

Explaining Clinical Decision Support Systems in Medical Imaging using Cycle-Consistent Activation Maximization

Alexander Katzmann^{a,b,*}, Oliver Taubmann^a, Stephen Ahmad^a, Alexander Mühlberg^a, Michael Sühling^a, Horst-Michael Groß^b

^a*Siemens Healthineers, Computed Tomography, 91301 Forchheim, Germany*

^b*Ilmenau, University of Technology, Neuroinformatics and Cognitive Robotics Lab, 98693 Ilmenau, Germany*

Abstract

Clinical decision support using deep neural networks has become a topic of steadily growing interest. While recent work has repeatedly demonstrated that deep learning offers major advantages for medical image classification over traditional methods, clinicians are often hesitant to adopt the technology because its underlying decision-making process is considered to be intransparent and difficult to comprehend. In recent years, this has been addressed by a variety of approaches that have successfully contributed to providing deeper insight. Most notably, additive feature attribution methods are able to propagate decisions back into the input space by creating a saliency map which allows the practitioner to “see what the network sees.” However, the quality of the generated maps can become poor and the images noisy if only limited data is available—a typical scenario in clinical contexts. We propose a novel decision explanation scheme based on CycleGAN activation maximization which generates high-quality visualizations of classifier decisions even in smaller data sets. We conducted a user study in which these visualizations significantly outperformed existing methods on the LIDC dataset for lung lesion malignancy classification. With our approach we make a significant contribution to a better understanding

*Corresponding author

Email address: alexander.katzmann@siemens-healthineers.com (Alexander Katzmann)

of clinical decision support systems based on deep neural networks and thus aim to foster overall clinical acceptance.

Keywords: Medical Imaging, Deep Neural Networks, Decision Explanation, CycleGANs, Saliency Maps

1. Introduction

Within the last years, clinical decision support using deep neural networks (DNNs) has become a topic of steadily growing interest. This includes applications in microscopy and histopathology [1, 2], time-continuous biosignal analysis [3, 4], and, quite prominently, medical image analysis for volumetric imaging data as generated by computed tomography [5, 6], positron emission tomography [7, 8] or magnetic resonance imaging [9, 10, 11]. In the field of medical imaging, recent work has demonstrated a variety of applications for DNNs, such as organ segmentation [12], anomaly detection [13], lesion detection [14], segmentation [15] and assessment [16], providing major advantages and even repeatedly outperforming gold-standard human assessment [17].

A nearby field of similarly growing research interest established with the publications of Kumar et al. and Aerts et al. [18, 19] is „Radiomics” using traditional machine learning (ML) techniques. Compared to deep learning techniques, traditional ML methods like random forests and support vector machines have a largely transparent decision-making process, which is generally easier to comprehend and/or depict – a clear argument for their preference in clinical practice. Many publications have shown the advantages of DNNs in comparison to traditional machine learning techniques, such as the ability to learn descriptive features from data instead of a complex and expensive handcrafted feature design, as well as an improved classification performance on medical imaging tasks [20, 21], with some architectures being on par with gold-standard human assessment [17]. However, as DNNs learn features from the given data, the semantic of these features is in general not immediately evident. Thus, clinicians understandably approach these methods with a high degree of skepticism.

1.1. Related Work

A variety of methods have been proposed for decision explanation in DNNs, ranging from model-dependent analysis methods, over model-agnostic methods, to image synthesis-based approaches. Common to these approaches is a visualization of the decision process, or part of it, by highlighting regions in the input image which the algorithm identified as decision-relevant within the given image context. This allows a human observer to visually inspect the classifiers's regions of interest and whether these match with expected image areas.

Model-dependent analysis methods like DeconvNet [22], (Grad)CAM(++) [23, 24, 25] and Attention Gated Networks [26] work by explicitly modifying the network structure in such a way that there is an immediate spatial representation of relevance as a condition for a successful classification, e.g. by predicting attention areas and masking activations in regions which are expected by the network to be non-relevant. Typically, methods like these result in blob-like structures (cf. [26], [25]) and, while giving a coarse idea of decision-relevant image areas, still provide only little insight into the decision-making process.

Model-agnostic methods like LRP [27], DeepLIFT [28], LIME [29] and SHAP [30], as demonstrated by Lundberg and Lee, utilize a mechanism called *additive feature attribution* and therefore provide comparable results [30]. The basic idea of these methods is to use a backward pass through the network for splitting "relevance" starting from the target output, and attributing it to the previous layers with respect to their respective contribution to the output activation. This procedure is done recursively until the input layer is reached. While employing an equal mechanism, these methods differ in the quality of their generated explanations, with user studies indicating the most intuitive results for DeepSHAP, a deep learning based approximation of SHAP inspired by DeepLIFT. However, as we will show in Sec. 4, the quality of visualization can be highly dependent on the classifier used, and may be poor if the classifier is trained on only few samples.

Regarding image synthesis-based approaches, recent methods attempted to maximize the class output probability of a given class using a gradient ascent

optimization on the image input space, e. g. Activation Maximization (AM) [31]. Other work used a more complex approach to achieve this goal, for instance by using a generative adversarial network (GAN)-based activation maximization [32]. GANs are a method for synthesizing images based on a zero-sum game-like training process [33], and are able to produce images of high quality. AM, however, often results in images of low realism [34], and while this can be improved with GAN-based AM approaches, recent work focused on visualizing the information within the network, rather than the decision-relevant areas for the given image context. For clinical acceptance in a medical environment, however, intuitive visualizations with *realistic images* that *fit the image context* are crucial, meaning that the images should reflect the original radiological image as best as possible. While having significantly contributed to the scientific State of the Art (SoA) on network visualization, current image synthesis-based approaches therefore do not satisfyingly address the task of decision explanation.

1.2. Outline

We propose a method which follows the notion of activation maximization and combines it with a cycle-consistent GAN (CycleGAN) to improve the overall realism of the generated decision-explanations. To the best of our knowledge, we are the first to propose a combination of CycleGANs and activation maximization for this purpose. In detail:

- We will discuss the steps needed for such an algorithm and explain how each of its parts participates in creating a realistic decision-explanation in Sec. 2.1-2.6.
- In Sec. 3 we are constructing an experimental setup employing both, a quantitative, as well as a qualitative evaluation technique to evaluate our method in comparison to SoA approaches.
- We will present and discuss the results achieved within both experiments in Sec. 4

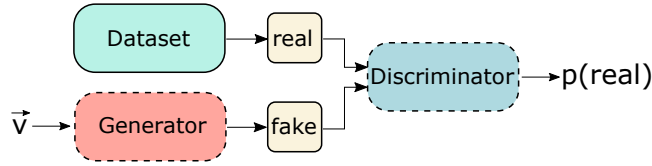


Figure 1: Basic GAN architecture. A generator generates a fake image from a random vector \vec{v} . A discriminator is trained to identify fakes, while the generator is then trained by backpropagating the inverted loss of the discriminator.

- In Sec. 5 an outlook is given, showing up future directions as well as limitations of the study at hand.

2. Material and Methods

Our proposed method is based on training a cycle-consistent generative adversarial network (CycleGAN [35], see Sec. 2.1-2.3) which is used for activation maximization (Sec. 2.4). Within this CycleGAN, each generator is trained to generate images which maximize the activation of one of the output neurons, i. e. class probabilities, of a given two-class classifier for which a decision-explanation should be given. After training (Sec. 2.5), a difference image of the results of both generators is created, emphasizing decision-relevant regions for arbitrary image inputs (Sec. 2.6).

2.1. Basic Architecture

Generally, a CycleGAN consists of a pair of generative adversarial networks (GANs), each consisting of one generator and one discriminator network. In a GAN, a generator network is trained to create synthetic images for the discriminator's image domain. Simultaneously, the discriminator is trained to distinguish real from fake images (see Fig. 1).

A CycleGAN employs this process cyclically using image-to-image translations *between* two domains. The discriminator networks are analogously trained to identify *fake images*, i. e. images which originate from the respective other

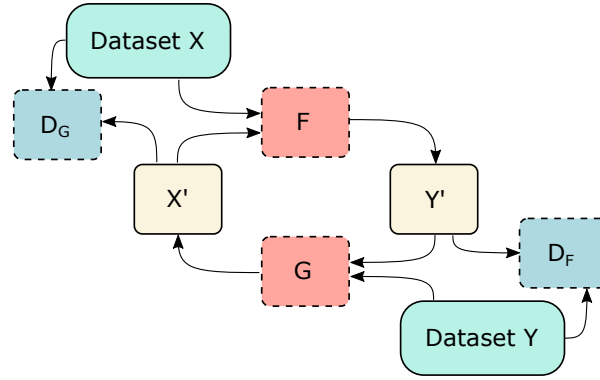


Figure 2: Basic CycleGAN architecture. A pair of GANs (F, D_F) and (G, D_G) is cyclically connected to form a domain transfer between domains X and Y by creating fake images X', Y' using the GAN training concept.

domain, resulting in a zero-sum game like behavior with increasingly growing image realism over the course of the training [35] (see Fig. 2).

In our approach we use the same architecture, but in contrast to the original approach train both generators with the *same* image domain and add an additional loss term to each of the generators to maximize one of the classifier's output activations. An overview of the overall architecture is depicted in Fig. 3 and will be the basis of the explanations in the following subsections.

2.2. Cycle Consistency

CycleGANs are trained to preserve *cycle consistency*, i.e. $G^+(G^-(x)) \approx G^-(G^+(x)) \approx x$ for image inputs x and generators G^+, G^- , meaning that a successive mapping through both generators should approximately reconstruct the original input. To enforce cycle consistency, we use the multiscale structural dissimilarity loss (*MS-DSSIM*) as proposed by Isola et al. [36]. Structural dissimilarity considers the image context, luminance, contrast and structure, and—from a perceptual point of view—is thus preferable to other losses such as L1 or binary cross entropy [37]. The corresponding loss function reads:

$$\mathcal{L}_{\text{DSSIM}}(x, y) = 1 - \frac{(2\mu_x\mu_y + c_1) \cdot (2\sigma_{xy} + c_2)}{2 \cdot (\mu_x^2 + \mu_y^2 + c_1) \cdot (\sigma_x^2 + \sigma_y^2 + c_2)}, \quad (1)$$

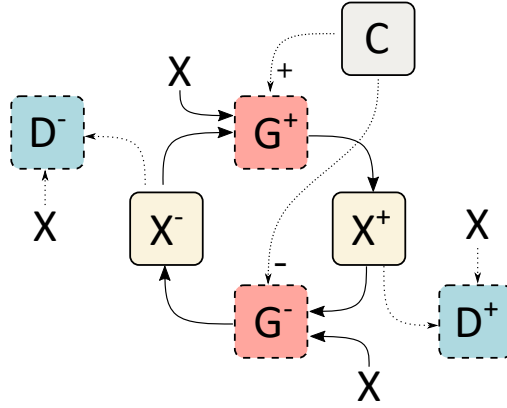


Figure 3: Overview on the general concept of our approach while training with randomly chosen images from image domain X . Generators G^+ and G^- each maximize one output of the classifier C by generating fake images X^+ and X^- , while discriminators D^+ , D^- are used to discriminate real from fake images.

with stabilizing parameters $c_1 = .01$ and $c_2 = .03$. As we furthermore want to preserve the average intensity, and DSSIM mainly accounts for covariances, we add an L1 term with a weight of .5, resulting in a cycle consistency loss of

$$\mathcal{L}_{\text{cycle}}(x, x') = \frac{1}{2}(|x - x'| + \mathcal{L}_{\text{DSSIM}}(x, x')) \quad (2)$$

This cycle consistency loss is applied to both generator networks G^+ and G^- as $\mathcal{L}_{\text{cycle}}(x, G^+(G^-(x)))$ and $\mathcal{L}_{\text{cycle}}(x, G^-(G^+(x)))$, respectively. Additionally, changes to the original image should be kept as small as possible, as only directly decision-relevant regions should be highlighted in the final visualization. To this end, we add additional similarity losses $\mathcal{L}_{\text{sim.}}(x, G^+(x))$ and $\mathcal{L}_{\text{sim.}}(x, G^-(x))$ using the same loss definition as above with $\mathcal{L}_{\text{sim.}} \equiv \mathcal{L}_{\text{DSSIM}}$.

2.3. Domain Transfer

We chose the Markovian discriminator approach commonly known as *PatchGAN* [36]. In a PatchGAN, the discriminator network's final output is not reduced to a binary (or categorical) variable, but rather preserved as an image-like output, with each output neuron representing the likelihood of its receptive field, i. e. patch, to stem from a real image. We will call these outputs *likelihood maps*.

As shown by Isola et al. [36], the PatchGAN structure prevents the discriminator from focusing on image artifacts, such as the ones created by convolutions at image corners, as otherwise it might not provide useful feedback for the generator anymore. We use multiple intermediate scales, i. e. receptive field sizes, of the same discriminator to further improve the image quality (cf. [38, 39]).

The PatchGAN structure expects varying difficulties of identifying fake images, depending on the image region, with some easily assessable, e. g. at image borders (see above), and some rather hard. Difficult regions therefore contribute less to improvements as they provide less information for the generator. Given *two* images, one real and one fake, the task of deciding which one is real becomes easier, as it reduces to a categorical decision, creating a more stable feedback and thus facilitates learning. To achieve this, our discriminator is given the likelihood maps of two images, one fake and one real, with image patches extracted using weight sharing. At each image position, a two-neuron dense layer with softmax activation is appended. The patches are shuffled and reordered after classification to prevent the emergence of a *real*- and a *fake*-neuron. The discriminators are trained using categorical cross entropy:

$$\mathcal{L}_{\text{CE}}(y, \hat{y}) = y \cdot \log(\hat{y}) + (1 - y) \cdot \log(1 - \hat{y}) \quad (3)$$

where y indicates whether a sample is real or generated and \hat{y} being the discriminator's prediction. The generators are trained by back-propagating the discriminator loss with inverted label information while freezing the discriminator weights. It is noteworthy that none of the generators receive label information *at any time*. However, while both generators are trained *using the same images*, they aim for maximizing opposite class activations (see Fig. 3).

2.4. Activation Maximization

As pointed out in Sec. 2.1, each generator is trained to maximize one of the classifier's output activations using Activation Maximization (AM). The AM is realized using an additional loss term while training the generators. We again use categorical cross entropy $\mathcal{L}_{\text{AM}}(l, \hat{l}) = \mathcal{L}_{\text{CE}}(l, \hat{l})$, where l denotes the label

the generator is trained to maximize, e. g. $l = 1$ for G^+ and $l = 0$ for G^- , and \hat{l} denotes the classifier's prediction for the current sample. The *actual* label of the sample is not used in this process.

2.5. Training Overview

The generators and discriminators are trained alternately using random pairs of images (x_a, x_b) from the same domain as the classifier input. First the discriminators, afterwards both generators are trained on one batch at a time. Combining the loss functions from above, the overall generator loss for G^+ reads:

$$\begin{aligned} \mathcal{L}_{\text{gen}}(x_a, x_b, l, G^+, G^-, D^+, C) = & \mathcal{L}_{\text{cycle}}(x_a, G^-(G^+(x_a))) + \mathcal{L}_{\text{sim.}}(x_a, G^+(x_a)) + \\ & \mathcal{L}_{\text{CE}}(D^+(G^+(x_a), x_b), 0) + \mathcal{L}_{\text{AM}}(l, C(x_a)) \end{aligned} \quad (4)$$

with x_a, x_b being images in the classifier image domain X , C being the classifier to be explained, the label l which should be maximized by the generator G^+ , its respective discriminator D^+ , and G^- being the opposite generator. The loss function is analogously applied to generator G^- with inverted $+$ and $-$. The procedure is repeated until the overall loss function of both generators converges. Using this loss definition, the CycleGAN is trained to

1. maximize the output probability of the class assigned to each of the generators (**AM loss**) with only small changes (**similarity loss**),
2. generate images indistinguishable from real images (**discriminator loss**),
3. be cycle-consistent (**cycle-consistency loss**)

2.6. Relevancy Visualization

Based on the loss function above, the generators are trained to provide two slightly modified versions $G^+(x), G^-(x)$ of the original images x for which the decision of the classifier should be explained, with each of them maximizing one of the classifier's output activations. As the generators are trained to introduce as few changes as possible (*similarity loss*), the modified regions are expected to be immediately decision-relevant within the given image context. Assuming

that the generators G^+ and G^- are trained on the same data set and realize the transitions $C(G^+(x)) \rightarrow 1$ and $C(G^-(x)) \rightarrow 0$, respectively, differences between x , $G^+(x)$, and $G^-(x)$ can be attributed to semantic differences with respect to the classifier's decision strategy. Let us define the unadjusted decision relevance per class as:

$$\Delta^+ = G^+(x) - x, \quad \Delta^- = G^-(x) - x. \quad (5)$$

We define overall relevance R as:

$$\begin{aligned} R &= \Delta^+ - \Delta^- \\ &= (G^+(x) - x) - (G^-(x) - x) \\ &= G^+(x) - G^-(x). \end{aligned} \quad (6)$$

Assuming a successful training of G^+ and G^- according to the above-mentioned targets, R has the following properties: it shows low absolute magnitudes in areas for which a modification would not result in an activation modification for the classifier, and high magnitudes in areas either only relevant for G^+ , or G^- , or both, but with different signs, as they would result in an output modification for the classifier. Therefore, R is an immediate indicator of decision relevance.

3. Experimental

For the experimental evaluation, we chose a two-step experimental setup. First, we train a classification network on the well-known LIDC-IDRI computed tomography dataset for lung lesion malignancy estimation [40, 41, 42]. After validating its ability to separate malignant from benign cases, our method for decision explanation is trained on the created classifier as described in Sec. 2.1-2.6. For ensuring a successful training of the CycleGAN, a quantitative evaluation is done, analyzing whether the generated samples are indeed modifying the class output probabilities of the classifier in the expected way. In a second step, a double-blind user study is employed to compare our method to various SoA approaches for decision explanation. The generated results are comprehensively evaluated and statistically tested.

3.1. Dataset

We used the LIDC-IDRI dataset for lung lesion malignancy classification. The data was preprocessed in accordance with prior work by Nibali et al. [43]. Lesions were extracted as axial slices at 64 x 64 pixels using bilinear interpolation and represented a window of 45 x 45 mm in world coordinates around the lesion's midpoint. Only lesions with at least three radiological annotations were used, assuming a lesion with a median malignancy rating above three to be malignant and below benign. As in Nibali et al. [43], samples with a borderline malignancy rating of three were discarded to enforce separability. The resulting data set contained 772 samples with 348 positive and 424 negative samples. 70 % of the data were used for training (236/301), 30 % as test set (112/123).

3.2. Architecture and Training

A ConvNet was used for malignancy classification, built of four convolutional blocks using strided convolutions with a stride of 1 in the first and 2 for the subsequent layers, followed by batch normalization and ReLU activation. In order, the blocks had 32, 64, 128 and 256 kernels. After flattening, a softmax layer was appended. The classifier was trained until convergence using weighted categorical cross entropy. Afterwards, our method for decision explanation was trained. The CycleGAN's generator networks were based on U-Net [44], using a depth of 3, with 3 convolutions per stage. The convolutional layers had 48, 96, 192 and 384 kernels at stages 0, 1, 2 and 3, respectively. The discriminators were based on the same architecture, but clipped to the outputs of stages 2 and 3 (cf. PatchGAN, see Sec. 2.3). The architecture of the generator and discriminator networks are depicted in Figs. 4 and 5, respectively.

3.3. Quantitative Evaluation

For the quantitative evaluation, we first evaluate whether the trained classifier is adequate for the task at hand. Therefore we evaluate the classifier's performance for lung lesion malignancy classification on the test set using various classification metrics, such as accuracy, sensitivity, specificity, positive and

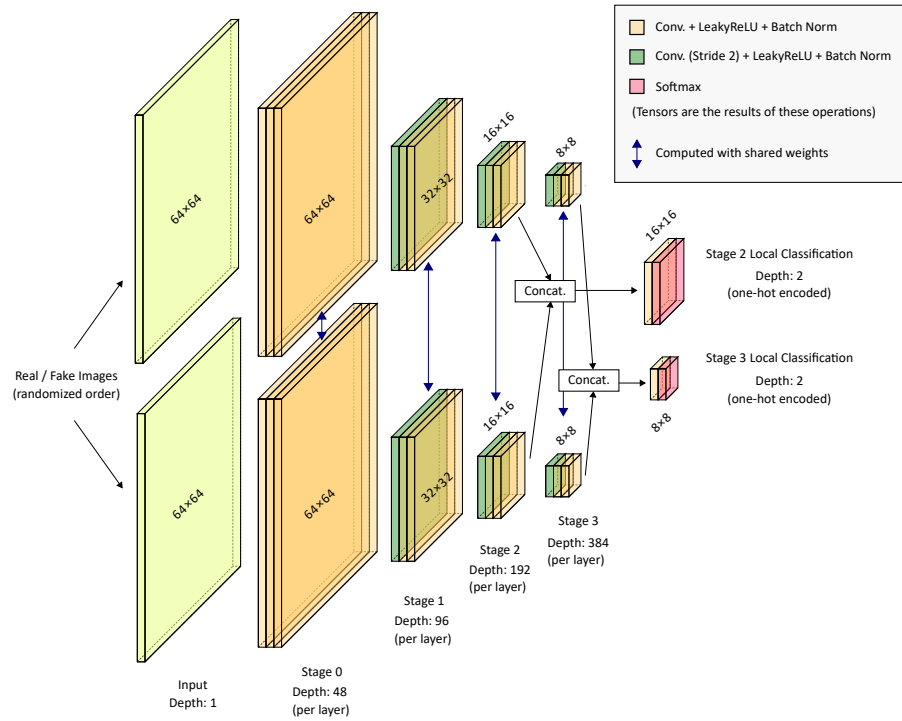


Figure 5: Overview on the discriminator network architecture. The architecture resembles the generator architecture from Fig. 4, but is cut off after stages 2 and 3 for multiscale patch classification (cf. [38, 39]). Each output neuron represents a receptive field of the input space (PatchGAN). Feature extraction is done on two images simultaneously, one real and one fake with randomized order, using weight sharing. The final classification is done by concatenating the features and appending a 2-neuron softmax output for each pair of patches (see Sec. 2.3).

negative predictive value, informedness, markedness, Matthews correlation coefficient and area under the curve (AUC). After ensuring the classifier is applicable, we evaluate whether our proposed CycleGAN network is capable of a successful domain transfer by evaluating the classifier’s output class probabilities on the original data, as well as the modified versions, i. e. the outputs of G^+ and G^- , and thus whether a change in the indicated regions is associated with an actual modification of the class output probabilities in the expected way.

3.4. User Study

For the qualitative evaluation, we asked $N = 8$ medical engineers with multiple years of experience in the field of medical algorithm development (average: 6,5 years) to evaluate the results in terms of a) *intuitive validity* of the visualization („Does it look reasonable?”), b) *semantic meaningfulness* in the context of lung lesion malignancy classification („Does it make sense?”), and c) overall *image quality* („Does it look good?”). For each criterion we analyzed the inter-observer reliability using pairwise Pearson product-moment correlation ρ over all participants with questionnaires Q as:

$$\rho = \frac{1}{\binom{N}{2}} \sum_{i=1}^{N-1} \sum_{j=i+1}^N r(Q_i, Q_j), \quad (7)$$

We double-blindly evaluated our method against DeepSHAP [30] using the reference implementation from Lundberg et al., as well as DeepTaylor [45] and LRP [27] using the iNNvestigate library from Alber et al. [46]. Each participant received images of 24 lesions (12 benign, 12 malignant) with the classifier’s decisions being visualized by each of the methods, depicted as a colored overlay over the original lesion image, resulting in a total of 192 answers per criterion and algorithm, 576 samples for each algorithm, or 2,304 answers in total. The image order was randomized for each question by a computer program to avoid an identification of the decision-explanation algorithm which generated the image. Two questionnaire variants (A and B) have been generated to avoid order effects [47]. The images were randomly drawn from the correctly classified test

Metric	Result
Accuracy	.809 [.757,.860]
F1	.801 [.737,.856]
Sensitivity	.813 [.737,.884]
Specificity	.805 [.729,.866]
PPV	.790 [.715,.861]
NPV	.826 [.757,.894]
Informedness	.618 [.514,.719]
Markedness	.616 [.511,.718]
MCC	.617 [.511,.718]
AUC	.809 [.757,.860]

Table 1: Test set performance of the used ConvNet classifier for malignancy estimation with 95 % CI.

set samples to avoid mixing up classification and visualization errors. The participants were asked to assign an integral score from -4 (low) to 4 (high) for each criterion to each image. Afterwards, the scores were collected and de-randomized by the computer program again. All responses were adjusted by z-normalizing over all answers for the same lesion and criterion. The results of the different decision explanation methods were compared against each other with respect to their average adjusted scores as well as their average rank scores and statistically tested using a two-tailed t -test with $t(N - 2) = t(6)$. Finally, after verifying its applicability using the Kaiser-Meyer-Olkin measure [48, 49], a principal component analysis (*PCA*) was employed to analyze whether the *preferability* of an image can be described by a general factor Φ that covers most of the variance.

4. Results

4.1. Quantitative Results

Classifier Performance. The employed classifier achieves a sensitivity of .813 95% CI [.737,.856] and a specificity of .805 [.729,.866] with an area under the curve (AUC) of .800 [.741, .847], giving a reasonable baseline for decision explanation. More detailed performance results can be found in Tab. 1.

Domain transfer. The classifier's average malignancy estimates on the test data were .470, 95 % CI [.418,.522] before modification, .289 [.250,.330] after negative transfer, and .820 [.788, .848] after positive transfer, i.e. x , $G^-(x)$ and $G^+(x)$. The pairwise differences between the transferred and the original domains were each highly significant with $p \ll 10^{-5}$ (two-tailed t -test, $t(233) = 17.9$ for $x/G^+(x)$, 14.5 for $x/G^-(x)$), implying a substantial modification of the classifier's class output probabilities using the proposed CycleGAN architecture.

4.2. User Study

The results of the user study are depicted in Tab. 2 - 4 and visualized in Fig. 7. A qualitative comparison of our approach to DeepSHAP, DeepTaylor and LRP is depicted in Fig. 6. Regarding the tested criteria, our method outperformed all other tested methods with average raw values of 1.92, 95 % CI [0.50, 3.13], 1.76 [0.25, 3.13], and 1.04 [-0.63, 2.50] for intuitive validity, semantic meaningfulness and image quality, respectively. Average values for the second best method for each criterion were -1.52 (DeepTaylor, $t(6) = 3.07, p = .018$, two-tailed t -test), -1.54 (DeepSHAP, $t(6) = 3.26, p = .014$), and -.85 (DeepTaylor, $t(6) = 1.71, n.s.$). Regarding intuitive validity and semantic meaningfulness, our method was significantly superior to all other tested methods with $p < .05$. Regarding image quality, significant superiority could be shown to DeepSHAP and LRP ($p = .015/.015$), while for DeepTaylor the result was non-significant with $p = .132$.

The z-adjusted values of our method were 2.84, 95 % CI [1.50,3.98], 2.78 [1.39,4.10] and 2.04 [0.50,3.41] for intuitive validity, semantic meaningfulness

	Intuitivity	Semantics	Quality
Ours	1.92 [0.50, 3.13]	1.76 [0.25, 3.13]	1.04 [-0.63, 2.50]
DeepSHAP	-1.56* [-3.13, 0.13]	-1.54* [-3.13, 0.25]	-2.02* [-3.25,-0.63]
DeepTaylor	-1.52* [-3.25, 0.38]	-1.72* [-3.38, 0.13]	-0.85 [-2.50, 0.75]
LRP	-2.53** [-3.63,-1.25]	-2.57** [-3.75,-1.13]	-2.16* [-3.38,-0.75]

Table 2: Raw questionnaire average results per algorithm with 95 % CI (higher is better). p -values for two-tailed t -test with $t(6)$ are indicated as $p < .05^*$, $p < .01^{**}$. The best result is marked bold.

	Intuitivity	Semantics	Quality
Ours	2.84 [1.50, 3.98]	2.78 [1.39, 4.10]	2.04 [0.50, 3.41]
DeepSHAP	-0.64* [-2.00, 0.86]	-0.52* [-1.97, 1.05]	-1.02* [-2.12, 0.16]
DeepTaylor	-0.60* [-2.16, 1.10]	-0.71* [-2.23, 1.00]	0.14 [-1.32, 1.57]
LRP	-1.61** [-2.72,-0.33]	-1.55** [-2.73,-0.11]	-1.16* [-2.24,-0.02]

Table 3: Questionnaire results after z-adjustment with 95 % CI (higher is better). p -values for two-tailed t -test with $t(6)$ are indicated as $p < .05^*$, $p < .01^{**}$. The best result is marked bold.

	Intuitivity	Semantics	Quality
Ours	1.28 [1.00,1.75]	1.36 [1.00,1.88]	1.52 [1.06,2.19]
DeepSHAP	2.78** [2.19,3.38]	2.75* [2.13,3.31]	3.04* [2.50,3.50]
DeepTaylor	2.71* [2.06,3.31]	2.75* [2.06,3.38]	2.36 [1.75,3.00]
LRP	3.23*** [2.75,3.63]	3.14** [2.63,3.56]	3.07* [2.50,3.56]

Table 4: Average rank per criterion and question with 95 % CI (lower is better). p -values for two-tailed t -test with $t(6)$ are indicated as $p < .05^*$, $p < .01^{**}$, $p < .001^{***}$. The best result is marked bold.

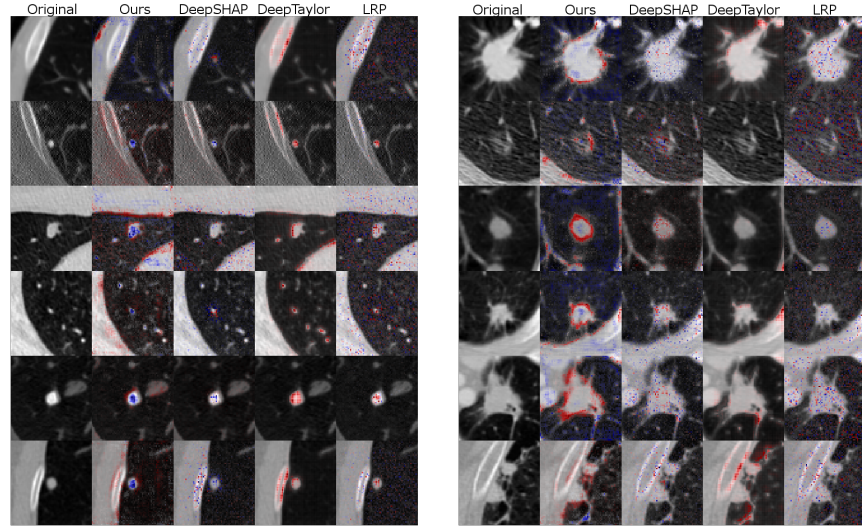


Figure 6: Example images for benign (left) and malignant (right) lung lesions with blue/red indicating benignity/malignancy-indicating regions for multiple methods (f.l.t.r.: Original, Ours, DeepSHAP, DeepTaylor, LRP).

and image quality. Significant superiority could be shown analogously to the unadjusted values to each method and each criterion with p -values between .002 – .019 with the exception of image quality for the DeepTaylor algorithm ($t(6) = 1.71, p = .131$).

The average rank of our method in direct comparison was 1.28, 1.36 and 1.52, for intuitive validity, semantic meaningfulness and image quality, compared with 2.78, 2.71, 3.23 (DeepSHAP/DeepTaylor/LRP) for intuitive validity, 2.75, 2.75, 3.41 for semantic meaningfulness and 3.04, 2.36, 3.07 for image quality. p -values ranged from $< .001$ – .019, again with the exception of image quality for the DeepTaylor algorithm ($p = .152$).

The pairwise Pearson correlations between the tested criteria over all questionnaires were between $r = .689$ (image quality and intuitive validity) and $r = .862$ (intuitive validity and semantic meaningfulness), implying an expected association between the tested criteria (see Fig. 8). The inter-observer reliability was $\rho = .639/.590/.557$ for intuitive validity, semantic meaningfulness and image

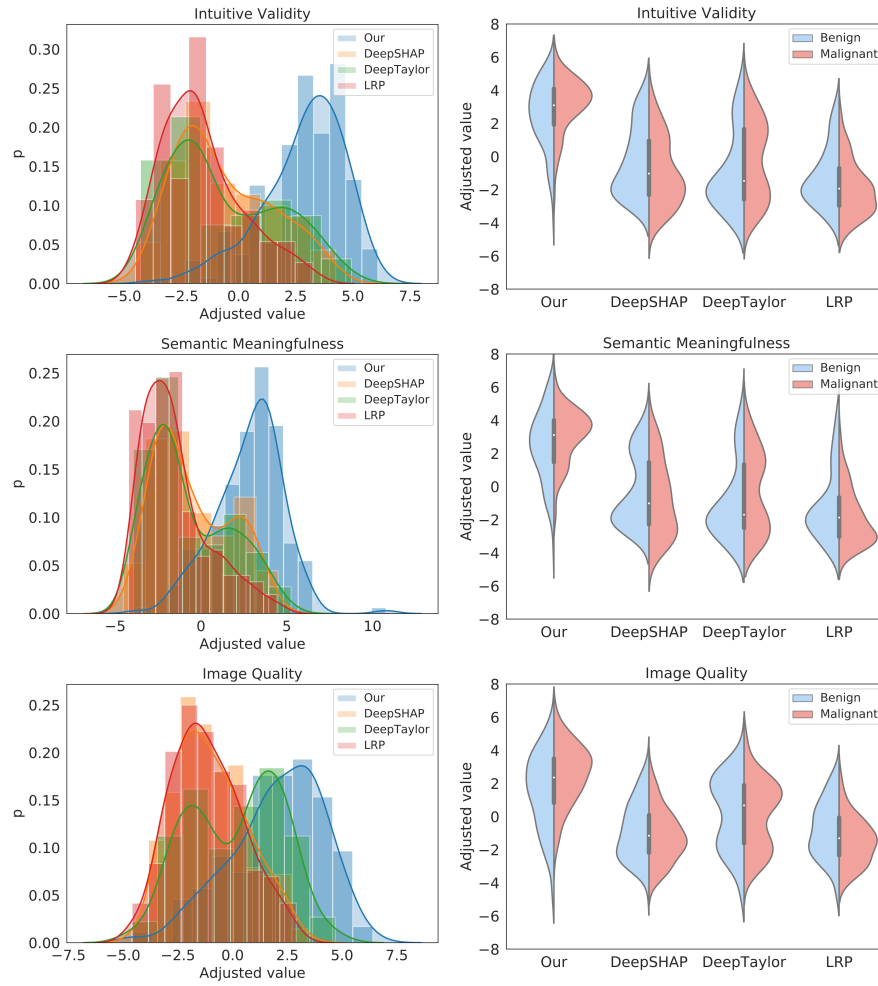


Figure 7: Comparison of the z-adjusted questionnaire results for each criterion and method cumulative (left) and by lesion type (right) for benign (blue) and malignant (red) lesions.

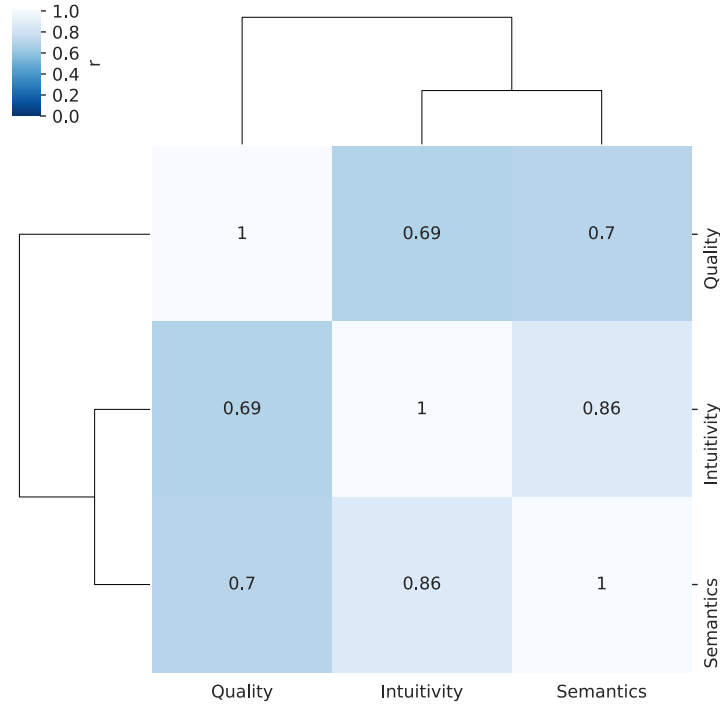


Figure 8: Dendrogram of the pairwise Pearson correlation matrix for the analyzed criteria.

quality, respectively, implying a substantial to moderate inter-observer agreement [50]. The Kaiser-Mayer-Olkin criterion for our sample was $KMO = .716$, implying the adequacy of our distribution for principal component analysis [49]. Applying the PCA, it was possible to extract a general factor of preferability Φ , which accounted for 84.6 % of the observed variance (see Fig. 9).

The differences between the questionnaire sheets A and B for each method and criterion were all non-significant with p -values between .133 and .957 (avg: .649, $.002 < t(6) < 1.68$), implying that there was no systematic difference between the questionnaires and no order effects occurred.

4.3. Discussion

Based on the results of the user study, our method was able to outperform DeepSHAP, DeepTaylor and LRP in terms of intuitive validity, semantic mean-

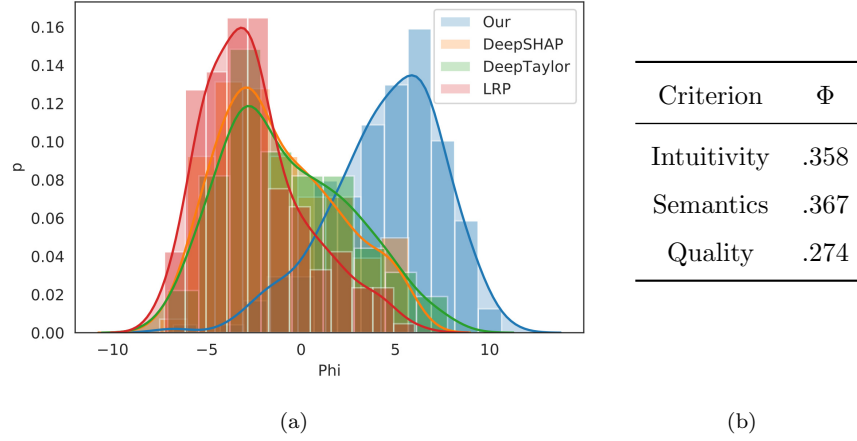


Figure 9: Comparison of the analyzed approaches (a) when extracting a general factor Φ using PCA and (b) its L1-normalized eigenvector. Φ accounts for 84.6 % of the observed variance.

ingfulness and image quality on the given data. While for DeepSHAP and LRP the results were significant for all tested criteria, the image quality could not be shown to be significantly different for the DeepTaylor algorithm due to the limited amount of degrees of freedom of the used statistical test. Nevertheless, a large effect size was observed for image quality, too. We demonstrated our method to be able to semantically modify the image with respect to the classification context. The user study indicates that the visualizations generated by our method emphasize regions which are perceived to be decision-relevant in the given images, while preserving a high image quality as well as intuitive validity. The extraction of a general factor of preferability Φ , as well as the inter-criteria correlations showed that the measured performance criteria are correlated and influence each other, expressing the need for future work on decision explanation to account for them. With an average rank of 1.28, 1.36 and 1.52 over all questions, our method markedly ranked best in comparison to the analyzed SoA methods, followed by DeepTaylor, DeepSHAP and LRP. The survey evaluation revealed no significant differences between the different questionnaire variants (cf. Sec. 3.4) and showed a substantial to moderate correlation between the raters, which suggests a reasonable questionnaire design.

5. Conclusion

Deep learning is an important tool for medical image analysis due to its ability to analyze vast amounts of data autonomously. Its low transparency and the resulting lack of understanding of its decision-making process, however, pose a major obstacle for its application in clinical routine. In medical imaging, it is particularly important to understand the inner workings of an algorithm with respect to issues such as algorithm validation, product quality, and finally liability. For this reason, there is a need for methods which allow clinicians as well as engineers to intuitively visualize the network decision process, i.e. to “see what the network sees.” As shown in Sec. 4, existing methods reach this aim only to a limited degree, and there is a lot of room for improvement.

With our method we were able to generate high-quality decision explanations for a trained classifier, which are both intuitively valid as well as semantically correct, and could clearly improve on the tested SoA approaches. As shown in Sec. 4, our method focuses on decision-relevant areas, and changes in these regions were demonstrated to be associated with changes in the classifier output probabilities. Regarding the analyzed criteria, it could be shown that intuitive validity, semantic meaningfulness and image quality go hand-in-hand with each other, and that therefore algorithms for decision explanation should account for all of them equally. With only 772 samples, the medical data set we used is significantly smaller than comparable computer vision benchmark data sets, such as the CIFAR-10/CIFAR-100 (60,000 samples) [51], MNIST (70,000 samples) [52] or cityscapes dataset (20,000 samples) [53]. Nevertheless, our method was able to successfully visualize network decisions, making it a candidate for explaining deep learning-based models for medical image analysis.

With our work, we contribute a significant step towards a better understanding of DNN-based classifier decisions, which in the future could help both engineers as well as clinical practitioners to better understand and hence develop medical algorithms more effectively, which might ultimately lead to an improved overall clinical acceptance.

5.1. Limitations

The quality of our method is determined by the quality of the trained CycleGAN. While recent work suggested approaches for dealing with non-convergence in GANs [54, 55], there is currently no well-established measure of convergence, especially regarding perceptual quality of decision explanations. While in our experiments the model did not appear to behave chaotically, it was not analyzed whether this is a side-effect of the model structure, or was specific to our problem. As our model is based on GAN training, creating a decision explanation network requires a high amount of computational resources, which becomes more relevant with increasing input image sizes. Thus, our approach is applicable to the explanation of fully-trained networks, yet it may be less suitable for immediate visualization and rapid prototyping. Additionally, future work should investigate the applicability to further problems, including tasks other than binary classification. While we could show significant improvements over the analyzed SoA methods, the average raw questionnaire results for our method were between 1.04 and 1.92 on a scale from -4 to 4, indicating that there is still room for improvement. Compared to -1.52 to -0.85 for the second ranked method DeepTaylor, our method was able to cover some of this potential. However, it is desirable that future research can further improve on that.

With only $N = 8$ participants, our user study was rather small. Each of our $N = 8$ participants evaluated 24 lesions for 3 criteria each, resulting in a total of 576 samples per algorithm or 2,304 answers in total. All statistical tests were done using a two-tailed t -test with $t(N - 2) = t(6)$. However, due to the large effect sizes, statistical significance could be shown for most comparisons. Nonetheless, based on the promising results of this study, a more comprehensive user survey in future work is highly desirable.

Vitae



Alexander Katzmann, is a research scientist at Siemens Healthineers. He holds a M.Sc. degree in computer science from Ilmenau, University of Technology. After researching in the field of neuroinformatics and cognitive robotics, his current work focuses on deep learning-based medical image analysis. He is based in the CT Image Analytics group in Forchheim, Germany, and currently working on his Ph. D. in collaboration with Ilmenau, University of Technology.



Dr. Oliver Taubmann is a research scientist with Siemens Healthineers. As a member of the CT Image Analytics group based in Forchheim, Germany, he develops clinical prototypes for learning-based image analysis to better support radiologists. Oliver holds a Ph. D. degree in computer science from FAU Erlangen-Nürnberg, for which he devised novel techniques for time-resolved cardiac image reconstruction from rotational angiography in the interventional suite.



Stephen Ahmad is a software developer at Siemens Healthineers. He holds a M.Sc. degree in engineering informatics and is currently based in the department of Digital Health Imaging Decision Support, Erlangen, Germany. He is focused on AI-related algorithm engineering and currently a developer for the Siemens Healthineers "AI-Rad Companion".



Dr. Michael Sühling is head of the CT Image Analytics group at the Siemens Healthineers CT headquarter in Forchheim, Germany. His research interests are mainly in the area of image analysis technologies for CT scan automation and clinical decision making support. He has been reviewer for the IEEE Transactions on Medical Imaging and the IEEE Transactions on Image Processing.



Horst-Michael Groß is full professor for Computer Science at Ilmenau University of Technology and head of the research lab (chair) for Neuroinformatics and Cognitive Robotics. From methodological view, his research is focused on real-time approaches to person detection, tracking, and re-identification in real-world video data streams and deep learning approaches for image-based object recognition, semantic segmentation, pose recognition, grasp pose estimation, and action recognition. Amongst others, he is a member of the IEEE and the International Neural Networks Society (INNS).

Disclaimer

The concepts and information presented in this article are based on research and are not commercially available.

References

- [1] N. Coudray, P. S. Ocampo, T. Sakellaropoulos, N. Narula, M. Snuderl, D. Fenyő, A. L. Moreira, N. Razavian, A. Tsirigos, Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning, *Nature medicine* 24 (10) (2018) 1559–1567.

- [2] N. Bayramoglu, J. Kannala, J. Heikkilä, Deep learning for magnification independent breast cancer histopathology image classification, in: 2016 23rd International conference on pattern recognition (ICPR), IEEE, 2016, pp. 2440–2445.
- [3] N. Ganapathy, R. Swaminathan, T. M. Deserno, Deep learning on 1-d biosignals: a taxonomy-based survey, *Yearbook of medical informatics* 27 (1) (2018) 98.
- [4] Y. Yuan, G. Xun, Q. Suo, K. Jia, A. Zhang, Wave2vec: Learning deep representations for biosignals, in: 2017 IEEE International Conference on Data Mining (ICDM), IEEE, 2017, pp. 1159–1164.
- [5] K.-L. Hua, C.-H. Hsu, S. C. Hidayati, W.-H. Cheng, Y.-J. Chen, Computer-aided classification of lung nodules on computed tomography images via deep learning technique, *OncoTargets and therapy* 8.
- [6] T. Würfl, F. C. Ghesu, V. Christlein, A. Maier, Deep learning computed tomography, in: *International conference on medical image computing and computer-assisted intervention*, Springer, 2016, pp. 432–440.
- [7] X. Dong, Y. Lei, T. Wang, K. Higgins, T. Liu, W. J. Curran, H. Mao, J. A. Nye, X. Yang, Deep learning-based attenuation correction in the absence of structural information for whole-body positron emission tomography imaging, *Physics in Medicine & Biology* 65 (5) (2020) 055011.
- [8] Y. Ding, J. H. Sohn, M. G. Kawczynski, H. Trivedi, R. Harnish, N. W. Jenkins, D. Lituiev, T. P. Copeland, M. S. Aboian, C. Mari Aparici, et al., A deep learning model to predict a diagnosis of alzheimer disease by using 18f-fdg pet of the brain, *Radiology* 290 (2) (2019) 456–464.
- [9] D. Shen, G. Wu, H.-I. Suk, Deep learning in medical image analysis, *Annual review of biomedical engineering* 19 (2017) 221–248.
- [10] S. Trebeschi, J. J. van Griethuysen, D. M. Lambregts, M. J. Lahaye, C. Parmar, F. C. Bakers, N. H. Peters, R. G. Beets-Tan, H. J. Aerts, Deep learn-

ing for fully-automated localization and segmentation of rectal cancer on multiparametric mr, *Scientific reports* 7 (1) (2017) 1–9.

- [11] F. Liu, H. Jang, R. Kijowski, T. Bradshaw, A. B. McMillan, Deep learning mr imaging-based attenuation correction for pet/mr imaging, *Radiology* 286 (2) (2018) 676–684.
- [12] E. Gibson, F. Giganti, Y. Hu, E. Bonmati, S. Bandula, K. Gurusamy, B. Davidson, S. P. Pereira, M. J. Clarkson, D. C. Barratt, Automatic multi-organ segmentation on abdominal ct with dense v-networks, *IEEE transactions on medical imaging* 37 (8) (2018) 1822–1834.
- [13] R. Chalapathy, S. Chawla, Deep learning for anomaly detection: A survey, *arXiv preprint arXiv:1901.03407*.
- [14] K. Yan, X. Wang, L. Lu, R. M. Summers, Deeplesion: automated mining of large-scale lesion annotations and universal lesion detection with deep learning, *Journal of Medical Imaging* 5 (3) (2018) 036501.
- [15] K. Kamnitsas, C. Ledig, V. F. Newcombe, J. P. Simpson, A. D. Kane, D. K. Menon, D. Rueckert, B. Glocker, Efficient multi-scale 3d cnn with fully connected crf for accurate brain lesion segmentation, *Medical image analysis* 36 (2017) 61–78.
- [16] A. Hosny, C. Parmar, T. P. Coroller, P. Grossmann, R. Zeleznik, A. Kumar, J. Bussink, R. J. Gillies, R. H. Mak, H. J. Aerts, Deep learning for lung cancer prognostication: A retrospective multi-cohort radiomics study, *PLoS medicine* 15 (11) (2018) e1002711.
- [17] X. Liu, L. Faes, A. U. Kale, S. K. Wagner, D. J. Fu, A. Bruynseels, T. Mahendiran, G. Moraes, M. Shamdas, C. Kern, et al., A comparison of deep learning performance against health-care professionals in detecting diseases from medical imaging: a systematic review and meta-analysis, *The lancet digital health* 1 (6) (2019) e271–e297.

- [18] V. Kumar, Y. Gu, S. Basu, A. Berglund, S. A. Eschrich, M. B. Schabath, K. Forster, H. J. Aerts, A. Dekker, D. Fenstermacher, et al., Radiomics: the process and the challenges, *Magnetic resonance imaging* 30 (9) (2012) 1234–1248.
- [19] H. J. Aerts, E. R. Velazquez, R. T. Leijenaar, C. Parmar, P. Grossmann, S. Carvalho, J. Bussink, R. Monshouwer, B. Haibe-Kains, D. Rietveld, et al., Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach, *Nature communications* 5 (1) (2014) 1–9.
- [20] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. Van Der Laak, B. Van Ginneken, C. I. Sánchez, A survey on deep learning in medical image analysis, *Medical image analysis* 42 (2017) 60–88.
- [21] A. Katzmann, A. Muehlberg, M. Sühling, D. Noerenberg, J. W. Holch, V. Heinemann, H.-M. Groß, Predicting lesion growth and patient survival in colorectal cancer patients using deep neural networks.
- [22] H. Noh, S. Hong, B. Han, Learning deconvolution network for semantic segmentation, in: *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1520–1528.
- [23] B. Zhou, A. Khosla, L. A., A. Oliva, A. Torralba, Learning Deep Features for Discriminative Localization., *CVPR*.
- [24] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra, Grad-cam: Visual explanations from deep networks via gradient-based localization, in: *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [25] A. Chattopadhyay, A. Sarkar, P. Howlader, V. N. Balasubramanian, Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks, in: *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE, 2018, pp. 839–847.

- [26] J. Schlemper, O. Oktay, M. Schaap, M. Heinrich, B. Kainz, B. Glocker, D. Rueckert, Attention gated networks: Learning to leverage salient regions in medical images, *Medical image analysis* 53 (2019) 197–207.
- [27] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, *PloS one* 10 (7).
- [28] A. Shrikumar, P. Greenside, A. Kundaje, Learning important features through propagating activation differences, in: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, JMLR. org, 2017, pp. 3145–3153.
- [29] M. T. Ribeiro, S. Singh, C. Guestrin, ” why should i trust you?” explaining the predictions of any classifier, in: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1135–1144.
- [30] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: *Advances in neural information processing systems*, 2017, pp. 4765–4774.
- [31] D. Erhan, Y. Bengio, A. Courville, P. Vincent, Visualizing higher-layer features of a deep network, *University of Montreal* 1341 (3) (2009) 1.
- [32] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, J. Clune, Synthesizing the preferred inputs for neurons in neural networks via deep generator networks, in: *Advances in neural information processing systems*, 2016, pp. 3387–3395.
- [33] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: *Advances in neural information processing systems*, 2014, pp. 2672–2680.

- [34] A. Nguyen, J. Yosinski, J. Clune, Understanding neural networks via feature visualization: A survey, in: *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, Springer, 2019, pp. 55–76.
- [35] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, in: *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [36] P. Isola, J.-Y. Zhu, T. Zhou, A. A. Efros, Image-to-image translation with conditional adversarial networks, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.
- [37] Z. Wang, E. P. Simoncelli, A. C. Bovik, Multiscale structural similarity for image quality assessment, in: *The Thrity-Seventh Asilomar Conference on Signals, Systems & Computers*, 2003, Vol. 2, Ieee, 2003, pp. 1398–1402.
- [38] T. Park, M.-Y. Liu, T.-C. Wang, J.-Y. Zhu, Semantic image synthesis with spatially-adaptive normalization, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2337–2346.
- [39] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, B. Catanzaro, High-resolution image synthesis and semantic manipulation with conditional gans, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8798–8807.
- [40] S. G. Armato III, G. McLennan, L. Bidaut, M. F. McNitt-Gray, C. R. Meyer, A. P. Reeves, L. P. Clarke, et al., Data from lidc-idri. the cancer imaging archive, DOI [http://doi.org/10.7937/K9\(7\)](http://doi.org/10.7937/K9(7)).
- [41] S. G. Armato III, G. McLennan, L. Bidaut, M. F. McNitt-Gray, C. R. Meyer, A. P. Reeves, B. Zhao, D. R. Aberle, C. I. Henschke, E. A. Hoffman, et al., The lung image database consortium (lidc) and image database resource initiative (idri): a completed reference database of lung nodules on ct scans, *Medical physics* 38 (2) (2011) 915–931.

- [42] K. Clark, B. Vendt, K. Smith, J. Freymann, J. Kirby, P. Koppel, S. Moore, S. Phillips, D. Maffitt, M. Pringle, et al., The cancer imaging archive (tcia): maintaining and operating a public information repository, *Journal of digital imaging* 26 (6) (2013) 1045–1057.
- [43] A. Nibali, Z. He, D. Wollersheim, Pulmonary nodule classification with deep residual networks, *International journal of computer assisted radiology and surgery* 12 (10) (2017) 1799–1808.
- [44] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: *International Conference on Medical image computing and computer-assisted intervention*, Springer, 2015, pp. 234–241.
- [45] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, K.-R. Müller, Explaining nonlinear classification decisions with deep taylor decomposition, *Pattern Recognition* 65 (2017) 211–222.
- [46] M. Alber, S. Lapuschkin, P. Seegerer, M. Hägele, K. T. Schütt, G. Montavon, W. Samek, K.-R. Müller, S. Dähne, P.-J. Kindermans, investigate neural networks, *Journal of Machine Learning Research* 20 (93) (2019) 1–8.
- [47] J. A. Krosnick, Questionnaire design, in: *The Palgrave handbook of survey research*, Springer, 2018, pp. 439–455.
- [48] H. F. Kaiser, A second generation little jiffy, *Psychometrika* 35 (4) (1970) 401–415.
- [49] C. D. Dziuban, E. C. Shirkey, When is a correlation matrix appropriate for factor analysis? some decision rules., *Psychological bulletin* 81 (6) (1974) 358.
- [50] J. R. Landis, G. G. Koch, The measurement of observer agreement for categorical data, *biometrics* (1977) 159–174.

- [51] A. Krizhevsky, G. Hinton, et al., Learning multiple layers of features from tiny images.
- [52] Y. LeCun, C. Cortes, C. Burges, Mnist handwritten digit database, ATT Labs [Online]. Available: <http://yann.lecun.com/exdb/mnist> 2.
- [53] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, B. Schiele, The cityscapes dataset for semantic urban scene understanding, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 3213–3223.
- [54] D. Berthelot, T. Schumm, L. Metz, Began: Boundary equilibrium generative adversarial networks, arXiv preprint arXiv:1703.10717.
- [55] S. Liu, O. Bousquet, K. Chaudhuri, Approximation and convergence properties of generative adversarial learning, in: Advances in Neural Information Processing Systems, 2017, pp. 5545–5553.