

Uncovering convolutional neural network decisions for diagnosing multiple sclerosis on conventional MRI using layer-wise relevance propagation

Fabian Eitel^{a,b}, Emily Soehler^{a,b}, Judith Bellmann-Strobl^{d,e}, Alexander U. Brandt^{c,d,h}, Klemens Ruprecht^c, René M. Giess^{c,d}, Joseph Kuchling^{c,d,e}, Susanna Asseyer^{c,d,e}, Martin Weygandt^{c,d}, John-Dylan Haynes^{b,g}, Michael Scheel^{c,d,f}, Friedemann Paul^{c,d,e,g,1}, Kerstin Ritter^{a,b,1,*}

^aCharité – Universitätsmedizin Berlin, corporate member of Freie Universität Berlin, Humboldt-Universität zu Berlin, and Berlin Institute of Health (BIH); Department of Psychiatry and Psychotherapy; 10117 Berlin, Germany.

^bCharité – Universitätsmedizin Berlin, corporate member of Freie Universität Berlin, Humboldt-Universität zu Berlin, and Berlin Institute of Health (BIH); Berlin Center for Advanced Neuroimaging, Bernstein Center for Computational Neuroscience; 10117 Berlin, Germany.

^cCharité – Universitätsmedizin Berlin, corporate member of Freie Universität Berlin, Humboldt-Universität zu Berlin, and Berlin Institute of Health (BIH); Department of Neurology; 10117 Berlin, Germany.

^dCharité – Universitätsmedizin Berlin, corporate member of Freie Universität Berlin, Humboldt-Universität zu Berlin, and Berlin Institute of Health (BIH); NeuroCure Clinical Research Center; 10117 Berlin, Germany.

^eCharité – Universitätsmedizin Berlin, corporate member of Freie Universität Berlin, Humboldt-Universität zu Berlin, and Berlin Institute of Health (BIH); Experimental and Clinical Research Center, Max Delbrück Center for Molecular Medicine; 10117 Berlin, Germany.

^fCharité – Universitätsmedizin Berlin, corporate member of Freie Universität Berlin, Humboldt-Universität zu Berlin, and Berlin Institute of Health (BIH); Department of Neuroradiology; 10117 Berlin, Germany.

^gEinstein Center for Digital Future Berlin, Germany.

^hDepartment of Neurology, University of California, Irvine, California, USA.

arXiv:1904.08771v1 [cs.CV] 18 Apr 2019

Abstract

Machine learning-based imaging diagnostics has recently reached or even superseded the level of clinical experts in several clinical domains. However, classification decisions of a trained machine learning system are typically non-transparent, a major hindrance for clinical integration, error tracking or knowledge discovery. In this study, we present a transparent deep learning framework relying on convolutional neural networks (CNNs) and layer-wise relevance propagation (LRP) for diagnosing multiple sclerosis (MS), the most widespread autoimmune neuroinflammatory disease. MS is commonly diagnosed utilizing a combination of clinical presentation and conventional magnetic resonance imaging (MRI), specifically the occurrence and presentation of white matter lesions in T2-weighted images. We hypothesized that using LRP in a naive predictive model would enable us to uncover relevant image features that a trained CNN uses for decision-making. Since imaging markers in MS are well-established this would enable us to validate the respective CNN model. First, we pre-trained a CNN on MRI data from the Alzheimer’s Disease Neuroimaging Initiative ($n = 921$), afterwards specializing the CNN to discriminate between MS patients and healthy controls ($n = 147$). Using LRP, we then produced a heatmap for each subject in the holdout set depicting the voxel-wise relevance for a particular classification decision. The resulting CNN model resulted in a balanced accuracy of 87.04% and an area under the curve of 96.08% in a receiver operating characteristic curve. The subsequent LRP visualization revealed that the CNN model focuses indeed on individual lesions, but also incorporates additional information such as lesion location, non-lesional white matter or gray matter areas such as the thalamus, which are established conventional and advanced MRI markers in MS. We conclude that LRP and the proposed framework have the capability to make diagnostic decisions of CNN models transparent, which could serve to justify classification decisions for clinical review, verify diagnosis-relevant features and potentially gather new disease knowledge.

Keywords: convolutional neural networks, deep learning, multiple sclerosis, MRI, layer-wise relevance propagation, visualization, transfer learning

1. Introduction

Multiple Sclerosis (MS) is the most widespread autoimmune neuroinflammatory disease in young adults with 2.2 million cases reported worldwide [1]. The disease is mainly characterized by inflammation, demyelination and neurodegeneration in the central nervous system and often leads to substantial dis-

ability in patients [2]. The current quasi-standard for diagnosing MS, the McDonald criteria, relies on clinical presentation and the presence of lesions visible in conventional T2-weighted brain magnetic resonance imaging (MRI) data [3]. Most common are fluid-suppressed T2-weighted image sequences, which are sensitive towards MS-relevant white matter lesions, but also relatively unspecific with respect to underlying disease processes [4]. Several other imaging markers have been described including neurodegeneration, thalamic atrophy, cortical lesions, altered functional connectivity or central vein signs

*Corresponding author: kerstin.ritter@charite.de

¹Shared authorship

[5, 6, 7, 8, 9, 10], of which some are captured in conventional MRI and others require advanced MRI techniques such as diffusion weighted imaging or functional MRI.

In the last decade, a lot of research effort has been put on the automatic (i.e. data-driven) detection of neurological diseases based on neuroimaging data including MRI [11, 12]. Early approaches combined parameter-based machine learning algorithms, such as support vector machines, with carefully extracted features known or hypothesized to be relevant in the respective disease. In MS research, features ranging from T2 lesion characteristics to local intensity patterns or multi-scale information extracted from MRI data have been used in combination with standard machine learning analyses to either diagnose MS or predict disease progression [13, 14, 15, 16, 17]. While choosing features based on expert criteria reflects the current state of knowledge, it does not allow for finding new and potentially unexpected hidden data properties, which might also help in characterizing a certain disease. Deep learning techniques fill a gap here and allow for utilizing hierarchical information directly from raw or minimally processed data [18]. By being specifically tailored to image data, in particular convolutional neural networks (CNNs) have led to major breakthroughs in medical imaging [19, 20, 21, 22]. In neuroimaging, most CNN analyses so far focused on Alzheimer’s disease [23], but there are also some recent studies in MS. Given the importance of lesions in diagnosing MS and monitoring disease progression, most efforts have been put on the task of lesion segmentation [24, 25, 26]. Others used CNNs to diagnose MS based on 2-dimensional MRI slices [27] or to predict short-term disease activity based on binary lesion masks [28].

Despite their potential, deep learning methods are criticized for being non-transparent (such as a ‘black box’) due to the difficulty to retrace the classification decision in light of huge parameter spaces and highly non-linear interactions [29]. This is especially problematic in medical applications since understanding and explaining neural network decisions is required for clinical integration, error tracking or knowledge discovery. Explaining neural network decisions is an open research area in computer science and a number of suggestions have been made in recent years. Different directions for explanations include visualizing features [30], generating images that maximally activate a certain neuron [31] and creating heatmaps based on the input images indicating the relevance of each voxel for the final classification decision [32, 33, 34]. Heatmaps are in particular valuable in the medical context, since they allow for an easy and intuitive investigation of what the respective classifier found to be important directly in the input data. Besides understanding diagnostic decisions for individual patients, heatmaps might be useful in validating CNN models. Recently, we have shown the potential of transparent CNN applications for knowledge discovery in Alzheimer’s disease [35, 36].

The objective of the current study was to investigate whether a transparency approach can uncover decision processes in MRI-based diagnosis of MS, a disease with well-defined imaging markers, thereby supporting future clinical implementation and verification of machine learning-based diagnosis systems. We present a transparent CNN framework for the MRI-based

diagnosis of clinically definite MS relying on layer-wise relevance propagation (LRP, [33, 37]) – a heatmap method that has been shown to outperform previous approaches in terms of explainability and disease-specific evidence [37, 36]. Since the data set was rather small ($n = 147$), we investigated the effect of pre-training the CNN on data from the Alzheimer’s Disease Neuroimaging Initiative (ADNI, $n = 921$). Using LRP, individual heatmaps were generated for each subject and analyzed with respect to well-established imaging features in MS (e.g. white matter lesions or thalamus atrophy). By showing that LRP in combination with a naive CNN model (i.e. a model independent of MS-specific knowledge) indeed helps in uncovering relevant imaging features, we conclude that this framework is not only useful in justifying individual diagnostic decisions but also to validate CNN models (especially in light of small sample sizes).

2. Materials and methods

2.1. Subjects

In the present study, we retrospectively analyzed data collected by FP from Charité – Universitätsmedizin Berlin as part of the VIMS study: Follow-up examination of visual parameters for the creation of a database (neuro-ophthalmologic register) in patients with MS versus healthy subjects.² We enrolled 76 patients with clinically definite multiple sclerosis (MS) according to the McDonald criteria 2010 [38] and 71 healthy controls. Patients were excluded if they were outside the age range of 18 – 69 or did not have an MRI scan. All patients were examined under supervision of a board-certified neurologist at the NeuroCure Clinical Research Center (Charité – Universitätsmedizin Berlin) between January 2011 and July 2015. All participants provided written informed consent prior to their inclusion in the study. The study was approved by the local ethics committee and was performed in accordance with the 1964 Declaration of Helsinki in its currently applicable version. Part of this data has been used in previous studies (e.g. [39]). Demographical details of subjects can be found in Table 1.

2.2. MRI acquisition and preprocessing

All MRI data were acquired on the same 3 T scanner (Tim Trio Siemens, Erlangen, Germany) using a volumetric high-resolution T1 weighted magnetization prepared rapid acquisition gradient echo (MPRAGE) sequence (TR = 1900 ms, TE = 2.55 ms, TI = 900 ms, FOV=240x240 mm², matrix 240x240, 176 slices, slice thickness 1 mm) as well as a volumetric high-resolution fluid-attenuated inversion recovery sequence (FLAIR, TR = 6000 ms, TE = 388 ms, TI = 2100 ms; FOV=256x256 mm², slice thickness 1 mm). Individual lesion masks were generated based on FLAIR images by three expert raters under the supervision of a board-certified radiologist using ITK-SNAP³ [40]. The MRI data were preprocessed using

²<https://neurocure.de/en/clinical-center/clinical-studies/current-studies.html>

³www.itksnap.org

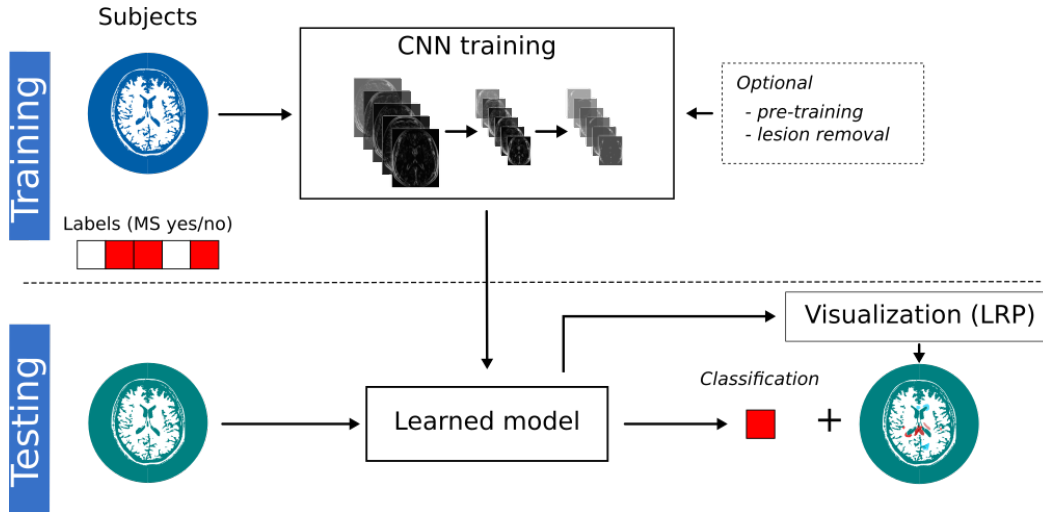


Figure 1: Illustration of the transparent CNN framework. In the training phase, the CNN model learns a non-linear relationship between the MRI data and the binary diagnostic labels (MS yes/no). Optionally, the CNN models are pre-trained on a substitute data set or lesions are filled in the MRI data. The learned CNN model is then tested on new subjects to predict the diagnostic label. By supplementing this label with a LRP heatmap, which indicates the relevance of each voxel for the respective label, this framework allows us to understand (at least to some extent) the classification decision in individual subjects. Additionally, the validity of the CNN models can be assessed by matching highlighted brain areas with domain knowledge.

	MS patients	Healthy controls
Subjects [n]	76	71
Female/Male , in %	55 % / 45 %	65 % / 35 %
Age , mean \pm std	43.32 (\pm 11.99)	38.23 (\pm 13.10)
Disease duration , mean \pm std	149.65 (\pm 123.35)	n.a.
EDSS , median, range	2.50 (0.00 - 6.50)	n.a.
Lesion volume , mean \pm std	7.28 (\pm 8.09)	0.57 (\pm 1.94)

Table 1: Demographics of MS patients and healthy controls. Disease duration is measured in months and lesion volume in ml. EDSS, expanded disability status scale; std, standard deviation.

a customized pipeline based on the software programs statistical parametric mapping (SPM) [41], Advanced Normalization Tools (ANTs) [42] and FMRIB Software Library (FSL) [43]. After a bias-field correction and field of view-cropping, MPRAGE images were linearly registered to the Montreal Neurological Institute (MNI)-template. FLAIR images were coregistered to the native MPRAGE images and then transformed to the MNI space using the transformation matrix estimated on the the MPRAGE images. Normalization parameters were estimated based on MPRAGE images because they provide a better tissue contrast than FLAIR. For each person, the skull was extracted using the Brain Extraction Tool (BET) of FSL. Please note that only FLAIR data entered the subsequent analyses and that they were preprocessed in that way to ensure that images are in relative realignment while preserving individual structural variations. For computational efficiency initial scan volumes (182x218x182) were down-sampled to 96x114x96 voxels and standardized for each subject using min-max scaling. To analyze the impact of white matter lesions, we generated an additional MRI data set, in which the lesions were filled. For this, we implemented a version of [44], in which lesion areas (according to the manually segmented lesion masks) have been replaced by local average intensities in normal-appearing white

matter. White matter maps were obtained from the SPM 12 tissue segmentation algorithm [45].

2.3. ADNI data for pre-training

Data used for pre-training were obtained from the Alzheimers Disease Neuroimaging Initiative (ADNI) database⁴. We have used subjects from ADNI phase 1 who were included in one of two standard MRI collections [46]. We only selected MRI data of Alzheimer’s disease patients and cognitive normal subjects, in total 921 MRI scans from 276 subjects (covering one to three time points). The MRI scans were acquired with 1.5 Tesla scanners at multiple sites and had already undergone gradient non-linearity, intensity inhomogeneity and phantom-based distortion correction. T1-weighted MPRAGE scans were downloaded and warped to MNI space with ANTs [42]. As for the MS data, the initial scan volumes were down-sampled to 96x114x96 voxels and standardized.

2.4. Classification and visualization analyses

Based on the preprocessed FLAIR data, we first trained several convolutional neural network (CNN) models (with and

⁴<http://adni.loni.usc.edu>, RRID:SCR.003007

without pre-training, with and without lesion-filling) to discriminate MS patients and healthy controls and then explained the model’s decisions for individual subjects in the test data using layer-wise relevance propagation (LRP). For the CNN models, we evaluated the effect of transfer learning by (1) training the model solely on MS data and (2) pre-training the model on ADNI data and fine-tuning it on MS data. To examine whether our pre-trained network can also learn from only normal-appearing brain matter (NABM), i.e. regions without hyper-intense lesions, we retrained the network on lesion-filled MS data. As baseline analyses, we included a support vector machine to classify based on (1) lesion volume and (2) preprocessed FLAIR data. Prior to training, the MS data set was randomly split into two sets: (1) a set for training and hyperparameter optimization (85 %) and (2) a holdout set used only for final model evaluation (15 %). The code for all models and also the lesion filling algorithm is available at <https://github.com/derEitel/explainableMS>. In the following subsections, we specify our parameter settings for CNNs, transfer learning and visualization techniques (in particular LRP).

2.4.1. Convolutional neural networks

In this study, we used a convolutional neural network (CNN) architecture consisting of four convolutional layers followed by exponential linear units (ELUs) activation functions and four max-pooling layers applied after the first, second and fourth ELU activation. For each convolutional layer, we learned 64 filters with a kernel size of 3x3x3. Finally, a linear layer with an output shape of 1 and a sigmoid activation returns the classification score. The rationale behind this architecture was mainly to avoid overfitting and therefore has a comparably low number of trainable parameters (namely 333,889). To improve generalization further, the model has been regularized using a dropout on the outputs of each max-pooling layer ($p = 0.3$), L2-regularization ($\lambda = 0.01$) using the weights of the third and fourth convolutional layer, and finally early-stopping the training after the validation loss has not improved for 15 epochs. Additionally, the data was augmented during training by flipping it along the sagittal axis with a probability of 50% and randomly translated between -2 and 2 pixels along the sagittal axis. We trained all models using the Adam optimizer [47]. Hyperparameters were optimized on 85% of the training data, leaving 15% for validation. After finding suitable hyperparameters, the model performance was tested out-of-sample on the holdout set. All CNN experiments were repeated 10 times, and thus reported metrics are an average over all 10 trials. We report balanced accuracy as a mean between sensitivity and specificity as well as area under the receiver operating characteristic curve (AUC). All code was implemented using Keras [48] with the TensorFlow [49] backend.⁵

2.4.2. Transfer learning

Due to the small sample size of the MS data set, we employed the principle of transfer learning [50, 51, 52], which

has been shown to improve performance in medical imaging including MRI data [53, 54, 55, 56, 57]. We pre-trained our CNN model on ADNI data to separate Alzheimer patients and healthy controls, and fine-tuned it on the MS data set to separate MS patients and healthy controls. Since the ADNI data set contains multiple scans for several subjects we ensured that validation and testing was done on disjoint subject sets. The average balanced accuracy over all trials was 78.47 %. For further analysis, we selected a model from the 10 trials based on its performance, and then picked its training checkpoint with the best validation accuracy of 82.50 %. Fine-tuning on the MS data set uses the same model architecture, which is initialized with the weights and biases of the selected pre-trained model instead of randomly distributed values. Please note that we here transferred a CNN model (1) across diseases (Alzheimer’s disease to MS) and (2) across MRI sequences (MPRAGE to FLAIR).

2.4.3. Visualization

Deep learning methods are often criticized for their lack of interpretability [29, 58, 59]. In contrast to the decision nodes of a decision tree, which can be easily followed, and standardized coefficients in regression analysis which can determine feature importance, the learned weights of a CNN and its non-linear combinations are harder to comprehend. Over the last years much research has focused on improving the interpretability of neural networks. While some work has focused on understanding class representations and functions of individual neurons, others have developed methods to generate heatmaps based on the input data that indicate the importance or relevance of each pixel or voxel for the final classification decision [33, 60, 34]. The latter approach is in particular promising in the medical field since it allows for explaining in a fast and intuitive way individual classification decisions without the need for delving deeply into the network structure [36]. However, heatmaps can be computed in different ways and it is important to understand their underlying principles. The most popular approach is sensitivity analysis [60], in which the norm $\left\| \cdot \right\|_{l_p}$ over the gradient of the classification score $f(x)$ with respect to each pixel or voxel x_i is calculated:

$$R_i = \left\| \frac{\partial}{\partial x_i} f(x) \right\|_{l_p} \quad (1)$$

This results in image-specific saliency maps attributing higher scores to voxels, for which small variations largely affect the classification score. There exist different variations and modifications of sensitivity analysis (e.g. [34]). However, sensitivity analyses and its variants are very susceptible to highlight regions which would change the classification score drastically if altered but might not be relevant to a human expert. Examples include pixels in which the main object is hidden by another object or possible artifacts in the data set. Furthermore sensitivity analysis highlights regions regardless of their importance for or against a class but for classification in general. When applied in neuroimaging this could cause the visualization to emphasize brain regions which have a general relevance for the given task

⁵Keras version 2.2; TensorFlow version 1.11

(e.g. hippocampus in classifying Alzheimer’s disease patients) but not an individual directed relevance for the specific subject [36].

Recently, layer-wise relevance propagation (LRP) has been introduced as an alternative method for producing heatmaps and shown to be superior to sensitivity analysis [33, 37, 59]. LRP uses the classification score $f(x)$ directly (and not the gradient) and propagates it through the network using the following rule

$$R_i = \sum_j \frac{x_i w_{ij}}{\sum_i x_i w_{ij} + \epsilon} R_j. \quad (2)$$

Here, the relevance from layer R_j is propagated to its previous layer R_i . The term ϵ can be set to a small value (in this study: 0.001) to avoid division by 0. By using both the activation x_i as well as the weights w_{ij} connecting layers i and j , LRP assigns a larger share to neurons that are more strongly activated and to connections which have been reinforced during training [61]. By decomposing the classification score $f(x)$ rather than the gradient and conserving the classification score during back-propagation, LRP overcomes the flaws of sensitivity analysis [61] and has been shown to provide evidence for Alzheimer’s disease in individual subjects [36].

In this study, we produced individual LRP heatmaps for every subject in the holdout set. We have used the iNNvestigate implementation of LRP [62].⁶ Besides qualitatively comparing individual heatmaps, we compared average heatmaps of MS patients and healthy controls. We evaluated the importance of different brain regions by computing the average relevance for each brain area in the (1) Neuromorphometrics atlas⁷ [63] and the (2) JHU DTI-based white-matter atlas⁸ [64]. To evaluate the effect of transfer learning on the heatmaps, we compare average heatmaps for MS patients before and after pre-training. To assess the relevance of normal-appearing brain areas in contrast to lesion areas, we computed relevance scores separately for the original MRI data set and the lesion filled MRI data set.

2.4.4. Baseline analyses

As a baseline we have trained a support vector machine (SVM) to classify between MS patients and healthy controls based on T2 lesion load. Hyperparameters were tuned on the training data set using grid search, nested within a 5-fold cross-validation (SVM kernel: linear and radial basis function [RBF], $C, \gamma = [0.001, 0.1, 1, 10]$). For completeness we have also trained a SVM on the preprocessed FLAIR images with the same cross-validation strategy, optionally together with a prior dimensionality reduction via principal component analysis.

3. Results

3.1. Classification performance

In Table 2, we depict the performance for the different classification models. As expected T2 lesion load – as one of the

core biomarkers in MS – in combination with a SVM led to a high balanced accuracy (88.46%) and a high area under the curve (AUC) of the receiver operating characteristic (94.62%). When instead of the T2 lesion load the entire FLAIR image is used as input to the SVM, the AUC dropped down to 66.92%. The CNN model solely trained on the MS data set resulted in a balanced accuracy of 71.23% and an AUC of 85.46%. When the network has been pre-trained on the ADNI data set and fine-tuned to the MS data set, the balanced accuracy increased by 16 percentage points to 87.04% and is therefore comparable to the performance of the baseline T2 lesion load model. Moreover, the pre-trained CNN model outperformed all other classifiers in terms of AUC (96.08%) and importantly also in terms of sensitivity (93.08%). The ROC curve for all 10 trials is shown in supplementary Figure 1. To assess the impact of normal-appearing brain matter, we trained the same CNN model on lesion-filled FLAIR data. Still, a reasonable balanced accuracy of 70.15 % and a relatively high AUC of 90.92 % has been achieved.

3.2. Visualization

After the CNN models have been trained, we used LRP to generate an individual heatmap for each subject in the hold-out data set indicating the relevance of each voxel for the respective classification decision. Unless otherwise stated, we restricted our analyses only to correctly classified examples (i.e. true positives and true negatives). In Figure 2, we show the individual heatmaps overlayed on the FLAIR data for four MS patients, who achieved the highest classification scores. High classification scores generally indicate a higher confidence of the model for the respective classification decision and thus the corresponding explanations are usually more pronounced and less diffuse as for cases with lower classification scores. All four patients have in common that high positive relevance is attributed around the occipital horn of both lateral ventricles and covers periventricular lesion areas as well as the body and splenium of the corpus callosum. Even though the images were clearly classified as MS, certain regions are assigned negative relevance, meaning that these areas speak against the MS diagnosis. Negative relevance can be found around the frontal horn of both ventricles, notably even in periventricular lesion areas (see for example subject 1). Interestingly, lesions not bordering the ventricles seem often to be ignored or are assigned negative relevance. Please note that – in contrast to classical multi-univariate studies – the CNN model in combination with LRP is invariant to translations and thus is capable of identifying specific tissue areas (e.g. white matter lesions) having a positive relevance for MS regardless of their position in the voxel space.

In Figure 3, we show average heatmaps for all correctly classified MS patients (top) and all correctly classified healthy controls (bottom) in the holdout set. In accordance with the heatmaps of the individual subjects in Figure 2, posterior periventricular white matter regions have a strong positive relevance for the MS diagnosis. This is true for both MS patients and healthy controls, but the effect is less pronounced for healthy controls. The reversed effect can be seen for clusters exhibiting negative relevance in white matter areas in the corpus callosum and close to occipital and parietal lobe. Over

⁶The implementation can be found at <https://github.com/albermax/investigate>

⁷Contained in the SPM12 software, <https://www.fil.ion.ucl.ac.uk/spm/software/spm12/>

⁸<https://fsl.fmrib.ox.ac.uk/fsl/fslwiki/Atlases>

Data	Pre-train.	Class.	Bal. acc.	Sens.	Spec.	AUC
T2 lesion load	-	SVM	88.46 %	76.92 %	100.00 %	94.62 %
FLAIR	-	SVM	66.92 %	53.85 %	80.00 %	66.92 %
FLAIR	no	CNN	71.23 %	68.46 %	74.00 %	85.46 %
FLAIR	yes	CNN	87.04 %	93.08 %	81.00 %	96.08 %
FLAIR - les. fill.	yes	CNN	70.15 %	92.31 %	48.00 %	90.92 %

Table 2: Performance (in %) for the different models on the holdout data set. Pre-train., pre-training; Class., classifier; Bal. acc., balanced accuracy; Sens., sensitivity; Spec., specificity; AUC, area under the curve of the receiver operating characteristic; les. fill., lesions filled.

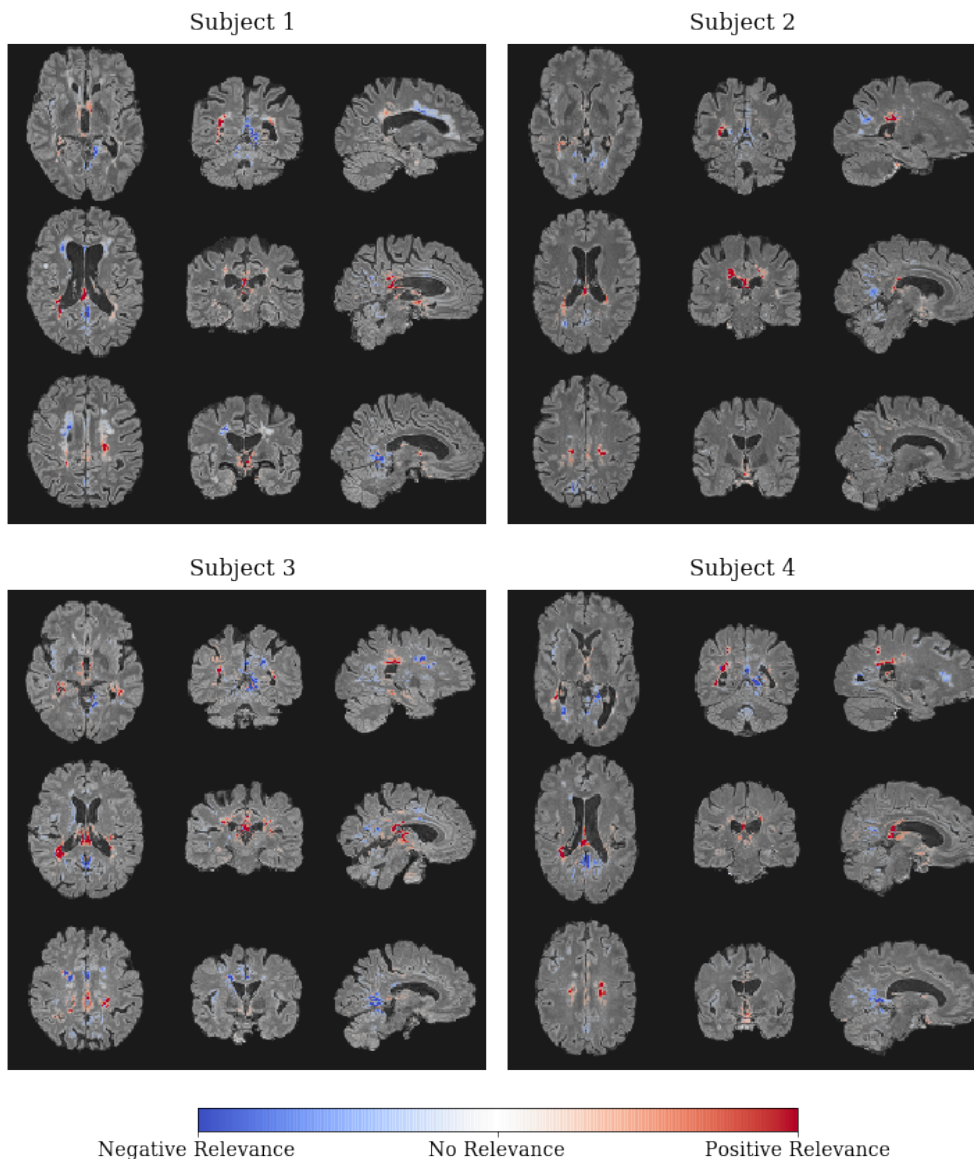


Figure 2: Individual LRP heatmaps (overlaid on the input FLAIR data) for the four MS patients with the highest classification score. Heatmap values are normalized in the range $[-0.03, 0.03]$. Colors indicate regions supporting (red) or rejecting (blue) the classification as a MS patient with respect to the underlying CNN model.

all voxels healthy controls typically obtain a negative relevance sum (mean±std: $-1.05e-6 \pm 0.0013$) as opposed to a positive relevance sum in MS patients ($3.07e-06 \pm 0.0014$).

To analyze the LRP heatmaps quantitatively with respect to different brain areas, we computed the relevance sum in (1) the Neuromorphometrics atlas mostly containing gray matter regions and (2) the JHU ICBM-DTI atlas containing white matter regions. Areas were aggregated between left and right hemisphere and certain substructures are combined into one region. For visualization of (1) we selected the 30 areas with the highest sum of absolute relevance means across MS patients and healthy controls in the test set, yielding areas with both the highest and lowest relevance. Please reconsider here that the MRI data have only been linearly registered and thus slight deviations from the anatomical locations stated in the atlases are conceivable. In Figure 4, we depict the region-wise LRP relevance for MS diagnosis, separately for MS patients and healthy controls. In the Neuromorphometrics atlas (see Figure 4a), most relevance is attributed to cerebral white matter, followed by thalamus, lateral ventricles and diencephalon. Negative relevance is strongest in the precuneus, followed by lingual gyrus, cuneus and insula. In the JHU white matter atlas (see Figure 4b), most positive relevance is attributed to posterior corona radiata and corpus callosum, followed by posterior thalamic radiation, tapetum, internal capsule and fornix. Notably, these areas are generally characterized by a high lesion density, which is also present in this MS data set (see supplementary Figures 2 and 3). Negative relevance has been found in the superior and anterior corona radiata. Generally, the relevance for MS patients is higher in white matter than in gray matter areas. Moreover, the differences between MS patients and healthy controls are more pronounced in white matter areas.

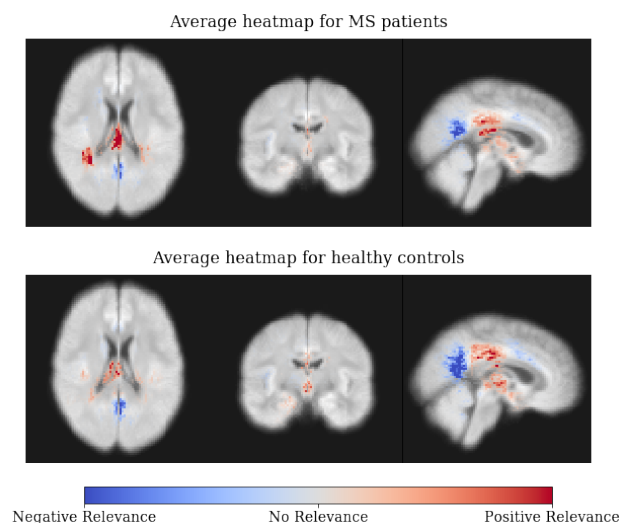


Figure 3: Average LRP heatmaps for all correctly classified MS patients (top) and all correctly classified healthy controls (bottom) in the holdout set. Values are normalized in the range $[-0.02, 0.02]$.

In Figure 5, we show the effects of transfer learning on the average relevance heatmaps for the MS patients in the hold-

out set. For the untrained model with random parameters (first row), only scarcely distributed individual voxels attain tiny relevance values. For the CNN model trained on ADNI and directly applied to MS patients (without fine-tuning; second row), more voxels are attributed relevance and are diffusely clustered. For the CNN model trained only on MS data (without pre-training; third row), strong relevance is projected to the ventricles and periventricular white matter. And finally, for the pre-trained model (transfer learning from ADNI to MS; last row), distinct clusters for both positive and negative relevance can be detected, which are more delineated than for the CNN model without pre-training.

To assess the contribution of normal-appearing brain matter, we compared the relevance maps between the CNN models trained on the original FLAIR data and the lesion-filled FLAIR data (for the performance see Table 2). In Figure 6, we depict the relevance for the 10 top-scored white matter regions, separately for both models. In general one can see that the relevance shifts from a distribution more evenly spread among multiple areas to a distribution with a prominent peak and otherwise low shares of relevance. Notably, relevance is shifted away from areas with large amounts of lesions such as posterior corona radiata, posterior thalamic radiata as well as tapetum towards mainly the corpus callosum and regions with very few lesions like fornix and external capsule (see supplementary Figure 2 for distribution of white matter lesions).

4. Discussion

4.1. Summary

In the present study, we introduced a transparent framework for analyzing neuroimaging data with CNNs that is able to explain individual classification decisions. By utilizing transfer learning we could further achieve good classification results from only a small data set of task-specific data. In combination with layer-wise relevance propagation (LRP), we could demonstrate the capacity of our framework to learn significant MS-relevant information from conventional MRI data. Notably, a pre-trained CNN was able to identify MS patients with an accuracy similar to a classical machine learning analysis, in which the T2 lesion load was used as input. This is quite remarkable, because the CNN model was considered to be naive by not being provided with any prior information on MS-relevant features such as hyperintense lesions. The subsequent visualization analysis, using heatmaps generated by LRP, revealed that the CNN model indeed uses (posterior) white matter lesions as primary information source. In addition, other information, e.g. in normal-appearing white and gray matter (e.g. the thalamus) have been found useful by the CNN model.

4.2. Key findings

CNNs learn to identify lesions as an important biomarker for MS. Although our pre-trained CNN model did not get any prior information about the relevance of hyperintense lesions for MS, it learned to successfully identify lesions as a primary

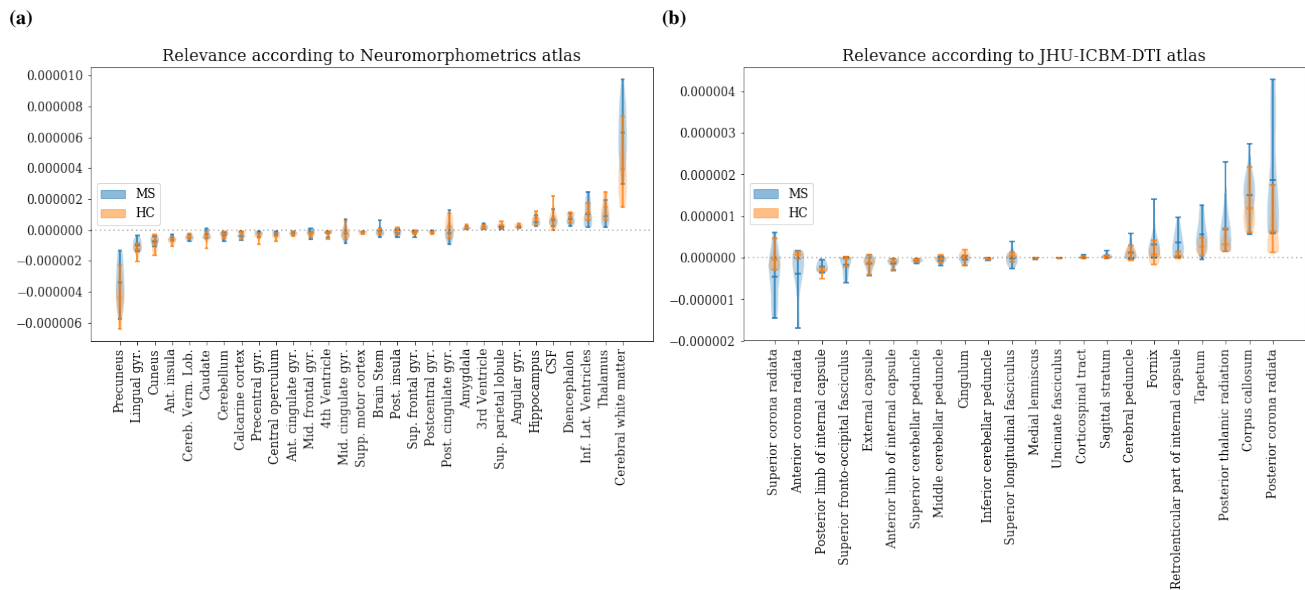


Figure 4: LRP relevance distribution over (a) 30 (mainly) gray matter areas from the Neuromorphometrics atlas and (b) 22 white matter areas from the JHU ICBM-DTI atlas, separately for MS patients and healthy controls in the holdout set.

information source. Notably, 9.71% of total relevance was attributed to lesion areas compared to a lesion coverage of only 0.44% in the training data set. We show that LRP heatmaps not only detect single lesions in individual patients but generally attributed most positive relevance to white matter areas around the posterior occipital horns. Importantly, the CNN model did not simply assign high relevance to hyperintense areas in the brain, but learned to distinguish between different lesion locations: while anterior periventricular lesions as well as lesions not bordering the lateral ventricles were assigned no or negative relevance, only posterior periventricular lesion areas were assigned positive relevance for MS. Strongest positive relevance was found for posterior corona radiata, corpus callosum and thalamic radiation, which are generally characterized by a high lesion density in MS patients (see [65] and supplementary Figures 2 and 3).

CNNs learn to identify relevant areas beyond lesions. The CNN model primarily focuses on lesions, but relevance has also been attributed to gray matter areas such as the thalamus, which is known to be affected in MS [6]. To further investigate what the CNN model learns beyond lesions, we repeated the analysis on lesion filled MRI data. As expected, the balanced accuracy as well as AUC decreased (by almost 17 and 6 percentage points respectively) and relevance has shifted away from regions which typically contain hyperintense lesions. The region that was assigned most relevance after lesion removal was the corpus callosum. While the corpus callosum is generally susceptible to demyelinating lesions [66, 67, 68] the literature also suggests further biomarkers such as axonal loss and diffuse atrophy [68, 69] or narrow T2 hyperintense bands along the callosal-septal interface [67]. The fornix, even though it contains a very small amount of lesions (see supplement

ary Figure 2 and [70]), is assigned positive relevance with lesions and an increased relevance without lesions. It has been shown that lower fractional anisotropy in the fornix is exhibited in MS subjects in comparison to healthy controls [71, 72]. Additionally, external capsule and superior cerebellar peduncle receive only positive relevance after lesion removal, which were found to be affected in MS patients [73, 74]. These results are generally in line with other machine learning studies finding differences in normal-appearing brain matter in MS patients [13, 14, 75]. It would be very interesting to further investigate whether our findings correlate with underlying pathological mechanisms only demonstrable by advanced MRI sequences such as diffusion weighted imaging or magnetization transfer imaging.

Transfer learning improves learning on small clinical neuroimaging cohorts, even across diseases and MRI sequences. In recent years, transfer learning has been successfully employed in brain lesion segmentation [55] and Alzheimer’s disease classification [53, 76, 56]. The latter studies used either autoencoders trained on MRI data or natural images [53, 76] or used one Alzheimer’s disease data set for pre-training and another Alzheimer’s disease data set for fine-tuning [56]. In the present study, we have shown that transfer learning can also help in learning (1) across diseases (Alzheimer’s disease to MS) and (2) across MRI sequences (MPRAGE to FLAIR). We demonstrated that not only the balanced accuracy increases drastically (about 16 percentage points), but also that LRP leads to much more focused heatmaps concentrating on (posterior) periventricular lesion areas. Given that our pre-trained model performed similar to a classical machine learning analysis using T2 lesion load as a classical biomarker in MS, we believe that larger data sets might allow for outperforming

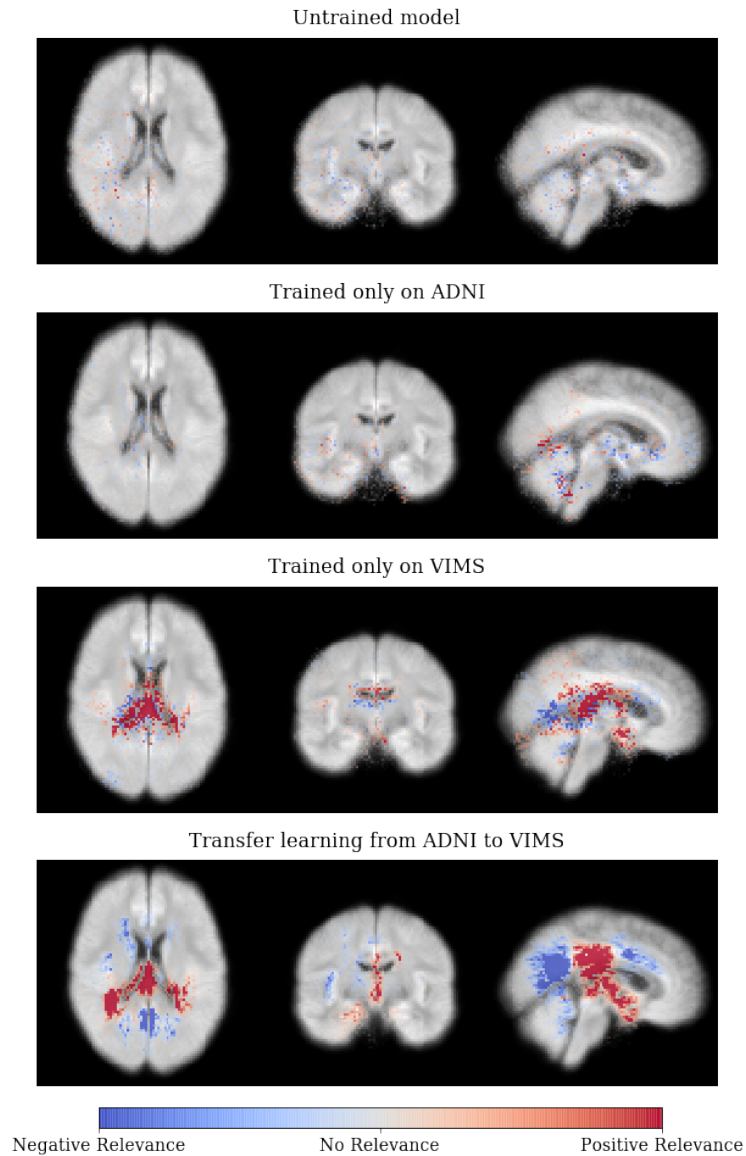


Figure 5: Average heatmaps for different CNN models applied to the MS (VIMS) cohort – starting from an untrained CNN model with random parameters over a CNN trained only on either ADNI or MS data to a CNN pre-trained on ADNI and fine-tuned on MS. As it can be seen, the fine-tuned model led to the most concise regions of positive and negative relevance. Please note that we averaged here the heatmaps over all (not only the correctly classified) MS patients in the holdout set and that the heatmap values here are not normalized to a fixed range but shown with respect to the minimum value of the untrained model.

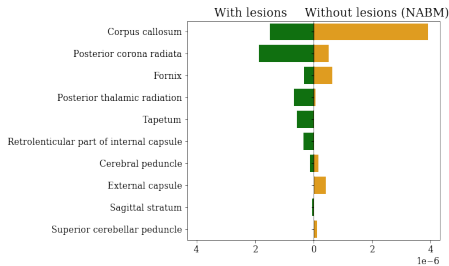


Figure 6: Comparison of average relevance distribution over white matter areas for a CNN model trained on original FLAIR data (left) and lesion-filled FLAIR data (right; NABM, normal-appearing brain matter). We calculated the relevance sum of both models (averaged over subjects) and show the 10 areas with the highest score.

models based on lesion masks in the future. Additionally, we are convinced that our approach – given a reasonable data basis – might also be very useful in answering more complex questions such as predicting disease progression.

4.3. Related work

Compared to other neurological diseases, in particular Alzheimer’s disease, only a few MS studies exist that employ machine learning methods outside the scope of lesion segmentation. We think that the main reasons are (1) the lack of easy accessible large open data bases such as the Alzheimer’s Neuroimaging Initiative (ADNI) data base and (2) the focus on white matter lesion volume as primary MRI-derived outcome measure in MS. Classical machine learning methods in combination with more or less sophisticated feature extraction methods, from both conventional and advanced MRI data, have been used to (1) diagnose MS [13, 15, 77] (2) decode symptom severity [14] (3) identify clinical subtypes [78, 79, 80] and (4) predict conversion from clinically isolated syndrome to MS [17]. Deep learning architectures have so far been implemented for lesion segmentation [24, 25, 26], predicting MS based on binary lesion masks [28], modelling brain and lesion variability [81] and finding differences in normal-appearing brain matter based on T1-weighted and myelin images [75]. To the best of our best knowledge, the present study is the first study employing CNNs and advanced visualization techniques for diagnosing MS based on FLAIR data.

It is generally recognized that, especially in the medical field, it is very important that classification decisions are reasonably explained even in light of high accuracies (which are no guarantee for a – from a human perspective – sensible discrimination strategy [82, 59]). Although a number of methods exist that generate individual heatmaps [60, 34, 30, 83], we focused here on the LRP method [33, 58, 59] which has a solid theoretical framework and has been extensively validated (see e.g. [37, 61, 59]). Very recently, LRP has shown to be very helpful for explaining cognitive states or Alzheimer’s disease diagnosis in deep neural networks trained on either functional or structural MRI data [36, 84]. In the present study, we demonstrated that LRP is capable of identifying reasonable areas supporting a MS diagnosis in addition to features needing further

clinical validation. By this, we have shown that those heatmaps can be very valuable in explaining decisions of neural networks trained on small sample sizes and to verify whether an algorithm has learned something meaningful (i.e. matching domain knowledge) or just spotted biases or artifacts in the data (see also [33, 59]).

4.4. Limitations

The main limitation of this study is the limited sample size. Although a sample size of $n = 147$ is comparable with other deep learning studies in the neuroimaging field [23], it is generally considered to be too low to learn robust representations from the data and to generalize to other data sets. To partly alleviate this problem, we pre-trained our network on ADNI data ($n = 921$) and fine-tuned it on the MS data. By visualizing the average heatmaps for MS patients, we show in addition to a balanced accuracy of 87.04 % that the CNN captures MS-relevant information by focusing on posterior ventricular regions usually characterized by a high rate of MS lesion incidences. Nevertheless, future studies should verify our results in larger data sets, preferably coming from different sites. Another limitation, related to the first one, is that we were limited in the choice of architecture used for the CNN analysis. To avoid overfitting, we have chosen a relatively simple CNN architecture and included different methods for regularization (drop out, L2-regularization and early stopping). Moreover, by registering the MRI data only linearly to MNI space, the regions contained in both atlases only roughly correspond to individual anatomical locations. On the other hand, non-linear registration can lead to strong deformations, in particular in patients, and we show here that our CNN model can also operate on a more native level (in accordance with [85]). To be able to make more specific anatomical claims in individual subjects, future studies might use individual atlases. And finally, heatmaps do neither allow to determine the underlying pathological mechanism (e.g. atrophy, demyelination or axonal loss) resulting in assigning a voxel to be relevant or to assess interactions between voxels. For this, one would have to take a deeper look into the specific filters that have been learned throughout the training process in combination with MR sequences more sensitive for certain tissue damage (e.g. diffusion weighted imaging). Nevertheless, we still believe that heatmaps can be very helpful in supplementing individual disease diagnoses by providing a simple and intuitive explanation.

5. Conclusion

In conclusion, we have shown that CNN models are capable of learning MS-relevant information from a typical-sized neuroimaging data set. In particular, we demonstrated that pre-training on additional data substantially increases the prediction performance (even across diseases and MRI sequences) and that the LRP method is very valuable not only in explaining individual network’s decisions, but also in generally helping to assess whether CNN models have learned significant features. Notably, our CNN models focus on hyperintense lesions

as primary information source, but also incorporates information from lesion location and normal-appearing brain areas. We see a high potential in the combination of CNNs, transfer learning and LRP heatmaps and are convinced that our framework might not only be helpful in other disease decoding studies, but also for answering more complex questions such as predicting disease progression or treatment response in individual subjects.

6. Funding

We acknowledge support from the German Research Foundation (DFG, 389563835), the Manfred and Ursula-Müller Stiftung and Charité – Universitätsmedizin Berlin (Rahel-Hirsch scholarship).

7. References

References

- [1] T. Mitchell, W. J. Culpepper, E. Nichols, Z. A. Bhutta, T. T. Gebrehiwot, S. I. Hay, I. A. Khalil, K. J. Krohn, X. Liang, M. Naghavi, A. H. Mokdad, M. R. Nixon, R. C. Reiner, B. Sartorius, M. Smith, R. Topor-Madry, A. Werdecker, T. Vos, V. L. Feigin, C. J. L. Murray, Global, regional, and national burden of multiple sclerosis 1990-2016: a systematic analysis for the Global Burden of Disease Study 2016., *The Lancet. Neurology* 18 (2019) 269–285.
- [2] D. S. Reich, C. F. Lucchinetti, P. A. Calabresi, *Multiple Sclerosis*, *New England Journal of Medicine* 378 (2018) 169–180.
- [3] A. J. Thompson, S. E. Baranzini, J. Geurts, B. Hemmer, O. Ciccarelli, *Multiple sclerosis*, *The Lancet* 391 (2018) 1622–1636.
- [4] R. Gheraldes, O. Ciccarelli, F. Barkhof, N. De Stefano, C. Enzinger, M. Filippi, M. Hofer, F. Paul, P. Preziosa, A. Rovira, G. C. DeLuca, L. Kappos, T. Yousry, F. Fazekas, J. Frederiksen, C. Gasperini, J. Sastre-Garriga, N. Evangelou, J. Palace, The current role of MRI in differentiating multiple sclerosis from its imaging mimics, *Nature Reviews Neurology* 14 (2018) 199–213.
- [5] M. J. Lowe, M. D. Phillips, J. T. Lurito, D. Mattson, M. Dzemidzic, V. P. Mathews, *Multiple Sclerosis: Low-Frequency Temporal Blood Oxygen Level-Dependent Fluctuations Indicate Reduced Functional Connectivity – Initial Results*, *Radiology* 224 (2002) 184–192.
- [6] C. J. Azevedo, S. Y. Cen, S. Khadka, S. Liu, J. Kornak, Y. Shi, L. Zheng, S. L. Hauser, D. Pelletier, Thalamic atrophy in multiple sclerosis: A magnetic resonance imaging marker of neurodegeneration throughout disease, *Annals of Neurology* 83 (2018) 223–234.
- [7] M. Absinta, P. Sati, D. S. Reich, *Advanced MRI and staging of multiple sclerosis lesions*, *Nature Reviews Neurology* 12 (2016) 358–368.
- [8] M. Filippi, M. A. Rocca, O. Ciccarelli, N. De Stefano, N. Evangelou, L. Kappos, A. Rovira, J. Sastre-Garriga, M. Tintorè, J. L. Frederiksen, C. Gasperini, J. Palace, D. S. Reich, B. Banwell, X. Montalban, F. Barkhof, MRI criteria for the diagnosis of multiple sclerosis: MAGNIMS consensus guidelines, *The Lancet Neurology* 15 (2016) 292–303.
- [9] P. Sati, J. Oh, R. T. Constable, N. Evangelou, C. R. G. Guttmann, R. G. Henry, E. C. Klawiter, C. Mainero, L. Massacesi, H. McFarland, F. Nelson, D. Ontaneda, A. Rauscher, W. D. Rooney, A. P. R. Samaraweera, R. T. Shinohara, R. A. Sobel, A. J. Solomon, C. A. Treaba, J. Wuerfel, R. Zivadinov, N. L. Sicotte, D. Pelletier, D. S. Reich, o. b. o. t. N. Cooperative, The central vein sign and its clinical evaluation for the diagnosis of multiple sclerosis: a consensus statement from the North American Imaging in Multiple Sclerosis Cooperative, *Nature Reviews Neurology* 12 (2016) 714–722.
- [10] Y. Backner, J. Kuchling, S. Massarwa, T. Oberwahrenbrock, C. Finke, J. Bellmann-Strobl, K. Ruprecht, A. U. Brandt, H. Zimmermann, N. Raz, F. Paul, N. Levin, Anatomical Wiring and Functional Networking Changes in the Visual System Following Optic Neuritis, *JAMA Neurology* 75 (2018) 287.
- [11] G. Orrù, W. Pettersson-Yeo, A. F. Marquand, G. Sartori, A. Mechelli, Using support vector machine to identify imaging biomarkers of neurological and psychiatric disease: a critical review, *Neuroscience & Biobehavioral Reviews* 36 (2012) 1140–1152.
- [12] C. W. Woo, L. J. Chang, M. A. Lindquist, T. D. Wager, Building better biomarkers: Brain models in translational neuroimaging, *Nature Neuroscience* 20 (2017) 365–377.
- [13] M. Weygandt, K. Hackmack, C. Pfüller, J. Bellmann-Strobl, F. Paul, F. Zipp, J. D. Haynes, MRI pattern recognition in multiple sclerosis normal-appearing brain areas, *PLoS ONE* 6 (2011) e21138.
- [14] K. Hackmack, M. Weygandt, C. F. Pfueller, J. Bellmann-Strobl, J. Wuerfel, J.-D. Haynes, F. Paul, Can we overcome the clinico-radiological paradox’ in multiple sclerosis?, *Journal of Neurology* 259 (2012) 2151–2160.
- [15] K. Hackmack, F. Paul, M. Weygandt, C. Allefeld, J. D. Haynes, Multi-scale classification of disease using structural MRI and wavelet transform, *NeuroImage* 62 (2012) 48–58.
- [16] M. Weygandt, H.-M. Hummel, K. Schregel, K. Ritter, C. Allefeld, E. Dommers, P. Huppke, J. Haynes, J. Wuerfel, J. Gärtner, MRI-based diagnostic biomarkers for early onset pediatric multiple sclerosis, *NeuroImage: Clinical* 7 (2015) 400–408.
- [17] V. Wotschel, D. C. Alexander, P. P. Kwok, D. T. Chard, M. L. Stromillo, N. De Stefano, A. J. Thompson, D. H. Miller, O. Ciccarelli, Predicting outcome in clinically isolated syndrome using machine learning., *NeuroImage: Clinical* 7 (2015) 281–7.
- [18] Y. Lecun, Y. Bengio, G. Hinton, Deep learning, *Nature* 521 (2015) 436–444.
- [19] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. W. M. van der Laak, B. van Ginneken, C. I. Sánchez, A Survey on Deep Learning in Medical Image Analysis, Preprint in arxiv: <http://arxiv.org/abs/1702.05747> (2017).
- [20] P. Rajpurkar, J. Irvin, A. Bagul, D. Ding, T. Duan, H. Mehta, B. Yang, K. Zhu, D. Laird, R. L. Ball, C. Langlotz, K. Shpanskaya, M. P. Lungren, A. Y. Ng, MURA: Large Dataset for Abnormality Detection in Musculoskeletal Radiographs, Preprint in arxiv: <http://arxiv.org/abs/1712.06957> (2017).
- [21] P. Rajpurkar, J. Irvin, K. Zhu, B. Yang, H. Mehta, T. Duan, D. Ding, A. Bagul, C. Langlotz, K. Shpanskaya, M. P. Lungren, A. Y. Ng, CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning, Preprint in arxiv: <http://arxiv.org/abs/1711.05225> (2017).
- [22] J. De Fauw, J. R. Ledsam, B. Romera-Paredes, S. Nikolov, N. Tomasev, S. Blackwell, H. Askham, X. Glorot, B. O’Donoghue, D. Visentin, G. van den Driessche, B. Lakshminarayanan, C. Meyer, F. Mackinder, S. Bouton, K. Ayoub, R. Chopra, D. King, A. Karthikesalingam, C. O. Hughes, R. Raine, J. Hughes, D. A. Sim, C. Egan, A. Tufail, H. Montgomery, D. Hassabis, G. Rees, T. Back, P. T. Khaw, M. Suleyman, J. Cornebise, P. A. Keane, O. Ronneberger, Clinically applicable deep learning for diagnosis and referral in retinal disease, *Nature Medicine* 24 (2018) 1342–1350.
- [23] S. Vieira, W. H. Pinaya, A. Mechelli, Using deep learning to investigate the neuroimaging correlates of psychiatric and neurological disorders: Methods and applications, *Neuroscience & Biobehavioral Reviews* 74 (2017) 58–75.
- [24] S. Valverde, M. Cabezas, E. Roura, S. González-Villà, D. Pareto, J. C. Vilanova, L. Ramió-Torrentà, À. Rovira, A. Oliver, X. Lladó, Improving automated multiple sclerosis lesion segmentation with a cascaded 3D convolutional neural network approach, *NeuroImage* 155 (2017) 159–168.
- [25] D. K. B. Li, T. Brosch, L. Y. W. Tang, A. Traboulsee, R. Tam, Y. Yoo, Deep 3D Convolutional Encoder Networks With Shortcuts for Multiscale Feature Integration Applied to Multiple Sclerosis Lesion Segmentation, *IEEE Transactions on Medical Imaging* 35 (2016) 1229–1239.
- [26] H. Khastavaneh, H. Ebrahimpour-Komleh, Neural Network-Based Learning Kernel for Automatic Segmentation of Multiple Sclerosis Lesions on Magnetic Resonance Images., *Journal of biomedical physics & engineering* 7 (2017) 155–162.
- [27] S.-H. Wang, C. Tang, J. Sun, J. Yang, C. Huang, P. Phillips, Y.-D. Zhang, Multiple sclerosis identification by 14-layer convolutional neural network with batch normalization, dropout, and stochastic pooling, *Frontiers in neuroscience* 12 (2018) 818.
- [28] Y. Yoo, L. W. Tang, T. Brosch, D. K. B. Li, L. Metz, A. Traboulsee,

- R. Tam, Deep Learning of Brain Lesion Patterns for Predicting Future Disease Activity in Patients with Early Symptoms of Multiple Sclerosis Deep Learning of Lesion Patterns for Early MS Activity Prediction, LNCS 10008 (2016) 86–94.
- [29] D. Castelvocchi, Can we open the black box of AI?, *Nature* 538 (2016) 20–23.
- [30] M. Zeiler, R. Fergus, Visualizing and Understanding Convolutional Networks, in: D. Fleet, T. Pajdla, B. Schiele, T. Tuytelaars (Eds.), *Computer Vision ECCV 2014*, volume 8689 of *Lecture Notes in Computer Science*, Springer International Publishing, 2014, pp. 818–833.
- [31] C. Olah, A. Mordvintsev, L. Schubert, Feature visualization, *Distill* (2017).
- [32] K. Simonyan, A. Zisserman, Two-Stream Convolutional Networks for Action Recognition in Videos, in: *Advances in Neural Information Processing Systems*, pp. 568–576.
- [33] 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 (2015) e0130140.
- [34] J. T. Springenberg, A. Dosovitskiy, T. Brox, M. Riedmiller, Striving for Simplicity: The All Convolutional Net, *ICLR* (2015).
- [35] J. Rieke, F. Eitel, M. Weygandt, J.-D. Haynes, K. Ritter, Visualizing convolutional networks for mri-based diagnosis of alzheimers disease, in: *Understanding and Interpreting Machine Learning in Medical Image Computing Applications*, Springer, 2018, pp. 24–31.
- [36] M. Boehle, F. Eitel, M. Weygandt, K. Ritter, Visualizing evidence for Alzheimer’s disease in deep neural networks trained on structural MRI data, Preprint in arxiv: <http://arxiv.org/abs/1903.07317> (2019).
- [37] W. Samek, A. Binder, G. Montavon, S. Lapuschkin, K. Müller, Evaluating the visualization of what a deep neural network has learned, *IEEE Transactions on Neural Networks and Learning Systems* 28 (2017) 2660–2673.
- [38] C. H. Polman, S. C. Reingold, B. Banwell, M. Clanet, J. A. Cohen, M. Filippi, K. Fujihara, E. Havrdova, M. Hutchinson, L. Kappos, F. D. Lublin, X. Montalban, P. O’Connor, M. Sandberg-Wollheim, A. J. Thompson, E. Waubant, B. Weinshenker, J. S. Wolinsky, Diagnostic criteria for multiple sclerosis: 2010 revisions to the McDonald criteria., *Annals of neurology* 69 (2011) 292–302.
- [39] J. Kuchling, Y. Backner, F. C. Oertel, N. Raz, J. Bellmann-Strobl, K. Ruprecht, F. Paul, N. Levin, A. U. Brandt, M. Scheel, Comparison of probabilistic tractography and tract-based spatial statistics for assessing optic radiation damage in patients with autoimmune inflammatory disorders of the central nervous system., *NeuroImage: Clinical* 19 (2018) 538–550.
- [40] P. A. Yushkevich, J. Piven, H. C. Hazlett, R. G. Smith, S. Ho, J. C. Gee, G. Gerig, User-guided 3D active contour segmentation of anatomical structures: Significantly improved efficiency and reliability, *NeuroImage* 31 (2006) 1116–1128.
- [41] K. J. Friston, J. Ashburner, S. J. Kiebel, T. E. Nichols, W. D. Penny (Eds.), *Statistical Parametric Mapping: The Analysis of Functional Brain Images*, Academic Press, 2007.
- [42] B. B. Avants, N. J. Tustison, G. Song, P. A. Cook, A. Klein, C. Gee, A Reproducible Evaluation of ANTs Similarity Metric Performance in Brain Image Registration, *Neuroimage* 54 (2011) 2033–2044.
- [43] S. M. Smith, M. Jenkinson, M. W. Woolrich, C. F. Beckmann, T. E. J. Behrens, H. Johansen-Berg, P. R. Bannister, M. De Luca, I. Drobnjak, D. E. Flitney, Others, Advances in functional and structural MR image analysis and implementation as FSL., *Neuroimage* 23 (2004) S208—S219.
- [44] S. Valverde, A. Oliver, X. Lladó, A white matter lesion-filling approach to improve brain tissue volume measurements, *NeuroImage: Clinical* 6 (2014) 86–92.
- [45] J. Ashburner, K. J. Friston, Image segmentation, in: R. S. J. Frackowiak, K. J. Friston, C. Frith, R. Dolan, K. J. Friston, C. J. Price, S. Zeki, J. Ashburner, W. D. Penny (Eds.), *Human Brain Function*, Academic Press, 2nd edition, 2003.
- [46] B. T. Wyman, D. J. Harvey, K. Crawford, M. A. Bernstein, O. Carmichael, P. E. Cole, P. K. Crane, C. DeCarli, N. C. Fox, J. L. Gunter, D. Hill, R. J. Killiany, C. Pachai, A. J. Schwarz, N. Schuff, M. L. Senjem, J. Suhy, P. M. Thompson, M. Weiner, C. R. Jack, Alzheimer’s Disease Neuroimaging Initiative, Standardization of analysis sets for reporting results from ADNI MRI data, *Alzheimer’s & Dementia* 9 (2013) 332–337.
- [47] D. P. Kingma, J. Ba, Adam: A Method for Stochastic Optimization, Preprint in arxiv: <http://arxiv.org/abs/1412.6980> (2014).
- [48] F. Chollet, Others, Keras, <https://github.com/fchollet/keras>, 2015.
- [49] M. Abadi, Others, Tensorflow: Large-scale machine learning on heterogeneous systems, 2015.
- [50] K. Crammer, M. Kearns, J. Wortman, Learning from Multiple Sources, *Journal of Machine Learning Research* 9 (2008) 1757–1774.
- [51] L. Duan, I. W. Tsang, D. Xu, T.-S. Chua, Domain adaptation from multiple sources via auxiliary classifiers, in: *Proceedings of the 26th Annual International Conference on Machine Learning - ICML ’09*, ACM Press, New York, New York, USA, 2009, pp. 1–8.
- [52] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, J. W. Vaughan, A theory of learning from different domains, *Machine Learning* 79 (2010) 151–175.
- [53] A. Gupta, M. Ayhan, A. Maida, Natural image bases to represent neuroimaging data, in: S. Dasgupta, D. Mallester (Eds.), *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, volume 28, *JMLR Workshop and Conference Proceedings*, 2013, pp. 987–994.
- [54] N. Tajbakhsh, J. Y. Shin, S. R. Gurudu, R. T. Hurst, C. B. Kendall, M. B. Gotway, J. Liang, Convolutional Neural Networks for Medical Image Analysis: Full Training or Fine Tuning?, *IEEE Transactions on Medical Imaging* 35 (2016) 1299–1312.
- [55] M. Ghafoorian, A. Mehrtash, T. Kapur, N. Karssemeijer, E. Marchiori, M. Pesteie, C. R. G. Guttman, F.-E. de Leeuw, C. M. Tempny, B. van Ginneken, A. Fedorov, P. Abolmaesumi, B. Platel, W. M. Wells, Transfer Learning for Domain Adaptation in MRI: Application in Brain Lesion Segmentation, in: M. Descoteaux, L. Maier-Hein, A. Franz, P. Jannin, D. L. Collins, S. Duchesne (Eds.), *Medical Image Computing and Computer Assisted Intervention - MICCAI 2017*, Springer International Publishing, Cham, 2017, pp. 516–524.
- [56] E. Hosseini-Asl, M. Ghazal, A. Mahmoud, A. Aslantas, A. M. Shalaby, M. F. Casanova, G. N. Barnes, G. Gimel’farb, R. Keynton, A. El-Baz, Alzheimer’s disease diagnostics by a 3D deeply supervised adaptable convolutional network., *Frontiers in bioscience (Landmark edition)* 23 (2018) 584–596.
- [57] S. Basaia, F. Agosta, L. Wagner, E. Canu, G. Magnani, R. Santangelo, M. Filippi, Automated classification of Alzheimer’s disease and mild cognitive impairment using a single MRI and deep neural networks, *NeuroImage: Clinical* 21 (2019) 101645.
- [58] G. Montavon, W. Samek, K.-R. Müller, Methods for interpreting and understanding deep neural networks, *Digital Signal Processing* 73 (2018) 1–15.
- [59] S. Lapuschkin, S. Wäldchen, A. Binder, G. Montavon, W. Samek, K.-R. Müller, Unmasking Clever Hans Predictors and Assessing What Machines Really Learn, *Nature Communications* 10 (2019) 1096.
- [60] K. Simonyan, A. Vedaldi, A. Zisserman, Deep inside convolutional networks: Visualising image classification models and saliency maps, Preprint in arxiv: <https://arxiv.org/abs/1312.6034> (2013).
- [61] W. Samek, T. Wiegand, K.-R. Müller, Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models, Preprint in arxiv: <https://arxiv.org/abs/1708.08296> (2017).
- [62] M. Alber, S. Lapuschkin, P. Seegerer, M. Hägele, K. T. Schütt, G. Montavon, W. Samek, K.-R. Müller, S. Dähne, P.-J. Kindermans, iNNvestigate neural networks!, *CoRR* abs/1808.0 (2018).
- [63] R. Bakker, P. Tiesinga, R. Kötter, The Scalable Brain Atlas: Instant Web-Based Access to Public Brain Atlases and Related Content, *Neuroinformatics* 13 (2015) 353–366.
- [64] S. S. Mori, B. J. Crain, MRI atlas of human white matter, Elsevier, 2005.
- [65] A. Gass, E.-W. Radue, T. E. Nichols, F. Barkhof, H. Vrenken, S. Traud, L. Kappos, C. Polman, Y. Naegelin, T. Sprenger, P. Kuster, K. Bendfeldt, N. Mueller-Lenke, L. Filli, L. Hofstetter, S. J. Borgwardt, Spatiotemporal distribution of white matter lesions in relapsingremitting and secondary progressive multiple sclerosis, *Multiple Sclerosis Journal* 18 (2012) 1577–1584.
- [66] R. O. Barnard, M. Triggs, Corpus callosum in multiple sclerosis., *Journal of neurology, neurosurgery, and psychiatry* 37 (1974) 1259–64.
- [67] N. Garg, S. W. Reddel, D. H. Miller, J. Chataway, D. S. Riminton, Y. Barnett, L. Masters, M. H. Barnett, T. A. Hardy, The corpus callosum in the diagnosis of multiple sclerosis and other CNS demyelinating and inflammatory diseases, *Journal of Neurology, Neurosurgery & Psychiatry* 86

- (2015) 1374–1382.
- [68] D. Renard, G. Castelnovo, C. Campello, S. Bouly, A. Le Floch, E. Thouvenot, A. Waconge, G. Taieb, An MRI review of acquired corpus callosum lesions, *Journal of Neurology, Neurosurgery & Psychiatry* 85 (2014) 1041–1048.
- [69] N. Evangelou, D. Konz, M. M. Esiri, S. Smith, J. Palace, P. M. Matthews, Regional axonal loss in the corpus callosum correlates with cerebral white matter lesion volume and distribution in multiple sclerosis, *Brain* 123 (2000) 1845–1849.
- [70] A. G. Thomas, P. Koumellis, R. A. Dineen, The Fornix in Health and Disease: An Imaging Review, *RadioGraphics* 31 (2011) 1107–1121.
- [71] S. D. Roosendaal, J. J. G. Geurts, H. Vrenken, H. E. Hulst, K. S. Cover, J. A. Castelijns, P. J. W. Pouwels, F. Barkhof, Regional DTI differences in multiple sclerosis patients, *NeuroImage* 44 (2009) 1397–1403.
- [72] K. C. Kern, A. D. Ekstrom, N. A. Suthana, B. S. Giesser, M. Montag, A. Arshanapalli, S. Y. Bookheimer, N. L. Sicotte, Fornix damage limits verbal memory functional compensation in multiple sclerosis, *NeuroImage* 59 (2012) 2932–2940.
- [73] V. M. Anderson, C. A. Wheeler-Kingshott, K. Abdel-Aziz, D. H. Miller, A. Toosy, A. J. Thompson, O. Ciccarelli, A comprehensive assessment of cerebellar damage in multiple sclerosis using diffusion tractography and volumetric analysis, *Multiple Sclerosis Journal* 17 (2011) 1079–1087.
- [74] C. Zhang, Y. Liu, X.-m. Han, J.-b. Gu, R. Bakshi, Z. Han, H.-j. Tian, X. Cao, Correlation between white matter damage and gray matter lesions in multiple sclerosis patients, *Neural Regeneration Research* 12 (2017) 787.
- [75] Y. Yoo, L. Y. Tang, T. Brosch, D. K. Li, S. Kolind, I. Vavasour, A. Rauscher, A. L. MacKay, A. Traboulsee, R. C. Tam, Deep learning of joint myelin and T1w MRI features in normal-appearing brain tissue to distinguish between multiple sclerosis patients and healthy controls, *NeuroImage: Clinical* 17 (2018) 169–178.
- [76] A. Payan, G. Montana, Predicting Alzheimer’s disease: a neuroimaging study with 3D convolutional neural networks, *CoRR abs/1502.0* (2015).
- [77] M. Zurita, C. Montalba, T. Labbé, J. P. Cruz, J. Dalboni da Rocha, C. Tejos, E. Ciampi, C. Cárcamo, R. Sitaram, S. Uribe, Characterization of relapsing-remitting multiple sclerosis patients using support vector machine classifications of functional and diffusion MRI data, *NeuroImage: Clinical* 20 (2018) 724–730.
- [78] A. Eshaghi, S. Riyahi-Alam, R. Saeedi, T. Roostaei, A. Nazeri, A. Aghsaei, R. Doosti, H. Ganjgahi, B. Bodini, A. Shakourirad, M. Pakravan, H. Ghanaati, K. Firouznia, M. Zarei, A. R. Azimi, M. A. Sahraian, Classification algorithms with multi-modal data fusion could accurately distinguish neuromyelitis optica from multiple sclerosis, *NeuroImage: Clinical* 7 (2015) 306–314.
- [79] A. Eshaghi, R. V. Marinescu, A. L. Young, N. C. Firth, F. Prados, M. Jorge Cardoso, C. Tur, F. De Angelis, N. Cawley, W. J. Brownlee, N. De Stefano, M. Laura Stromillo, M. Battaglini, S. Ruggieri, C. Gasperini, M. Filippi, M. A. Rocca, A. Rovira, J. Sastre-Garriga, J. J. Geurts, H. Vrenken, V. Wottschel, C. E. Leurs, B. Uitdehaag, L. Pirpamer, C. Enzinger, S. Ourselin, C. A. Gandini Wheeler-Kingshott, D. Chard, A. J. Thompson, F. Barkhof, D. C. Alexander, O. Ciccarelli, Progression of regional grey matter atrophy in multiple sclerosis, *Brain* 141 (2018) 1665–1677.
- [80] T. E. Nichols, S. J. Borgwardt, L. Kappos, P. Kuster, N. Mueller-Lenke, S. Traud, R. Smieskova, E.-W. Radue, Y. Naegelin, K. Bendfeldt, S. Klöppel, Multivariate pattern classification of gray matter pathology in multiple sclerosis, *NeuroImage* 60 (2012) 400–408.
- [81] T. Brosch, Efficient deep learning of 3D structural brain MRIs for manifold learning and lesion segmentation with application to multiple sclerosis (2016).
- [82] S. Lapuschkin, A. Binder, G. Montavon, K.-R. Müller, W. Samek, *Analyzing Classifiers: Fisher Vectors and Deep Neural Networks*, 2016.
- [83] L. M. Zintgraf, T. S. Cohen, T. Adel, M. Welling, *Visualizing Deep Neural Network Decisions: Prediction Difference Analysis*, ICLR (2017).
- [84] A. W. Thomas, H. R. Heekeren, K.-R. Müller, W. Samek, *Interpretable LSTMs For Whole-Brain Neuroimaging Analyses*, Preprint in arxiv: <http://arxiv.org/abs/1810.09945> (2018).
- [85] H.-I. Suk, S.-W. Lee, D. Shen, Hierarchical feature representation and multimodal fusion with deep learning for AD/MCI diagnosis, *NeuroImage* 101 (2014) 569–582.

1. Supplementary Materials

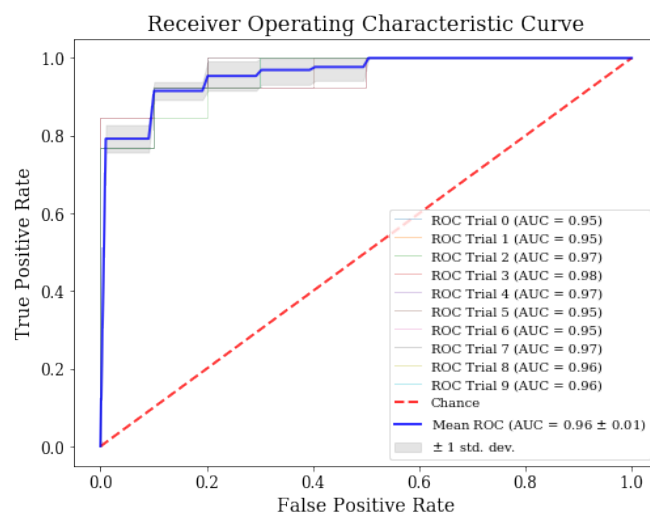


Figure 1: Performance of the pre-trained CNN model (fine-tuned on MS) on the holdout set.

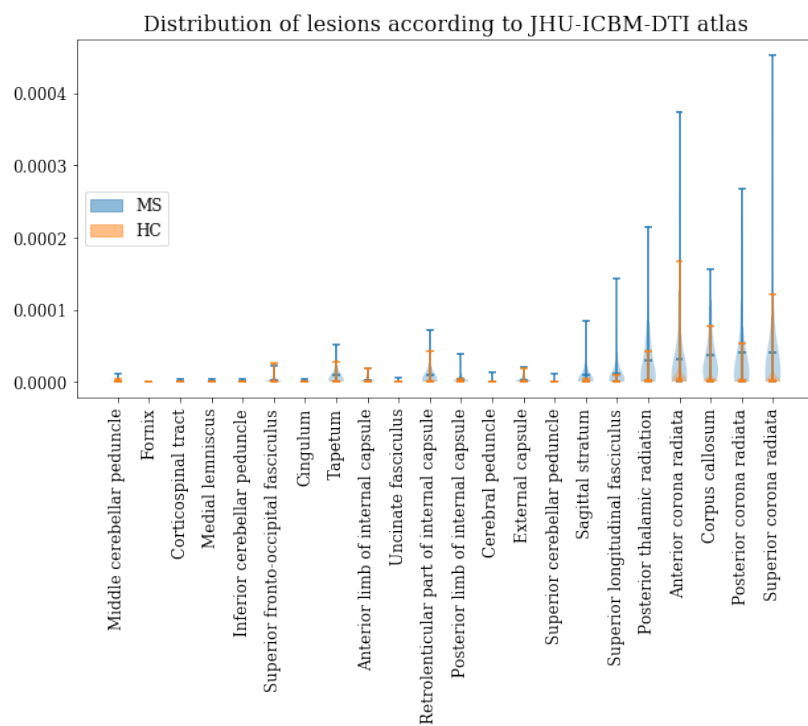


Figure 2: Lesion distribution over white matter areas from the JHU ICBM-DTI atlas, separately for MS patients and healthy controls (HC).

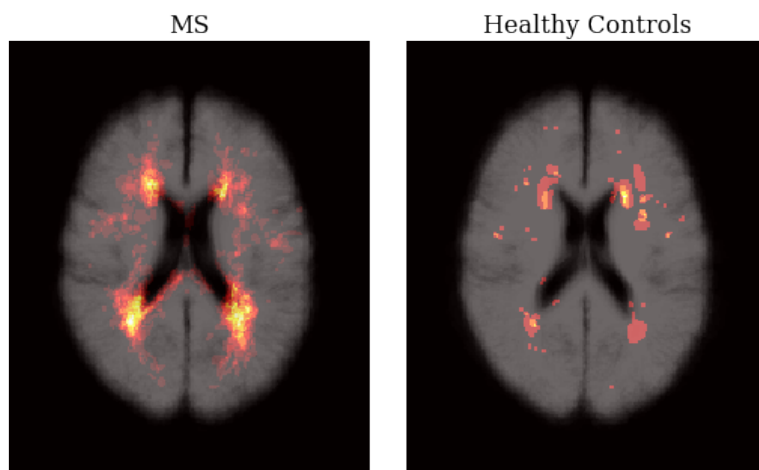


Figure 3: Lesion distribution in MS patients (left) and healthy controls (right).