Here one can see the created prototypes of the class goose, ostrich, limousine (left to right).



**Figure 11.** Gradient-weighted Class Activation Mapping (Grad-CAM) [76] explained the outcome decision cat or dog, respectively, of an input image using the gradient information to understand the importance of each neuron in the last convolutional layer of the CNN.

[77] showed that some convolutional layers behave as unsupervised object detectors. They use global average pooling and create heat maps of a pre-softmax layer that point out the regions of an image which is responsible for a prediction. The method is called Class Activation Mapping (**CAM**). Upon this was created Gradient-weighted Class Activation Mapping (**Grad-CAM**) [76], see Figure 11, which is applicable to several CNN model-families, classification, image captioning, visual question answering, reinforcement learning, or re-training. A outcome decision can be explained by Grad-CAM through using the gradient information to understand the importance of each neuron in the last convolutional layer of the CNN. The Grad-CAM localizations are combined with existing high-resolution visualizations to obtain high-resolution class-discriminative guided Grad-CAM visualizations as saliency masks. On the methods CAM and Grad-CAM were built **Grad-CAM++** [78] that gives human interpretable visual explanations

of CNN-based predictions of multiple tasks like classification, image captioning, or action recognition. Grad-CAM++ explains by regarding occurrences of multiple object instances in an image, combining the positive partial derivatives of feature maps of a convolutional rear layer with a weighted special class score.

To mark the most responsible pixels or areas of pixels of an image for a special class prediction it is a promising idea to increase human understanding. The approach of [79] focuses on single words of a caption generated by a RNN and highlights the region of the image which is most important for this word, see Figure 12.

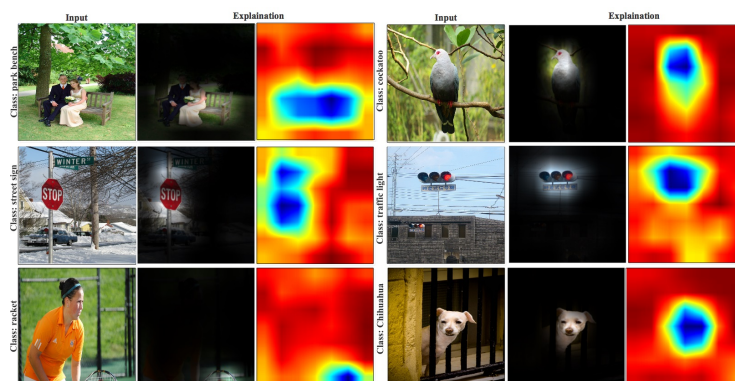


A **dog** is standing on a hardwood floor.

**Figure 12.** Deep Learning, [79]: This method highlights the region of an image (dog) which is most important for the part “dog” of the predicted output “A dog is standing on a hardwood floor” of a trained CNN.

Much more general is Local Interpretable Model-agnostic Explanations (**LIME**) presented by [7], which can explain the predictions of any agnostic black box classifier and any data, see Figure 6. It is a post-hoc, local model, interpretable and model-agnostic. LIME focuses on feature importance and gives outcome explanations: it highlights the super-pixels of the regions of the input image or the words from a text or tabular which are relevant for the given prediction. The disadvantage of this popular explainer is that it is itself a black box and in addition to this, [80] pointed to the poor performance of LIME during their proposed evaluation metrics correctness, consistency, and confidence in comparison to the other regarded explainers Grad-CAM, Smooth-Grad and IG (described in the following). Furthermore, one can explain only images which can be split in super-pixels. The authors do not describe how to explain video object detection or segmentation networks. In addition to this, the Submodular Pick (**SP-LIME**) algorithm judges whether you can trust the whole model or not. It selects a picked diverse set of representative instances with LIME explanations via submodular optimization. The user should evaluate the black box by regarding the feature words of the selected instances. With the knowledge it is also possible to improve a bad model. SP-LIME was researched with text data, but the authors claim that it can be transferred to any data type models.

Another approach that focuses on the most discriminative region in an image to explain an automatic decision is Deep Visual Explanation (**DVE**) [81], see Figure 13. They were inspired from CAM and Grad-CAM and tested the explainer to randomly chosen images from the COCO dataset [82], applied to the pre-trained neural network VGG-16 using Kullback Leibler (KL)-divergence [83]. They captured the discriminative areas of the input image by considering the activation of high and low spatial scales in Fourier space.

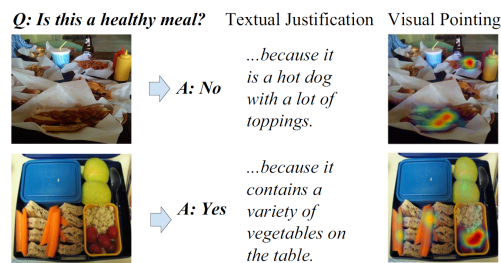


**Figure 13.** Deep Visual Explanation [81] highlights the most discriminative region in an image of six examples (park bench, cockatoo, street sign, traffic light, racket, chihuahua) to explain the decision made by the VGG-16.

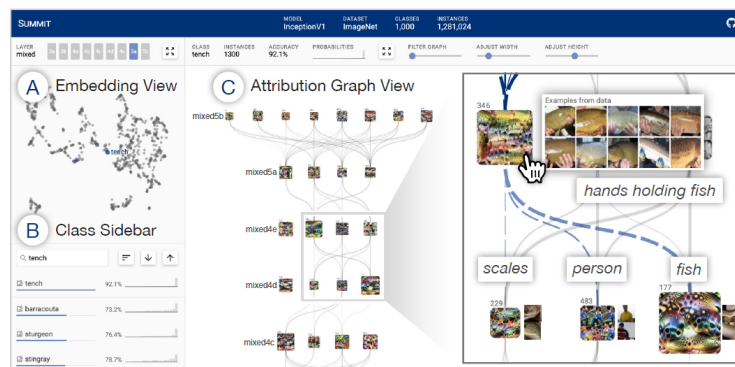
With their conditional multivariate model Prediction Difference Analysis (**PDA**) [84] have concentrated on explaining visualizations of natural and medical images in classification tasks. Their goal has been to improve and interpret DNNs. Their technique is based on the uni-variate approach of [85] and the idea that the relevance of an input feature with respect to a class can be estimated by measuring how the prediction changes if the feature is removed. Zintgraf et al. remove several features at one time using their knowledge about images by strategically choosing patches of connected pixels as feature sets. Instead of going through all individual pixels, they regard all patches of a special size implemented in a sliding window fashion. They visualize the effects of different window sizes and marginal versus conditional sampling and display feature maps of different hidden layers and top scoring classes.

[86] described **Smooth-Grad** that reduces visual noise and hence, improves visual explanations how a DNN is making a classification decision. Comparing their work to several gradient-based sensitivity map methods like LRP, DeepLift and Integrated Gradients (**IG**) [87] that estimate the global importance of each pixel and create saliency maps, shows that Smooth-Grad focuses on local sensitivity and calculates averaging maps with a smoothing effect made from several small perturbations of a input image. The effect is enhanced by further training with these noisy images and finally reaches an impact on the quality of sensitivity maps by sharpening them. The work of [80] evaluated the explainers LIME, Grad-CAM, Smooth-Grad and IG with regard to the properties correctness, consistency, and confidence and came to the result that Grad-CAM often performs better than others.

To improve and expend understanding Multimodal Explanation (**ME**) [88], a local, post-hoc approach gave visual and textual justifications of predictions with the help of two novel explanation datasets through crowd sourcing. The employed tasks were classification decision for activity recognition and visual question answering, see Figure 14. The visual explanation was created by an attention mechanism which conveys knowledge about what region of the image is important for the decision. This explanation guides to generate the textual justification out of a LSTM feature that is a prediction of a classification problem over all possible justifications.

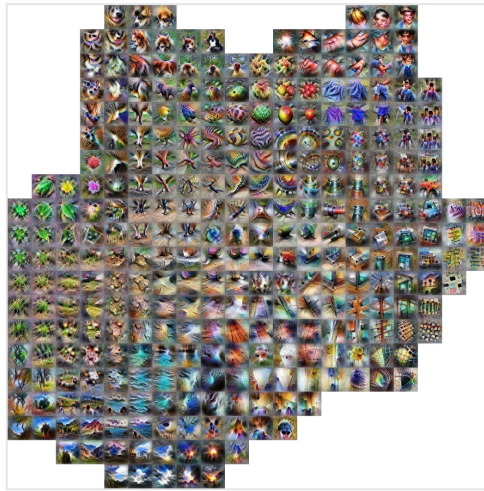


**Figure 14.** Multimodal Explanation (ME) [88] explains by two types of justifications of visual question answering task: The example shows two images with food and the question is, if they contain healthy meals or not. The explanations of the answers “Yes” or “No” are given textual in justifying in a sentence and visual in pointing out the most responsible areas of the image.



**Figure 15.** Summit [89] visualizes what features a Deep Learning model has learned and how those features are connected to make predictions. The Embedding View (A) shows which classes are related to each other, Class Sidebar (B) is linked to the embedding view listing all classes sorted in several ways, and the Attribution Graph (C) summarizes crucial neuron associations and substructures that contribute to a model’s prediction.

In this year a new and extensive visualized approach was created by [89], see Figure 15, showing what features a Deep Learning model has learned and how those features interact to make predictions. Their model is called **Summit** and combines two scalable tools: (1) activation aggregation discovers important neurons, and (2) neuron-influence aggregation identifies relationships among such neurons. An attribution graph that reveals and summarizes crucial neuron associations and substructures that contribute to a model’s outcomes is created. Summit combines famous methods like computing synthetic prototypes of features and showing examples from the dataset that maximize special neurons of different layers. Deeper in the graph is examined how the low-level-features combine to create high-level features. Also novel is exploiting neural networks with **activation atlases** [56], see Figure 16. This method uses feature inversion to visualize millions of activations from an image classification network to create an explorable activation atlas of features the network has learned. Their approach is able to reveal visual abstractions within a model and even high-level misunderstandings in a model that can be exploited. Activation atlases is a novel way to peer into convolutional vision networks and represents a global, hierarchical, and human-interpretable overview of concepts within the hidden layers.



**Figure 16.** Activation atlases with 100,000 activations, [56].

Just in the last months the importance of interpretability was growing, that is why there appeared several single studies, investigating the contribution of some aspects of a neural network like the impact of color [90], texture [16] etc. without explaining a whole model extensive, however, which all contribute to a deeper understanding.

In Table 1 we give an overview about the presented explainers of DNNs, sorted by data and year. The main techniques and properties are mentioned for a short comparison. The property model-agnostic is abbreviated with agn.



**Table 1.** Overview about some explainers for DNNs. Ordered by model or paper name, reference (author and year), data type, main method and main properties.

Model	authors	data	method	properties
Deep Inside	[61]	image	saliency mask	local, post-hoc
Deconvnet	[62]	image	gradients	global, post-hoc
All-CNN	[57]	image	gradients	global, post-hoc
Deep Visualization	[54]	image	neurons activation	global, ante-hoc
Deep Learning	[79]	image	visualization	local, post-hoc
Show, attend, tell	[91]	image	saliency mask	local, ante-hoc
LRP	[64]	image	decomposition	local, ante-hoc
CAM	[77]	image	saliency mask	local, post-hoc
Deep Generator	[60]	image	gradients, prototype	local, ante-hoc
Interpretable DNNs	[92]	image	saliency map	local, ante-hoc
VBP	[93]	image	saliency maps	local, post-hoc
DTD	[65]	image	decomposition	local, post-hoc
Meaningful	[94]	image	saliency mask	local, post-hoc, agn.
PDA	[84]	image	feature importance	local, ante-hoc
DVE	[83]	image	visualization	local, post-hoc
Grad-CAM	[76]	image	saliency mask	local, post-hoc
Grad-CAM++	[78]	image	saliency mask	local, post-hoc
Smooth-Grad	[86]	image	sensitivity analysis	local, ante-hoc
ME	[88]	image	visualization	local, post-hoc
Summit	[89]	image	visualization	local, ante-hoc
Activation atlases	[56]	image	visualization	local, ante-hoc
SP-LIME	[7]	text	feature importance	local, post-hoc
Rationalizing	[95]	text	saliency mask	local, ante-hoc
Generate reviews	[96]	text	neurons activation	global, ante-hoc
BRL	[66]	tabular	decision tree	global, ante-hoc
TreeView	[97]	tabular	decision tree	global, ante-hoc
IP	[98]	tabular	neurons activation	global, ante-hoc
KT	[70]	any	rule extraction	global, ante-hoc
Decicion Tree	[26]	any	rule extraction	global, ante-hoc, agn.
CIE	[71]	any	feature importance	local, post-hoc
DeepRed	[69]	any	rule extraction	global, ante-hoc
LIME	[7]	any	feature importance	local, post-hoc, agn.
NES	[99]	any	rule extraction	local, ante-hoc
BETA	[73]	any	decision tree	global, ante-hoc
PALM	[100]	any	decision tree	global, ante-hoc
DeepLift	[75]	any	feature importance	local, ante-hoc
IG	[87]	any	sensitivity analysis	global, ante-hoc
RETAIN	[74]	EHR	Reverse Time Atten.	global, ante-hoc

#### 2.4. Analysis of understanding and explaining methods

We still want to go into studies on the general analysis of explainability and machine understanding. To reach a better understanding of multilayer neural network [101] analyzed pre-training and fine-tuning of several classification and object recognition tasks. They found that some CNN-learned features are grandmother-cell-like, that are cells in the human brain which only respond to very specific and complex visual stimuli, such as the face of one's grand-mother [102] and that a longer pre-training significantly improves the performance. [53] investigated different methods to compute heat maps in Computer Vision application. They concluded that layer-wise relevance propagation algorithms (e.g. LRP) were qualitatively and quantitatively superior in explaining what made a DNN arrive at a particular classification decision to the sensitivity-based approaches or the deconvolution methods [62]. The inferior methods were much

noisier and less suitable for identifying the most important regions with respect to the classification task. Their work did not give an answer how to make a more detailed analysis of prioritization of image regions or even how to quantify the heatmap quality. [103] criticized the explaining methods Deconvnet, guided Backpropagation and LRP not to produce the theoretically correct explanation for a linear model and their contributions to understanding were scarce. Based on an analysis of linear models, see also [104], [65] they propose a generalization that yielded the two neuron-wise explanation techniques **PatternNet** (for signal visualization) and **PatternAttribution** (for decomposition methods) by taking the data distribution into account. They demonstrated that their methods were sound and constitute a theoretical, qualitative and quantitative improvement towards explaining and understanding DNNs.

**SHAP** (SHapley Additive exPlanations) [105] has been developed to interpret a models prediction by additive feature attribution methods. It unifies six previously existing explanations and assigns each feature an importance value for a particular prediction. Its new components include: (1) the identification of a new class of additive feature importance measures, and (2) theoretical results showing there is a unique solution in this class with a set of desirable properties. The framework SHAP reaches a higher performance and a better consistency with human intuition than previous approaches in this tasks, as the authors claim.

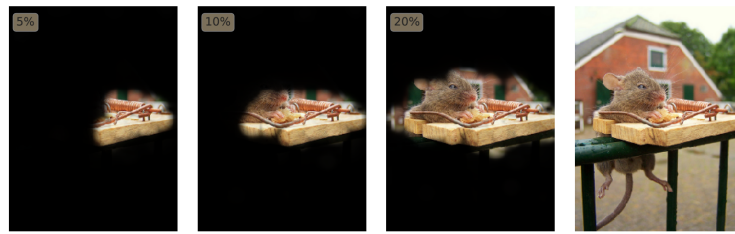
[106] made an approach to quantify interpretability of Deep Visual Representation by proposing a general framework called **Network Dissection** by evaluating the alignment between individual hidden units and a set of semantic concepts. The proposal raises how to detect and evaluate disentangled representations and investigates if there is a coherence between hidden units and a special alignment of feature space. Also the influence of training conditions on entanglement of explanations is regarded and confirmed to have a significant effect of the representation learned by the hidden units. Their framework Network Dissection comes to the conclusion that interpretability is not an axis-independent phenomenon, which is consistent with the hypothesis that interpretable units indicate a partially disentangled representation.

With an overview of interpretability of Machine Learning [49] try to explain explanations. They define some key terms and review a number of approaches towards classical explainable AI systems also focusing on Deep Learning tasks. Furthermore, they investigate the role of single layers, individual units and representation vectors in explanation of deep network representations. Finally, they present a taxonomy that examines what is being explained by these explanations. They summarize that it is not obvious what the best purpose or type of explanation metric is and should be and give the advice to combine explaining ideas from different fields.

Another approach towards understanding is to evaluate the human-interpretability of explanation [107]. They investigated the consistence of output of a ML-system with its input and the supposed rationale. Therefor they carried out user-studies to identify what kind of increases in complexity have the most dominant effect on the time humans need to take for verification the rationale, and which seem to be more insensitive. Their study quantifies what kind of explanation makes them to be most understandable by humans. As a main result they found out that in general, greater complexity results in higher response times and lower satisfaction.

Even simple interpretable explainers do mostly not quantify if the user can trust them. A study on trust in black box models and post-hoc explanations, [108] provides the problems in main problems in literature and their kind of black box systems. They evaluate three different explanation approaches: (1) based on the users' initial trust, (2) the users' trust in the provided explanation of three different post-hoc explanator approaches and (3) the established trust in the black box by a within-subject design study. The participants where asked if they trust that a special algorithm works well in the real world, if they suppose it to be able to distinguish between the classes and why. The results of their work led to the conclusion,

that although the black box prediction model and explanation are independent of each other, trusting the explanation approach is an important aspect of trusting the prediction model.



**Figure 17.** Extremal perturbations [109]: The example shows the regions of an image (boxed) that maximally affect the activation of a certain neuron in a DNN (“mousetrap” class score). For clarity, the masked regions are blacked out. In practice, the network sees blurred regions.

A discussion of some existing approaches to perturbation analysis is done in the study of [109], Figure 17. Their work, based on [94], found adversarial effect of perturbations on the network’s output. Extremal perturbations are regions of an input image that maximally affect the activation of a certain neuron in a DNN. Measuring the effects of perturbations of the image is an important family of attribution methods. Attribution aims at characterising the response of DNN by looking at which parts of the network’s input are the most responsible ones for determining its prediction, which is mostly done by several kinds of backpropagation. Fong et al. investigate their effect as a function of the area. In particular, they have visualized the difference between perturbed and unperturbed activations using a representation inversion technique. They introduced TorchRay [110], a PyTorch interpretability library.

### 2.5. Open problems of understanding DNNs

When summarizing the functionality of explainers, one notices that some facts are hard to measure: First, the time, which is needed to understand the decision is difficult to take. Local working explainers can deliver a root case for each prediction, but how many examples are necessary to look at, to be sure, that all results and thereby the black box is faithful? In addition to this, the model complexity of several explainers is different. The complexity is often calculated as an opposed term to interpretability. Complexity of a black box can be expressed for instance as the number of non-zero weights at neural networks or the depth of trees for decision. But the complexity of the explanation could depend of the complexity of the black box.

More work must be done also in data exploration: Interpretable data in the mentioned papers are mainly images, texts and tabular data, which is all easily interpretable data for humans. Missing is general data like vectors, matrices, or complex spatio-temporal numbers. Of course they have to be transformed before analyzing to be understandable from our brains. Sequences, networks etc. could be an input in a black box, but until now, such models are not explained.

There is no agreement how to quantify the grade of explanation. This is an open problem. Some metrics to evaluate explainers are proposed, e.g. [80], but unfortunately, many of them tend to suffer from major drawbacks like computational cost [111] or simply focusing on one desirable attribute of a good explainer [112]. But a definition with properties of a DNN model like reliability, soundness, completeness, compactness, comprehensibility, and the knowledge of breaking points of an algorithm are still missing. However, there is a need in focusing on uniform definitions. Important to regard could be the robustness of model or input data which indicates how robust the system is in small changes of the architecture or the test date. To compare models also reliability and fairness need to be quantified. Our further work will be done here.

### 3. Conclusion

In this paper we give reasons why it is necessary to comprehend black box Deep Neural Networks: Adversarial attacks that are not understandable by humans pose a serious threat to the robustness of DNNs. Furthermore, we expound why models learn prejudices through contaminated training data, hence, through their widespread application they are responsible for grown unfairness. On the other hand, novel laws demand for the right of explanation for users of intelligent decision-making systems. One first solution is the development of explainers. We give a taxonomic introduction in the main definitions, properties, and mechanisms of explaining systems. Unlike others – to our best knowledge – we focus only on state-of-the-art explainers for DNNs, especially on the area Computer Vision and compare their similar or different methods and representations of explanations. We differ the explainers with regard to technical mechanisms, application, data, and pros and cons. Finally, we introduce surveys and studies, that analyze or evaluate explaining systems and try to quantify machine understanding in general. In addition to this, summarizing open problems with an outlook of further ideas and specifying some missing definitions of explainers makes this work useful for Computer Vision scientists.

### References

1. Simonyan, K.; Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* **2014**.
2. He, K.; Zhang, X.; Ren, S.; Sun, J. Deep residual learning for image recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
3. Gunning, D. Explainable Artificial Intelligence (XAI). *DARPA* **2017**.
4. Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; Fergus, R. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199* **2013**.
5. Goodfellow, I.J.; Shlens, J.; Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* **2014**.
6. Freitas, A.A. Comprehensible Classification Models – a position paper **2014**.
7. Ribeiro, M.T.; Singh, S.; Guestrin, C. Why should i trust you?: Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, 2016.
8. Moosavi-Dezfooli, S.M. DeepFool: a simple and accurate method to fool deep neural networks **2016**.
9. Bose, A.J.; Aarabi, P. Adversarial attacks on face detectors using neural net based constrained optimization. *2018 IEEE 20th International Workshop on Multimedia Signal Processing (MMSP)*. IEEE, 2018.
10. Jia, R.; Liang, P. Adversarial examples for evaluating reading comprehension systems. *arXiv preprint arXiv:1707.07328* **2017**.
11. Eykholt, K.; Evtimov, I.; Fernandes, E.; Li, B.; Rahmati, A.; Xiao, C.; Prakash, A.; Kohno, T.; Song, D. Robust physical-world attacks on deep learning visual classification. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
12. Ranjan, A.; Janai, J.; Geiger, A.; Black, M.J. Attacking Optical Flow. *International Conference on Computer Vision*, 2019.
13. Sharif, M.; Bhagavatula, S.; Bauer, L.; Reiter, M.K. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2016.
14. Sharif, M.; Bhagavatula, S.; Bauer, L.; Reiter, M.K. A General Framework for Adversarial Examples with Objectives. *ACM Transactions on Privacy and Security (TOPS)* **2019**.
15. Nguyen. Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images **2015**.

16. Geirhos, R.; Rubisch, P.; Michaelis, C.; Bethge, M.; Wichmann, F.A.; Brendel, W. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness **2019**.
17. Goodman, F. European union regulations on algorithmic decision-making and a right to explanation. *ICML workshop on human interpretability in ML*, NY **2016**.
18. Holmes, J.; Meyerhoff, M. *The handbook of language and gender*; John Wiley & Sons, 2008.
19. Bolukbasi, T.; Chang, K.W.; Zou, J.Y.; Saligrama, V.; Kalai, A.T. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 2016.
20. Hirschberg, J.; Manning, C.D. Advances in natural language processing. *Science* **2015**.
21. Greenwald, A.G.; McGhee, D.E.; Schwartz, J.L. Measuring individual differences in implicit cognition: the implicit association test. *Journal of personality and social psychology* **1998**.
22. Chakraborty, T.; Badie, G.; Rudder, B. Reducing gender bias in word embeddings **2016**.
23. Font, J.E.; Costa-Jussa, M.R. Equalizing gender biases in neural machine translation with word embeddings techniques. *arXiv preprint arXiv:1901.03116* **2019**.
24. Hastie, T.; Tibshirani, R. Generalized additive model (GAM) **1986**.
25. Craven, M.; Shavlik, J.W. Extracting tree-structured representations of trained networks. *Advances in neural information processing systems*, 1996.
26. Friedman, J.; Hastie, T.; Tibshirani, R. *The elements of statistical learning*; Springer series in statistics New York, 2001.
27. Friedman, J.H. Greedy function approximation: a gradient boosting machine. *Annals of statistics* **2001**.
28. Goldstein, A.; Kapelner, A.; Bleich, J.; Pitkin, E. Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics* **2015**.
29. Van Lent, M.; Fisher, W.; Mancuso, M. An explainable artificial intelligence system for small-unit tactical behavior. *Proceedings of the national conference on artificial intelligence*, 2004.
30. Shortliffe, E.H. Mycin: A knowledge-based computer program applied to infectious diseases. *Proceedings of the Annual Symposium on Computer Application in Medical Care*. American Medical Informatics Association, 1977.
31. Baehrens, D.; Schroeter, T.; Harmeling, S.; Kawanabe, M.; Hansen, K.; Mäzler, K.R. How to explain individual classification decisions. *Journal of Machine Learning Research* **2010**.
32. Breiman, L. Random forests. *Machine learning* **2001**.
33. Kononenko, I. Inductive and Bayesian learning in medical diagnosis. *Applied Artificial Intelligence an International Journal* **1993**.
34. Becker, B.; Kohavi, R.; Sommerfield, D. Visualizing the simple Bayesian classifier. *Information visualization in data mining and knowledge discovery* **2001**.
35. Možina, M.; Demšar, J.; Kattan, M.; Zupan, B. Nomograms for visualization of naive Bayesian classifier. *European Conference on Principles of Data Mining and Knowledge Discovery*. Springer, 2004.
36. Poulet, F. Svm and graphical algorithms: A cooperative approach. *Fourth IEEE International Conference on Data Mining (ICDM'04)*. IEEE, 2004.
37. Hamel, L. Visualization of support vector machines with unsupervised learning. *2006 IEEE Symposium on Computational Intelligence and Bioinformatics and Computational Biology*. IEEE, 2006.
38. Breiman, L. Bagging predictors. *Machine learning* **1996**.
39. Ciresan, D.; Giusti, A.; Gambardella, L.M.; Schmidhuber, J. Deep neural networks segment neuronal membranes in electron microscopy images. *Advances in neural information processing systems*, 2012.
40. Ciresan, D.; Meier, U.; Schmidhuber, J. Multi-column deep neural networks for image classification. *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 2012.
41. Le, Q.V.; Zou, W.Y.; Yeung, S.Y.; Ng, A.Y. Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis. *CVPR 2011*. IEEE, 2011.
42. Krizhevsky, A.; Sutskever, I.; Hinton, G.E. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 2012.



43. Rumelhart, D.E.; Hinton, G.E.; Williams, R.J.; others. Learning representations by back-propagating errors. *Cognitive modeling* **1988**.
44. Collobert, R.; Weston, J.; Bottou, L.; Karlen, M.; Kavukcuoglu, K.; Kuksa, P. Natural language processing (almost) from scratch. *Journal of machine learning research* **2011**.
45. Hochreiter, S.; Schmidhuber, J. Long short-term memory. *Neural computation* **1997**.
46. Guidotti, R.; Monreale, A.; Ruggieri, S.; Turini, F.; Giannotti, F.; Pedreschi, D. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)* **2019**.
47. Goebel, R.; Chander, A.; Holzinger, K.; Lecue, F.; Akata, Z.; Stumpf, S.; Kieseberg, P.; Holzinger, A. Explainable AI: the new 42? International Cross-Domain Conference for Machine Learning and Knowledge Extraction. Springer, 2018.
48. Lipton, Z.C. The mythos of model interpretability. *arXiv preprint arXiv:1606.03490* **2016**.
49. Gilpin, L.H.; Bau, D.; Yuan, B.Z.; Bajwa, A.; Specter, M.; Kagal, L. Explaining explanations: An overview of interpretability of machine learning. 2018 IEEE 5th International Conference on data science and advanced analytics (DSAA). IEEE, 2018.
50. Ridgeway, G.; Madigan, D.; Richardson, T.; O’Kane, J. Interpretable Boosted Naïve Bayes Classification. 1998.
51. Doshi-Velez, F.; Kim, B. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608* **2017**.
52. Alexandrov, N. Explainable AI decisions for human-autonomy interactions. 17th AIAA Aviation Technology, Integration, and Operations Conference, 2017.
53. Samek, W.; Binder, A.; Montavon, G.; Lapuschkin, S.; Müller, K.R. Evaluating the visualization of what a deep neural network has learned. *IEEE transactions on neural networks and learning systems* **2016**.
54. Yosinski, J.; Clune, J.; Nguyen, A.; Fuchs, T.; Lipson, H. Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579* **2015**.
55. Dosovitskiy, A.; Brox, T. Inverting visual representations with convolutional networks. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016.
56. Carter, S.; Armstrong, Z.; Schubert, L.; Johnson, I.; Olah, C. Activation Atlas. *Distill* **2019**.
57. Springenberg, J.; Dosovitskiy, A.; Brox, T.; Riedmiller, M. Striving for Simplicity: The All Convolutional Net. ICLR (workshop track), 2015.
58. Zhang, R.; Isola, P.; Efros, A.A. Colorful image colorization. European conference on computer vision (ECCV). Springer, 2016.
59. Zhang, J.; Bargal, S.A.; Lin, Z.; Brandt, J.; Shen, X.; Sclaroff, S. Top-down neural attention by excitation backprop. *International Journal of Computer Vision* **2018**.
60. Nguyen, A.; Dosovitskiy, A.; Yosinski, J.; Brox, T.; Clune, J. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. Advances in Neural Information Processing Systems, 2016.
61. Simonyan, K.; Vedaldi, A.; Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034* **2013**.
62. Zeiler, M.D.; Fergus, R. Visualizing and understanding convolutional networks. Springer, 2014.
63. Mahendran, A.; Vedaldi, A. Understanding deep image representations by inverting them. Proceedings of the IEEE conference on computer vision and pattern recognition, 2015.
64. Bach, S.; Binder, A.; Montavon, G.; Klauschen, F.; Müller, K.R.; Samek, W. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one* **2015**.
65. Montavon, G.; Lapuschkin, S.; Binder, A.; Samek, W.; Müller, K.R. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition* **2017**.
66. Letham, B.; Rudin, C.; McCormick, T.H.; Madigan, D.; others. Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model. *The Annals of Applied Statistics* **2015**.
67. Ustun, B.; Rudin, C. Methods and models for interpretable linear classification. *arXiv preprint arXiv:1405.4047* **2014**.
68. Letham, B.; Rudin, C.; McCormick, T.H.; Madigan, D. An interpretable stroke prediction model using rules and Bayesian analysis. Workshops at the Twenty-Seventh AAAI Conference on Artificial Intelligence, 2013.

69. Zilke, J.R.; Mencía, E.L.; Janssen, F. DeepRED–Rule extraction from deep neural networks. *International Conference on Discovery Science*. Springer, 2016.
70. Fu, L. Rule generation from neural networks. *IEEE Transactions on Systems, Man, and Cybernetics* **1994**.
71. Bottou, L.; Peters, J.; Quiñonero-Candela, J.; Charles, D.X.; Chickering, D.M.; Portugaly, E.; Ray, D.; Simard, P.; Snelson, E. Counterfactual reasoning and learning systems: The example of computational advertising. *The Journal of Machine Learning Research* **2013**.
72. Hainmueller, J.; Hazlett, C. Kernel regularized least squares: Reducing misspecification bias with a flexible and interpretable machine learning approach. *Political Analysis* **2014**.
73. Lakkaraju, H.; Bach, S.H.; Leskovec, J. Interpretable decision sets: A joint framework for description and prediction. *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, 2016.
74. Choi, E.; Bahadori, M.T.; Sun, J.; Kulas, J.; Schuetz, A.; Stewart, W. Retain: An interpretable predictive model for healthcare using reverse time attention mechanism. *Advances in Neural Information Processing Systems*, 2016.
75. Shrikumar, A.; Greenside, P.; Kundaje, A. Learning important features through propagating activation differences. *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017.
76. Selvaraju, R.R.; Cogswell, M.; Das, A.; Vedantam, R.; Parikh, D.; Batra, D. Grad-cam: Visual explanations from deep networks via gradient-based localization. *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
77. Zhou, B.; Khosla, A.; Lapedriza, A.; Oliva, A.; Torralba, A. Learning deep features for discriminative localization. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
78. Chattopadhyay, A.; Sarkar, A.; Howlader, P.; Balasubramanian, V. Grad-CAM++: Generalized Gradient-based Visual Explanations for Deep Convolutional Networks **2017**.
79. LeCun, Y.; Bengio, Y.; Hinton, G. Deep learning. *nature* **2015**.
80. Anonymous. On Evaluating Explainability Algorithms. Submitted to International Conference on Learning Representations, 2020. under review.
81. Babiker, H.K.B.; Goebel, R. An introduction to deep visual explanation. *arXiv preprint arXiv:1711.09482* **2017**.
82. Lin, T.Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; Zitnick, C.L. Microsoft coco: Common objects in context. *European conference on computer vision*. Springer, 2014.
83. Babiker, H.K.B.; Goebel, R. Using KL-divergence to focus deep visual explanation. *arXiv preprint arXiv:1711.06431* **2017**.
84. Zintgraf, L.M.; Cohen, T.S.; Adel, T.; Welling, M. Visualizing deep neural network decisions: Prediction difference analysis. *arXiv preprint arXiv:1702.04595* **2017**.
85. Robnik-Šikonja, M.; Kononenko, I. Explaining classifications for individual instances. *IEEE Transactions on Knowledge and Data Engineering* **2008**.
86. Smilkov, D.; Thorat, N.; Kim, B.; Viégas, F.; Wattenberg, M. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825* **2017**.
87. Sundararajan, M.; Taly, A.; Yan, Q. Axiomatic attribution for deep networks. *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017.
88. Huk Park, D.; Anne Hendricks, L.; Akata, Z.; Rohrbach, A.; Schiele, B.; Darrell, T.; Rohrbach, M. Multimodal explanations: Justifying decisions and pointing to the evidence. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
89. Hohman, F.; Park, H.; Robinson, C.; Chau, D.H. Summit: Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations. *arXiv preprint arXiv:1904.02323* **2019**.
90. Buhrmester, V.; Münch, D.; Bulatov, D.; Arens, M. Evaluating the Impact of Color Information in Deep Neural Networks. *Iberian Conference on Pattern Recognition and Image Analysis (IbPRIA)*. Springer, 2019.
91. Xu, K.; Ba, J.; Kiros, R.; Cho, K.; Courville, A.; Salakhudinov, R.; Zemel, R.; Bengio, Y. Show, attend and tell: Neural image caption generation with visual attention. *International conference on machine learning*, 2015.

92. Sturm, I.; Lapuschkin, S.; Samek, W.; Müller, K.R. Interpretable deep neural networks for single-trial EEG classification. *Journal of neuroscience methods* **2016**.
93. Bojarski, M.; Choromanska, A.; Choromanski, K.; Firner, B.; Jackel, L.; Muller, U.; Zieba, K. Visualbackprop: visualizing cnns for autonomous driving. *arXiv preprint arXiv:1611.05418* **2016**.
94. Fong, R.C.; Vedaldi, A. Interpretable explanations of black boxes by meaningful perturbation. Proceedings of the IEEE International Conference on Computer Vision, 2017.
95. Lei, T.; Barzilay, R.; Jaakkola, T. Rationalizing neural predictions. *arXiv preprint arXiv:1606.04155* **2016**.
96. Radford, A.; Jozefowicz, R.; Sutskever, I. Learning to generate reviews and discovering sentiment. *arXiv preprint arXiv:1704.01444* **2017**.
97. Thiagarajan, J.J.; Kailkhura, B.; Sattigeri, P.; Ramamurthy, K.N. TreeView: Peeking into deep neural networks via feature-space partitioning. *arXiv preprint arXiv:1611.07429* **2016**.
98. Shwartz-Ziv, R.; Tishby, N. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810* **2017**.
99. Turner, R. A model explanation system. 2016 IEEE 26th International Workshop on Machine Learning for Signal Processing (MLSP). IEEE, 2016.
100. Krishnan, S.; Wu, E. Palm: Machine learning explanations for iterative debugging. Proceedings of the 2nd Workshop on Human-In-the-Loop Data Analytics. ACM, 2017.
101. Agrawal, P.; Girshick, R.; Malik, J. Analyzing the performance of multilayer neural networks for object recognition. European conference on computer vision (ECCV). Springer, 2014.
102. Quiroga, R.Q.; Reddy, L.; Kreiman, G.; Koch, C.; Fried, I. Invariant visual representation by single neurons in the human brain. *Nature* **2005**.
103. Kindermans, P.J.; Schütt, K.T.; Alber, M.; Müller, K.R.; Erhan, D.; Kim, B.; Dähne, S. Learning how to explain neural networks: Patternnet and patternattribution. *arXiv preprint arXiv:1705.05598* **2017**.
104. Haufe, S.; Meinecke, F.; Görgen, K.; Dähne, S.; Haynes, J.D.; Blankertz, B.; Bießmann, F. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *Neuroimage* **2014**.
105. Lundberg, S.M.; Lee, S.I. A unified approach to interpreting model predictions. Advances in Neural Information Processing Systems, 2017.
106. Bau, D.; Zhou, B.; Khosla, A.; Oliva, A.; Torralba, A. Network dissection: Quantifying interpretability of deep visual representations. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.
107. Narayanan, M.; Chen, E.; He, J.; Kim, B.; Gershman, S.; Doshi-Velez, F. How do humans understand explanations from machine learning systems? an evaluation of the human-interpretability of explanation. *arXiv preprint arXiv:1802.00682* **2018**.
108. El Bekri, N.; Kling, J.; Huber, M.F. A Study on Trust in Black Box Models and Post-hoc Explanations. International Workshop on Soft Computing Models in Industrial and Environmental Applications. Springer, 2019.
109. Fong, R.; Patrick, M.; Vedaldi, A. Understanding Deep Networks via Extremal Perturbations and Smooth Masks. *arXiv:1910.08485v1* **2019**.
110. Torchray. [github.com/facebookresearch/TorchRay](https://github.com/facebookresearch/TorchRay) **2019**.
111. Hooker, S.; Erhan, D.; Kindermans, P.J.; Kim, B. Evaluating feature importance estimates. *arXiv preprint arXiv:1806.10758* **2018**.
112. Yeh, C.K.; Hsieh, C.Y.; Suggala, A.S.; Inouye, D.; Ravikumar, P. How Sensitive are Sensitivity-Based Explanations? *arXiv preprint arXiv:1901.09392* **2019**.