


SCIENTIFIC REPORTS



OPEN

Location Sensitive Deep Convolutional Neural Networks for Segmentation of White Matter Hyperintensities

Mohsen Ghafoorian^{1,2}, Nico Karssemeijer², Tom Heskes¹, Inge W. M. van Uden³, Clara I. Sanchez², Geert Litjens², Frank-Erik de Leeuw³, Bram van Ginneken², Elena Marchiori¹ & Bram Platel^{1,2} 

The anatomical location of imaging features is of crucial importance for accurate diagnosis in many medical tasks. Convolutional neural networks (CNN) have had huge successes in computer vision, but they lack the natural ability to incorporate the anatomical location in their decision making process, hindering success in some medical image analysis tasks. In this paper, to integrate the anatomical location information into the network, we propose several deep CNN architectures that consider multi-scale patches or take explicit location features while training. We apply and compare the proposed architectures for segmentation of white matter hyperintensities in brain MR images on a large dataset. As a result, we observe that the CNNs that incorporate location information substantially outperform a conventional segmentation method with handcrafted features as well as CNNs that do not integrate location information. On a test set of 50 scans, the best configuration of our networks obtained a Dice score of 0.792, compared to 0.805 for an independent human observer. Performance levels of the machine and the independent human observer were not statistically significantly different (p -value = 0.06).

White matter hyperintensities (WMH), also known as leukoaraiosis or white matter lesions are a common finding on brain MR images of patients diagnosed with small vessel disease (SVD)¹, multiple sclerosis², Parkinsonism³, stroke⁴, Alzheimer's disease⁵ and Dementia⁶. WMHs often represent areas of demyelination found in the white matter of the brain, but they can also be caused by other mechanisms such as edema. WMHs are best observable in fluid-attenuated inversion recovery (FLAIR) MR images, as high value signals⁷. The prevalence of WMHs among SVD patients has been reported to reach up to 95% depending on the population studied and the imaging technique used⁸. Studies have reported a relationship between WMH severity and other neurological disturbances and symptoms including cognitive decline^{9,10}, gait dysfunction¹¹, hypertension¹² as well as depression¹³ and mood disturbances¹⁴. It has been shown that using a more accurate WMH volumetric assessment, a better association with clinical measures of physical performance and cognition is achieved¹⁵.

Accurate quantification of WMHs in terms of total volume and distribution is believed to be of clinical importance for prognosis, tracking of disease progression and assessment of the treatment effectiveness¹⁶. However, manual segmentation of WMHs is a laborious time consuming task that makes it infeasible for larger datasets and in clinical practice. Furthermore, manual segmentation is subject to considerable inter- and intra-rater variability¹⁷.

In the last decade, many automated and semi-automated algorithms have been proposed that can be classified into two general categories. Some methods use supervised machine learning algorithms, often using hand-crafted features^{18–29} or more recently with learned representations^{30–34}. This is while other methods use unsupervised approaches^{35–42} to cluster WMHs as outliers or model them with additional classes. Although a multitude of

¹Institute for Computing and Information Sciences, Radboud University, Nijmegen, The Netherlands. ²Diagnostic Image Analysis Group, Department of Radiology and Nuclear Medicine, Radboud University Medical Center, Nijmegen, The Netherlands. ³Donders Institute for Brain, Cognition and Behaviour, Department of Neurology, Radboud University Medical Center, Nijmegen, The Netherlands. Correspondence and requests for materials should be addressed to M.G. (email: mohsen.ghafoorian@radboudumc.nl)

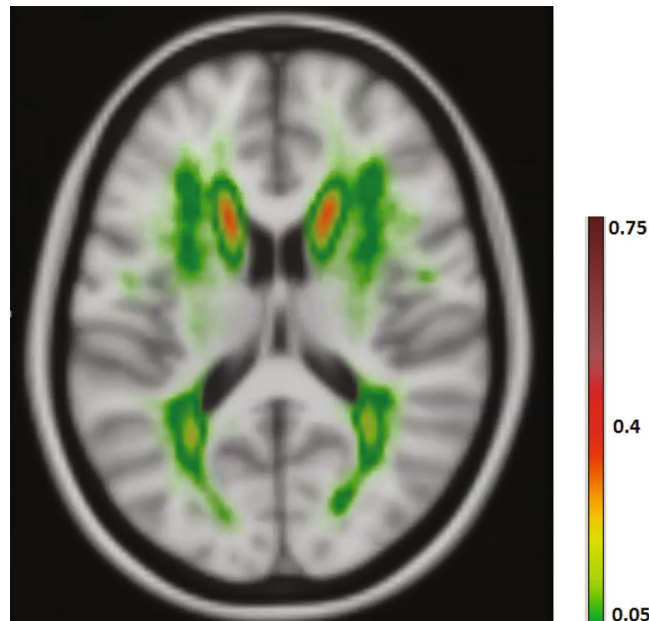


Figure 1. A pattern is observable in WMHs occurrence probability map.

approaches has been suggested for this problem, a truly reliable fully automated method that performs as good as human readers has not been identified^{43, 44}.

Deep neural networks^{45, 46} are biologically plausible learning structures, inspired by early neuroscience-related work^{47, 48} and have so far claimed human level or super-human performances in several different domains^{49–53}. Convolutional neural networks (CNN)⁵⁴, perhaps the most popular form of deep neural networks, have attracted enormous attention from the computer vision community since Alex Krizhevsky's network⁵⁵ won the Imagenet competition⁵⁶ by a large margin. Although the initial focus of CNN methods was concentrated on image classification, soon the framework was extended to cover segmentation as well. A natural way to apply CNNs to segmentation tasks is to train a network in a sliding-window setup to predict the label of each pixel/voxel considering a local neighborhood, which is usually referred to as a patch^{52, 57–59}. Later fully convolutional neural networks were proposed to computationally optimize the segmentation process^{60, 61}.

Deep neural networks have recently been widely used in many medical image analysis domains including lesion detection, image segmentation, shape modeling and image registration^{62, 63}. In particular on neuroimaging, several studies are proposed using CNNs for brain extraction⁶⁴, tissue and anatomical region segmentation^{65–70}, tumor segmentation^{71–74}, lacune detection⁷⁵, microbleed detection^{76, 77}, and brain lesion segmentation^{29, 31–34}.

In many bio-medical segmentation applications, including the segmentation of WMHs, anatomical location information plays an important role for an accurate classification of voxels^{18, 23, 43, 44, 78} (see Fig. 1). In contrast, in commonly used segmentation benchmarks in the computer vision community, such as general scene labeling and crowd segmentation, it is normally not a valid assumption to consider pixel/voxel spatial location as an important piece of information. This explains why the literature lacks enough studies investigating ways to integrate spatial information into CNNs.

In this study, we train a number of CNNs to build systems for an accurate fully-automated segmentation of WMHs. We train, validate and evaluate our networks with a large dataset of more than 500 patients, that enables us to learn optimal values for millions of weights in our deep networks. In order to feed the CNN with location information, it is possible to incorporate multi-scale patches or add an explicit set of spatial features to the network. We evaluate and compare three different strategies and network architectures for providing the networks with more context/spatial location information. Experimental results suggest not only our best performing network outperforms a conventional segmentation method with hand-crafted features with a considerable margin, but also its performance does not significantly differ from an independent human observer.

To summarize, the main contributions of the paper are the following: 1) Comparing and discussing the different strategies for fusing multi-scale information within a CNN on the WMH segmentation domain. 2) Integrating location features with the CNN in the same pass as the network is being trained. 3) Achieving results that are comparable to that of a human expert on a large set of independent test images.

Materials

Data. The research presented in this paper uses data from a longitudinal study called the Radboud University Nijmegen Diffusion tensor and Magnetic resonance imaging Cohort (RUN DMC)¹. Baseline scanning was performed in 2006. The patients were rescanned in 2011/2012 and currently a third follow-up is being acquired.

Subjects. Subjects for the RUN DMC study were selected at baseline based on the following inclusion criteria¹: (a) aged between 50 and 85 years (b) cerebral SVD on neuroimaging (appearance of WMHs and/or lacunes).

Modality	TR/TE/TI	Flip angle	Voxel size	Interslice gap
T1	2250/3.68/850 ms	15°	1.0 × 1.0 × 1.0	0
FLAIR	9000/84/2200 ms	15°	1.2 × 1.0 × 5.0	1 mm

Table 1. MR imaging protocol specification for the T1 and FLAIR modalities.

Exclusion criteria comprised: presence of (a) dementia (b) parkinson(-ism) (c) intracranial hemorrhage (d) life expectancy less than six months (e) intracranial space occupying lesion (f) (psychiatric) disease interfering with cognitive testing or follow-up (g) recent or current use of acetylcholine-esterase inhibitors, neuroleptic agents, L-dopa or dopa-a(NTA)gonists (h) non-SVD related WMH (e.g. MS) (i) prominent visual or hearing impairment (j) language barrier and (k) MRI contraindications. Based on these criteria, MRI scans of 503 patients were taken at baseline.

Magnetic resonance imaging. The machine used for the baseline was a single 1.5 Tesla scanner (Magnetom Sonata, Siemens Medical Solution, Erlangen, Germany). Details of the imaging protocol are listed in Table 1.

Reference annotations. Reference annotations were created in a slice by slice manner by two experienced raters, manually contouring hyperintense lesions on FLAIR MRI that did not show corresponding cerebrospinal fluid like hypo-intense lesions on the T1 weighted image. Gliosis surrounding lacunes and territorial infarcts were not considered to be WMH related to SVD⁷⁹. One of the observers (observer 1) manually annotated all of the cases. 50 of these 503 images were selected at random and were annotated also by another human observer (observer 2).

Preprocessing. Before supplying the data to our networks, we first pre-processed the data with the following four steps:

Multi-modal registration. Due to possible movement of patients during scanning, the image coordinates of the T1 and FLAIR modalities might not represent the same location. Thus we transformed the T1 image to align with the FLAIR image in the native space using FSL-FLIRT⁸⁰ implementation of rigid registration with trilinear interpolation and mutual information optimization criteria. Also to obtain a mapping between patient space and an atlas space, the ICBM152 atlas⁸¹ was non-linearly registered to each patient image using FSL-FNIRT⁸². The resulting transformations were used to bring x , y and z atlas space maps into the patient space.

Brain extraction. In order to extract the brain and exclude other structures, such as skull, eyes, etc., we apply FSL-BET⁸³ on T1 images, because this modality has the highest resolution. The resulting mask is then transformed using registration transformation and is applied to the FLAIR images.

Bias field correction. Bias field correction is another necessary step due to magnetic field inhomogeneity. We apply FSL-FAST⁸⁴, which uses a hidden Markov random field and an associated expectation-maximization algorithm to correct for spatial intensity variations caused by RF inhomogeneities.

Intensity normalization. Apart from intensity variations caused by the bias field, intensities can also vary between patients. Thus we normalize the intensities per patient to be within the range of [0, 1].

Training, validation and test sets. From the 503 RUN DMC cases, we removed a number of cases that were extremely noisy or had failed in some of the preprocessing steps including brain extraction and registration, which left us with 420 out of 453 cases with single annotations. From 420 cases annotated by one human observer, we select 378 cases for training the model and the remaining 42 cases for validation and parameter tuning purposes and the 50 cases that were annotated by both human observers as the independent test set. It should be noted that the set of 50 images used as the test set also contained low quality images or imperfect preprocessing, however we avoided filtering any of the images out so that the experimental evaluation would better reflect the performance of the proposed method on the real (often low quality) data.

Medical datasets usually suffer from the fact that pathological observations are significantly less frequent compared to healthy observations, which also holds for our dataset. Given this, a simple uniform sampling may cause serious problems for the learning process⁸⁵, as a classifier that labels all of the samples as normal, would achieve a high accuracy. To handle this, we undersample the negative samples to create a balanced dataset. We randomly select 50% of positive and select an equal number of negative samples from normal voxels of all cases. This sampling procedure resulted in datasets consisting of 3.88 million and 430 thousand samples for training and validation sets respectively.

Methods

Patch preparation. From each voxel neighborhood, we extract patches with three different sizes: 32×32 , 64×64 and 128×128 . To reduce the computational costs, we down sample the larger two scales to 32×32 . Resulting patches for this procedure are demonstrated in Fig. 2, for a negative and a positive sample, obtained from a FLAIR image. We included these three patches for both the T1 and FLAIR modalities for each sample. This results in a set of patches in three scales s_1 , s_2 and s_3 , each consisting of two patches from T1 and FLAIR, as depicted in Fig. 3.

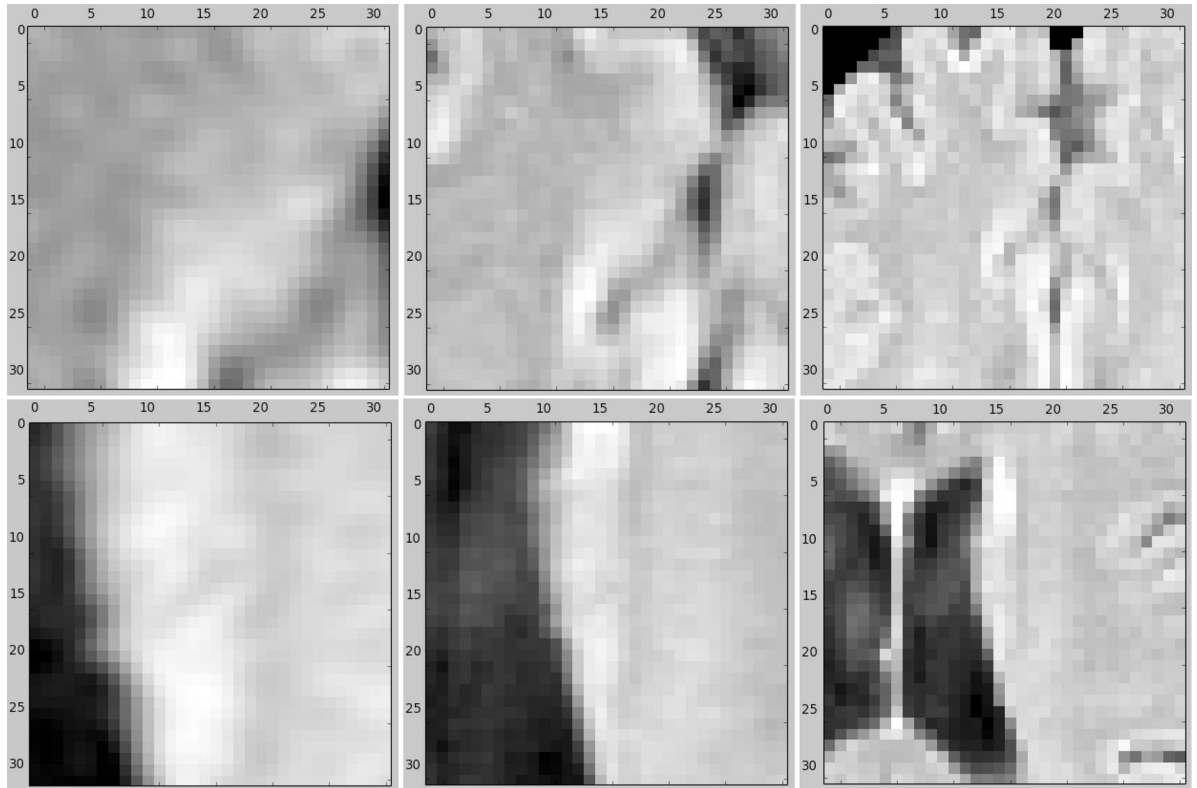


Figure 2. An example of negative (top row) and positive (bottom row) samples in three scales (from left to right) 32×32 , 64×64 and 128×128 on the FLAIR image. The two larger scales are down sampled to 32×32 .

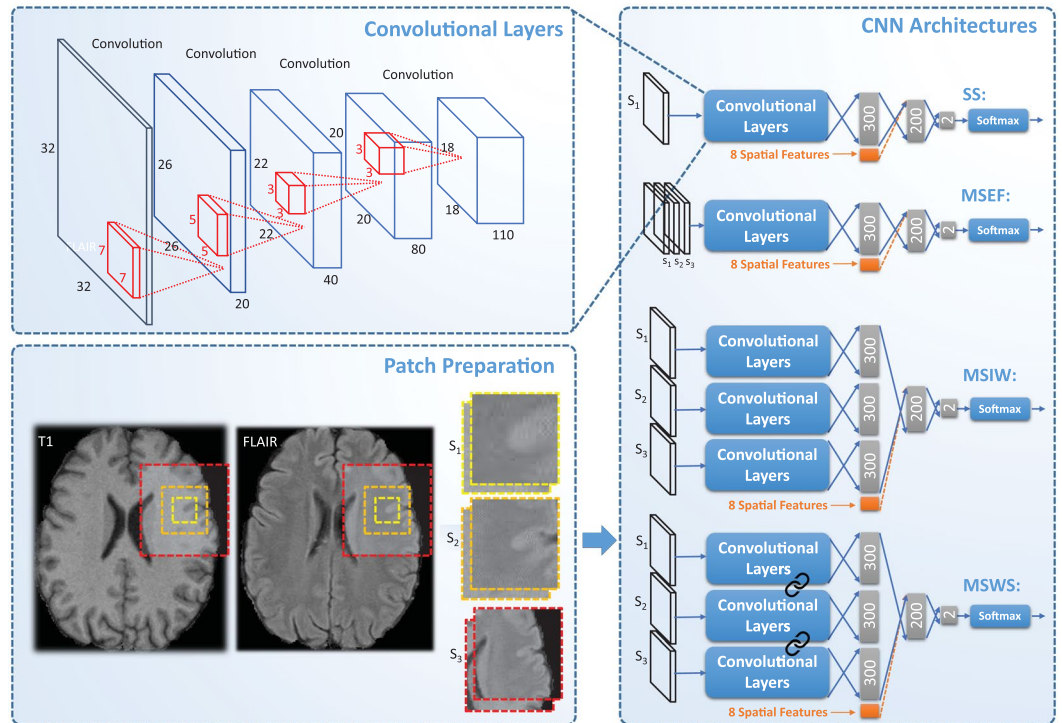


Figure 3. Patch preparation process and different proposed CNN architectures. The links between set of convolutional layers represent a weight sharing policy among the streams.

Network architectures. *Single-scale (SS) model.* The simplest CNN model we applied to our dataset was a CNN trained on patches from a single scale (with patches of 32×32). The top architecture in Fig. 3 shows the architecture of our single-scale deep CNN. This network, which is a basis for the other location sensitive architectures, consists of four convolutional layers that have 20, 40, 80 and 110 filters of size 7×7 , 5×5 , 3×3 , 3×3 respectively. We do not use pooling since it results in a shift-invariance property⁸⁶, which is not desired in segmentation tasks. Then we apply three layers of fully connected neurons of size 300, 200 and 2. Finally the resulting responses are turned into probability values using a softmax classifier.

Multi-scale early fusion (MSEF). In many cases, it is impossible to correctly classify a 32×32 patch just from its appearance. For instance, only looking at the small scale positive patch in Fig. 2, it is hard to distinguish it from cortex tissue. In contrast, given the two larger scale patches, it is fairly easy to identify it as WMH tissue near the ventricles. Furthermore there is a trade-off between context capturing and localization accuracy. Although more context information might be captured with a larger patch-size, the ability of the classifier to accurately localize the structure in the center of the patch is decreased⁶¹. This motivates a multi-scale approach that has the advantages of the smaller and larger size patches. A simple and intuitive way to train a multi-scale network is to accumulate the different scales as different channels of the input. This is possible since the larger scale patches were down sampled to 32×32 . The second top network in Fig. 3 illustrates this.

Multi-scale late fusion with independent weights (MSIW). Another possibility to create a model with multi-scale patches is to train independent convolutional layers for each scale, fusing the representations of each scale and taking them into more fully connected layers. As can be observed in Fig. 3, in this architecture each scale has its own fully connected layer. These are concatenated and fed into the joint fully connected layers. The main rationale behind giving each scale stream its own fully connected layer is that this incurs less weights compared to the approach that first merges the feature maps and then fully connects it to the first layer of neurons.

Multi-scale late fusion with weight sharing (MSWS). The first convolutional layers of a CNN typically detect various forms of edges, corners and other basic structuring elements. Since we do not expect that these basic building blocks differ much among the different scale patches, a considerable number of filters might be very similar in the three separate convolutional layers learned for different scales. Thus a potentially efficient strategy to reduce the number of weights and consequently to reduce the overfitting, is to share the convolutional filters among the different scales. As illustrated in Fig. 3, each of the scales from the different patches are separately passed through the same set of convolutional layers and each get described with separate feature maps. These feature maps are then connected to separate fully connected layers and are merged later, similar to the MSIW approach.

Integrating explicit spatial location features. The main aim for considering patches at different scales is to let the network learn about the spatial location of the samples it is observing. Alternatively we can provide the network with such information, by adding explicit features describing the spatial location. One possible place to add the location information is the first fully connected layer after the convolutional layers. All the location features are normalized per case to be within the range of $[0, 1]$. As the response of other neurons in the same layer that the location features are integrated with might have a different scale, all the eight features are scaled with a coefficient α as a parameter of the method. We tuned the best value for α as a parameter by validation. The possibility to add spatial location features is not restricted to the single-scale architecture. It is also feasible to integrate these features into the three possible architectures for multi-scale approaches. The orange parts in Fig. 3 illustrate this procedure.

There are eight features that we utilize to describe the spatial location: the x , y and z -coordinates of the corresponding voxel in the MNI atlas space, in-plane distances from the left ventricle, right ventricle, brain cortex and midsagittal brain surface as well as the prior probability of WMH occurring in that location²³.

Training procedure. For learning the network weights, we use the stochastic gradient descent algorithm⁸⁷, with mini-batch size of 128 and a cross-entropy cost. We also utilize the RMSPROP algorithm⁸⁸ to speed up the learning process by adaptively changing the learning rate for each parameter. The non-linearity applied to neurons is a rectified linear unit to prevent the vanishing gradient problem⁸⁹. As random weight initialization is important to break the symmetry between the units in the same hidden layer⁹⁰, the initial weights are drawn at random using the Glorot method⁹¹. Since CNNs are complex architectures, they are prone to overfit the data very early. Therefore we use drop-out regularization⁹² with 0.3 probability on all fully connected layers of the networks. We pick the resulting network from an epoch with the highest validation A_z as the final model.

Experimental Evaluation

For characterization of WMHs, several different methods have been proposed in this study, some of which only use patch appearance features, while others use multi-scale patches or explicit location features to the network or both. In order to obtain segmentations, we apply the trained networks to classify all the voxels inside the brain mask in a sliding window fashion. A comparison between the performance of the mentioned methods, together with a comparison to performance of an independent human observer and a conventional method with hand-crafted features would be insightful.

Integrating the location information into the first fully connected layer, as depicted in the architectures Fig. 3, is only one of the possibilities. We can alternatively add the spatial location features to one layer before or after, i.e. to the responses from the last convolutional layer and to the second fully connected layer. To evaluate the relative performance of each possibility, we also train single-scale networks with the two other possibilities and compare

Method	Without location features			With location features		
	Validation set A_z	Test set Dice (obs1)	Test set Dice (obs2)	Validation set A_z	Test set Dice (obs1)	Test set Dice (obs2)
SS	0.9939	0.731	0.729	0.9972	0.781	0.778
MSEF	0.9947	0.762	0.752	0.9966	0.777	0.769
MSIW	0.9966	0.778	0.768	0.9972	0.795	0.787
MSWS	0.9965	0.773	0.760	0.9973	0.792	0.783

Table 2. Performance comparison of different CNN architectures based on validation set A_z and test set Dice score considering observer 1 and observer 2 as the reference standard.

Method	Dice (obs1)	Dice (obs2)
Conventional	0.716	0.699
MSWS + Loc	0.792	0.783
observer 1	—	0.805
observer 2	0.805	—

Table 3. A performance comparison between conventional method, MSWS + Loc architecture, and human observers.

them to each other. In order to provide information on how much effect the dataset size has on the performance of the trained network, we present and compare the results of a MSWS + Loc network trained with 100%, 50%, 25%, 12.5% and 6.25% of the total training images.

Metrics. The Dice similarity index, also known as the Dice score, is the most widely used measure for evaluating the agreement between different segmentation methods and their reference standard segmentations^{43,44}. It is computed as

$$Dice = \frac{2 \times TP}{FP + FN + 2 \times TP} \quad (1)$$

where the value varies between 0 for complete disagreement, and 1 representing complete agreement between the reference standard and the evaluated segmentation. A Dice similarity index of 0.7 or higher is usually considered a good segmentation in the literature⁴⁵. To create binary masks out of probability maps resulting from CNNs, we find an optimal value as a threshold that maximizes the overall Dice score on the validation set. The optimal thresholds are computed separately for each method. We also present test set receiver operating characteristic (ROC) curves and validation set area under the ROC curve (A_z). For computing each of these measures, we only consider the voxels inside the brain mask, to avoid taking easy voxels belonging to the background into account.

For the statistical significance test, we created a 100 boot-straps by sampling 50 instances with replacement. Then the Dice scores were computed on each bootstrap for each of the two compared methods. Empirical p-values were reported as the proportion of bootstraps where the Dice score for method B was higher than A, when the null-hypothesis to reject was “method A is no better than B”. If no such bootstrap existed, the p-value < 0.01 was reported, representing a significant difference.

Conventional segmentation system. In order to evaluate the relative performance of the proposed deep learning systems, we also train a conventional segmentation system, using hand-crafted features²³. The set of hand-crafted features consists of 22 features in total: intensity features including FLAIR and T1 intensities, second order derivative features including multi-scale Laplacian of Gaussian ($\sigma = 1, 2, 4$ mm), multi-scale determinant of Hessian ($t = 1, 2, 4$ mm), vesselness filter ($\sigma = 1$ mm), a multi-scale annular filter ($t = 1, 2, 4$ mm), FLAIR intensity mean and standard deviation in a 16×16 neighborhood, as well as the same 8 location features that were used in the previous subsection. We use a random forest classifier with 50 subtrees to train the model.

Experimental Results

Table 2 represents a comparison on validation set A_z and test set Dice score, for each of the methods, once without and another time with addition of spatial location features, considering observer 1 as the reference standard. Table 3 compares the performance of the conventional segmentation method, our late fusion multi-scale architecture with weight sharing and location information (MSWS + Loc), and the two human observers on the independent test set, with each observer as the reference standard. P-values were computed as a result of patient-level boot-strapping on the test set and are presented in Table 4.

Regarding the different options for integration of the location information in the network, Table 5 compares the performance of these options on the validation and training sets.

Figure 4(a) shows the ROC curves for some of the trained CNN architectures and compares them to the conventional segmentation method and the independent human observer. The ROC curves have been cut to show only low false positive rates that are of interest for practical use. In order to preserve readability of the figures, we

Method	MSWS	SS + Loc	MSWS + Loc	Ind. Obs.
SS	<0.01	<0.01	<0.01	<0.01
MSWS	—	<0.01	<0.01	<0.01
SS + Loc	—	—	0.03	0.03
MSWS + Loc	—	—	—	0.06

Table 4. Statistical significance test for pairwise comparison of the methods Dice score. p_{ij} indicates the p-value for the null hypothesis that method i is better than method j .

Method	Validation set A_2	Test set Dice
LCL	0.9964	0.763
FFCL	0.9971	0.781
SFCL	0.9967	0.778

Table 5. A performance comparison of the single-scale architecture with different possible locations to add the spatial location information. Abbreviations: last convolutional layer (LCL), first fully connected layer (FFCL), second fully connected layer (SFCL).

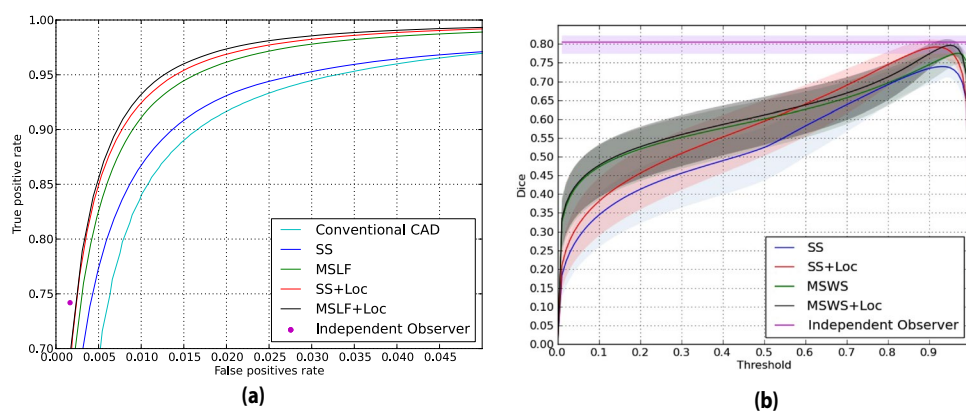


Figure 4. Integration of spatial location information fills the gap between performance of a normal CNN and human observer. (a) An ROC comparison of different CNN methods, a conventional segmentation method and independent human observer, considering observer 1 as the reference standard. (b) A comparison of different methods on Dice score as a function of binary masking threshold. The light shades around the curves indicate 95% confidence intervals with bootstrapping on patients.

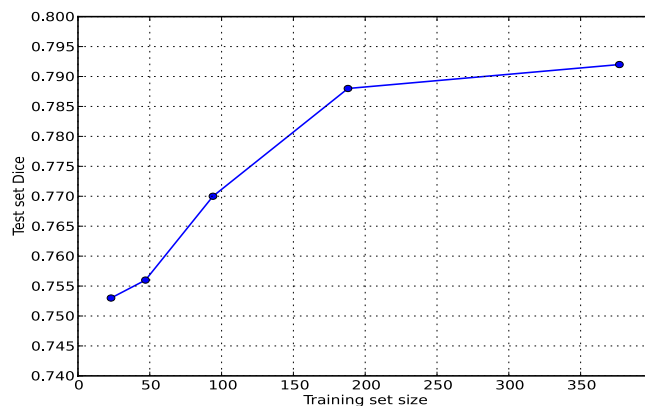


Figure 5. Test Dice as a function of training set size.

only compare the most informative methods. Figure 5(b) shows the Dice similarity scores as a function of the binary masking threshold. It also compares them to the Dice similarity measure between the two human observers. 95% confidence intervals are depicted for each curve, as a result of bootstrapping on patients. The effect of the training dataset size can be observed in Table 6 and Fig. 5.

Training set size	23	47	94	189	378
Test set Dice	0.753	0.756	0.770	0.788	0.792

Table 6. Test Dice as a function of training set size.

Discussion

Contribution of larger context and location information. Comparing the performance of the SS and SS + Loc approaches, as presented in the first row of Table 2, a significant difference in Dice score is observable (p -value < 0.01). This points us to the fact that a knowledge about where the input patch is located can substantially improve WMH segmentation quality of a CNN. A similar significant difference is observable when comparing performance measures of SS and MSWS methods (p -value < 0.01). This implies that by using a multi-scale approach, a CNN can learn about context information quite well. Considering the better performance of SS + Loc compared to MSWS, we can infer that the learning of location and large scale context from multi-scale patches is not as good as adding explicit location information to the architecture.

Early fusion vs. late fusion, independent weights vs. weight sharing. As the experimental results suggest, among the different multi-scale fusion architectures, early fusion shows the least improvement over the single-scale approach. The related patch voxels of different scales, do not have a meaningful correspondence. Given the fact that the convolution operation in the first convolutional layer sums up the responses on each scale, we assume that the useful information provided by different scales is washed out too early in the network. In contrast, the two late fusion architectures show comparable good performance, however in general, since the late fusion architecture with weight sharing is a simpler model with less parameters to be learned, one might prefer to use this model.

Comparison to human observer and a conventional method. Shown by Table 3, MSWS + Loc substantially outperforms a conventional segmentation method, with Dice score of 0.792 compared to 0.716 (p -value < 0.01). Furthermore, the Dice score of MSWS + Loc method closely resembles the inter-observer variability, which implies that the segmentation provided by MSWS + Loc approach is as good as the two human observers. Also the statistical test does not show a significant advantage of the independent observer compared to this method (p -value = 0.06).

A visual look into the results. Figures 6–8 show some qualitative examples. Figure 6 contains two sample cases, where the location and larger context information leads to a better segmentation. As evident from the first sample, the single-scale CNN falsely segments an area on septum pellucidum, which also appears as hyperintense tissue. These false positives can be avoided by considering location information. A second sample shows improvements on FNs of the single-scale method.

Figure 7 illustrates an instance of a prevalent class of false positives of the system, which are the hyperintense voxels around the lacunes. Since the model has not been trained on so many negative samples similar to this, the distinction between WMH and hyperintensities around lacunes is not well learned by the system. An obvious solution is to extensively include the lacunes surrounding voxels as negative samples in the training dataset.

As an example of missed lesions by human observers, Fig. 8 shows a small lesion on the right temporal lobe, missed by both human observers, where it is detected by MSWS + Loc method. Another sample of such missed lesions can be observed in the second sample of Fig. 6, on the right hemisphere frontal lobe. Based on similar observations, we can assume that some of the false positives are possibly small lesions missed by one or both of the observers. Therefore there may be a chance that the real performance of the system is better than reported, but it would require more research to investigate this.

Integration of location features. For integration of explicit spatial location information into the CNN, there are several possibilities that were investigated in this study. The results as represented in Table 3, suggest that adding the spatial location features to the first fully connected layer results in a significantly better performance. Adding them to around 35 K features as the responses of the last convolutional layer, almost makes the eight location features insignificant among so many representation features. At the other extreme, although integrating the location features into the second fully connected layer does not suffer from this problem, but leaves less flexibility for the network to consider location features for the discrimination to be learned. The first fully connected layer seems to be the best option, where the appearance features provided by the last convolutional layer are already considerably reduced, and at same time the more fully connected layer provides more flexibility for an optimal discrimination.

Two-stage vs. single-stage model. As shown in the results, integrating location information into a CNN can play an important role in obtaining an accurate segmentation. We integrate the features while we train our network to learn the representations. Another approach is to perform this task in two stages; first training an independent network that learns the representations, and later training a second classifier that takes the output features of the first network, integrated with location or other external features (as followed in ref. 93 for instance). The first approach, which is followed in this study, seems more reasonable as the set of learned filters without location information could differ from the optimal set of filters given the location information. The two-stage system lacks this information and might devote some of the filters for capturing of location that are redundant given the location features.

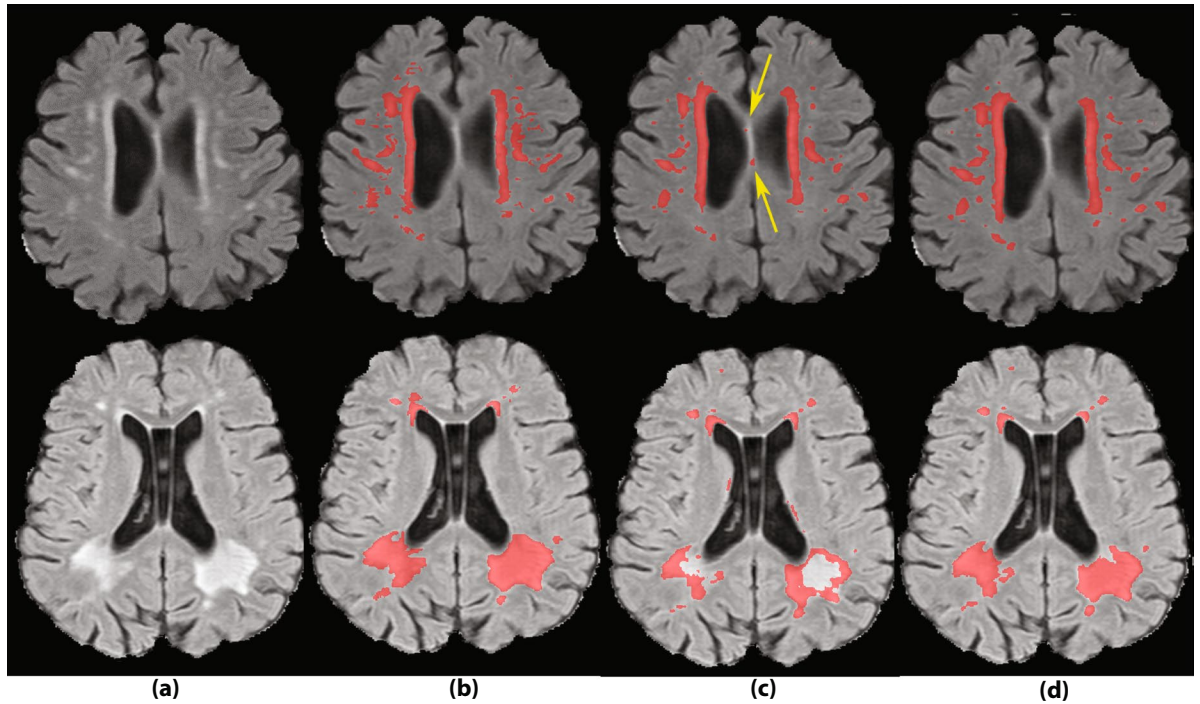


Figure 6. Two sample cases of segmentation improvement by adding location information to the network. (a) FLAIR images without annotations. (b) Segmentation by human observer 1. (c) Segmentation by SS method. (d) Segmentation by MSWS + Loc method.

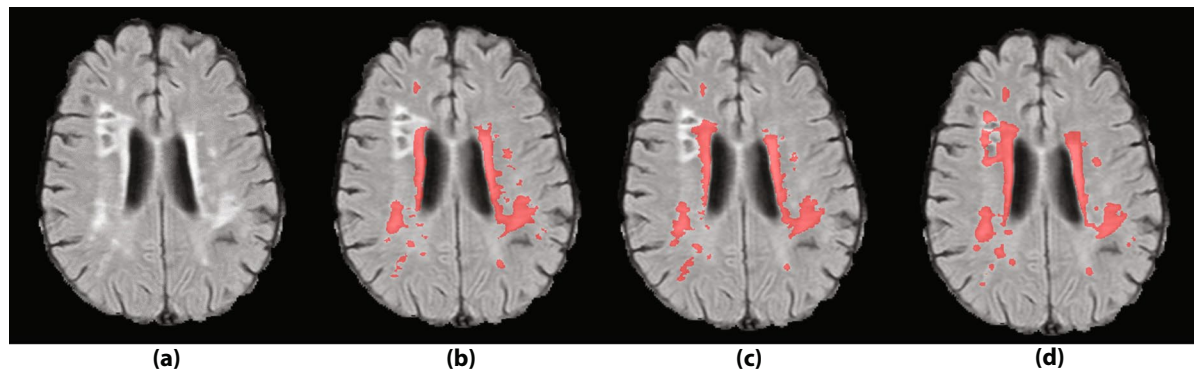


Figure 7. Gliosis around the lacunes is a prevalent type of false positive segmentation. (a) FLAIR images without annotations. (b) Segmentation by human observer 1. (c) Segmentation by human observer 2. (d) Segmentation by MSWS + Loc method.

2D vs. 3D patches. In this research, we sample 2D patches from each of the two modalities (T1 and FLAIR), while one might argue that considering consecutive slices and sampling 3D patches from each image modality could provide useful information. Given the slice thickness of 5 mm with a 1 mm inter-slice gap in our dataset, the consecutive slices do not highly correspond to each other. Furthermore incorporation of 3D patches extensively increases the computational costs at both the training and the segmentation time. These motivated us to use 2D patches. In contrast, for datasets with isotropic or thin slice FLAIR images, 3D patches might be very useful.

Fully convolutional segmentation network. Fully convolutional networks replace the fully connected layers with 1×1 convolutions that perform exactly the same functionality as the fully connected do, however implemented with convolutions⁶⁰. This would speed the segmentation up, since convolutions can get larger input images, make dense predictions for the whole input image and avoid repetitive computations. While we have trained our networks in a patch-based manner, it does not restrict us from reforming the fully connected layers of the trained network into convolutional layer counterparts at the segmentation time. The current implementation uses a patch-based segmentation, as we found it fast enough in the current experimental setup (~3 minutes for the multi-scale and ~1.5 minutes for the single-scale architectures per case on a Titan X card). It should also be noticed that a patch-based *training*, compared to the fully convolutional training, has the advantages that it can be

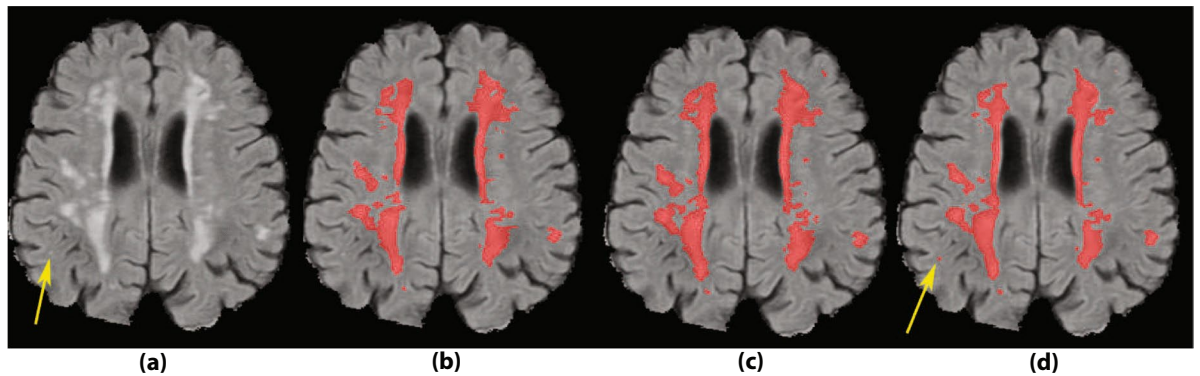


Figure 8. A sample case with a small lesion missed by the two human observers. (a) FLAIR image without annotations. (b) Segmentation by human observer 1. (c) Segmentation by human observer 2. (d) Segmentation by MSWS + Loc method.

much less memory demanding and is easier to optimize in highly imbalanced classification problems due to the possibility of the class-specific data augmentation.

Conclusions

In this study we showed that location information can have a significant added value when using CNNs for WMH segmentation. While for this task, making use of CNNs, not only a better performance compared to conventional segmentation method was achieved, we approached the performance level of an independent human observer with incorporation of location information.

References

- van Norden, A. G. *et al.* Causes and consequences of cerebral small vessel disease. The RUN DMC study: a prospective cohort study. Study rationale and protocol. *BMC Neurol* **11**, 29 (2011).
- Schoonheim, M. M. *et al.* Sex-specific extent and severity of white matter damage in multiple sclerosis: Implications for cognitive decline. *Human Brain Mapping* **35**, 2348–2358 (2014).
- Marshall, G., Shchelchkov, E., Kaufer, D., Ivanco, L. & Bohnen, N. White matter hyperintensities and cortical acetylcholinesterase activity in parkinsonian dementia. *Acta Neurologica Scandinavica* **113**, 87–91 (2006).
- Weinstein, G. *et al.* Brain imaging and cognitive predictors of stroke and alzheimer disease in the framingham heart study. *Stroke* **44**, 2787–2794 (2013).
- Hirono, N., Kitagaki, H., Kazui, H., Hashimoto, M. & Mori, E. Impact of white matter changes on clinical manifestation of alzheimer's disease a quantitative study. *Stroke* **31**, 2182–2188 (2000).
- Smith, C. D., Snowdon, D. A., Wang, H. & Markesbery, W. R. White matter volumes and periventricular white matter hyperintensities in aging and dementia. *Neurology* **54**, 838–842 (2000).
- Wardlaw, J. M. *et al.* Neuroimaging standards for research into small vessel disease and its contribution to ageing and neurodegeneration. *The Lancet Neurology* **12**, 822–838 (2013).
- De Leeuw, F. *et al.* Prevalence of cerebral white matter lesions in elderly people: a population based magnetic resonance imaging study. the rotterdam scan study. *Journal of Neurology, Neurosurgery & Psychiatry* **70**, 9–14 (2001).
- de Groot, J. C. *et al.* Cerebral white matter lesions and cognitive function: the rotterdam scan study. *Annals of Neurology* **47**, 145–151 (2000).
- Au, R. *et al.* Association of white matter hyperintensity volume with decreased cognitive functioning: the framingham heart study. *Archives of Neurology* **63**, 246–250 (2006).
- Whitman, G., Tang, T., Lin, A. & Baloh, R. A prospective study of cerebral white matter abnormalities in older people with gait dysfunction. *Neurology* **57**, 990–994 (2001).
- Firbank, M. J. *et al.* Brain atrophy and white matter hyperintensity change in older adults and relationship to blood pressure. *Journal of Neurology* **254**, 713–721 (2007).
- Herrmann, L. L., Le Masurier, M. & Ebmeier, K. P. White matter hyperintensities in late life depression: a systematic review. *Journal of Neurology, Neurosurgery & Psychiatry* **79**, 619–624 (2008).
- van Uden, I. W. *et al.* White matter integrity and depressive symptoms in cerebral small vessel disease: The run dmc study. *The American Journal of Geriatric Psychiatry* **23**, 525–535 (2015).
- Van Straaten, E. C. *et al.* Impact of white matter hyperintensities scoring method on correlations with clinical data the ladis study. *Stroke* **37**, 836–840 (2006).
- Polman, C. H. *et al.* Diagnostic criteria for multiple sclerosis: 2005 revisions to the “mcdonald criteria”. *Annals of Neurology* **58**, 840–846 (2005).
- Grimaud, J. *et al.* Quantification of mri lesion load in multiple sclerosis: a comparison of three computer-assisted techniques. *Magnetic Resonance Imaging* **14**, 495–505 (1996).
- Anbeek, P., Vincken, K. L., van Osch, M. J., Bisschops, R. H. & van der Grond, J. Probabilistic segmentation of white matter lesions in mr imaging. *NeuroImage* **21**, 1037–1044 (2004).
- Lao, Z. *et al.* Computer-assisted segmentation of white matter lesions in 3d mr images using support vector machine. *Academic Radiology* **15**, 300–313 (2008).
- Herskovits, E., Bryan, R. & Yang, F. Automated bayesian segmentation of microvascular white-matter lesions in the accord-mind study. *Advances in Medical Sciences* **53**, 182–190 (2008).
- Simões, R. *et al.* Automatic segmentation of cerebral white matter hyperintensities using only 3d flair images. *Magnetic Resonance Imaging* **31**, 1182–1189 (2013).
- Ithapu, V. *et al.* Extracting and summarizing white matter hyperintensities using supervised segmentation methods in alzheimer's disease risk and aging studies. *Human Brain Mapping* **35**, 4219–4235 (2014).
- Ghafoorian, M. *et al.* Small white matter lesion detection in cerebral small vessel disease. *SPIE Medical Imaging* **9414**, 941411–941411 (2015).

24. Klöppel, S. *et al.* A comparison of different automated methods for the detection of white matter lesions in mri data. *NeuroImage* **57**, 416–422 (2011).
25. Zijdenbos, A. P., Forghani, R. & Evans, A. C. Automatic “pipeline” analysis of 3-d mri data for clinical trials: application to multiple sclerosis. *Medical Imaging, IEEE Transactions on* **21**, 1280–1291 (2002).
26. Dyrby, T. B. *et al.* Segmentation of age-related white matter changes in a clinical multi-center study. *Neuroimage* **41**, 335–345 (2008).
27. Geremia, E. *et al.* Spatial decision forests for ms lesion segmentation in multi-channel magnetic resonance images. *NeuroImage* **57**, 378–390 (2011).
28. Ghafoorian, M. *et al.* Automated detection of white matter hyperintensities of all sizes in cerebral small vessel disease. *Medical Physics* **43** (2016).
29. Ghafoorian, M. *et al.* Transfer learning for domain adaptation in mri: Application in brain lesion segmentation. *arXiv preprint arXiv:1702.07841* (2017).
30. Vijverberg, K. *et al.* A single-layer network unsupervised feature learning method for white matter hyperintensity segmentation. In *SPIE Medical Imaging, 97851C–97851C* (International Society for Optics and Photonics, 2016).
31. Brosch, T. *et al.* Deep 3d convolutional encoder networks with shortcuts for multiscale feature integration applied to multiple sclerosis lesion segmentation. *IEEE transactions on medical imaging* **35**, 1229–1239 (2016).
32. Brosch, T. *et al.* Deep convolutional encoder networks for multiple sclerosis lesion segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015*, 3–11 (Springer, 2015).
33. Ghafoorian, M. *et al.* Non-uniform patch sampling with deep convolutional neural networks for white matter hyperintensity segmentation. In *International Symposium on Biomedical Imaging (ISBI)*, 1414–1417 (IEEE, 2016).
34. Kamnitsas, K. *et al.* Efficient multi-scale 3d cnn with fully connected crf for accurate brain lesion segmentation. *arXiv preprint arXiv:1603.05959* (2016).
35. Van Leemput, K., Maes, F., Vandermeulen, D., Colchester, A. & Suetens, P. Automated segmentation of multiple sclerosis lesions by model outlier detection. *Medical Imaging, IEEE Transactions on* **20**, 677–688 (2001).
36. Shi, L. *et al.* Automated quantification of white matter lesion in magnetic resonance imaging of patients with acute infarction. *Journal of Neuroscience Methods* **213**, 138–146 (2013).
37. Khademi, A., Venetsanopoulos, A. & Moody, A. R. Robust white matter lesion segmentation in flair mri. *Biomedical Engineering, IEEE Transactions on* **59**, 860–871 (2012).
38. Admiraal-Behloul, F. *et al.* Fully automatic segmentation of white matter hyperintensities in mr images of the elderly. *Neuroimage* **28**, 607–617 (2005).
39. de Boer, R. *et al.* White matter lesion extension to automatic brain tissue segmentation on mri. *Neuroimage* **45**, 1151–1161 (2009).
40. Jain, S. *et al.* Automatic segmentation and volumetry of multiple sclerosis brain lesions from mr images. *NeuroImage: Clinical* **8**, 367–375 (2015).
41. Shiee, N. *et al.* A topology-preserving approach to the segmentation of brain images with multiple sclerosis lesions. *NeuroImage* **49**, 1524–1535 (2010).
42. Schmidt, P. *et al.* An automated tool for detection of flair-hyperintense white-matter lesions in multiple sclerosis. *Neuroimage* **59**, 3774–3783 (2012).
43. Caligiuri, M. E. *et al.* Automatic detection of white matter hyperintensities in healthy aging and pathology using magnetic resonance imaging: A review. *Neuroinformatics* **13**, 1–16 (2015).
44. Garca-Lorenzo, D., Francis, S., Narayanan, S., Arnold, D. L. & Collins, D. L. Review of automatic segmentation methods of multiple sclerosis white matter lesions on conventional magnetic resonance imaging. *Medical Image Analysis* **17**, 1–18 (2013).
45. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
46. Schmidhuber, J. Deep learning in neural networks: An overview. *Neural Networks* **61**, 85–117 (2015).
47. Hubel, D. H. & Wiesel, T. N. Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex. *The Journal of Physiology* **160**, 106 (1962).
48. Fukushima, K. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological cybernetics* **36**, 193–202 (1980).
49. Cireşan, D., Meier, U., Masci, J. & Schmidhuber, J. Multi-column deep neural network for traffic sign classification. *Neural Networks* **32**, 333–338 (2012).
50. He, K., Zhang, X., Ren, S. & Sun, J. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *arXiv preprint arXiv:1502.01852* (2015).
51. Taigman, Y., Yang, M., Ranzato, M. & Wolf, L. Deepface: Closing the gap to human-level performance in face verification. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, 1701–1708 (2014).
52. Cireşan, D., Giusti, A., Gambardella, L. M. & Schmidhuber, J. Deep neural networks segment neuronal membranes in electron microscopy images. In *Advances in Neural Information Processing Systems*, 2843–2851 (2012).
53. Cireşan, D. & Schmidhuber, J. Multi-column deep neural networks for offline handwritten chinese character classification. *arXiv preprint arXiv:1309.0261* (2013).
54. LeCun, Y., Bottou, L., Bengio, Y. & Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**, 2278–2324 (1998).
55. Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, 1097–1105 (2012).
56. Deng, J. *et al.* Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, 248–255 (IEEE, 2009).
57. Farabet, C., Couprie, C., Najman, L. & LeCun, Y. Learning hierarchical features for scene labeling. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **35**, 1915–1929 (2013).
58. Gupta, S., Girshick, R., Arbeláez, P. & Malik, J. Learning rich features from rgb-d images for object detection and segmentation. In *Computer Vision–ECCV 2014, Lecture Notes in Computer Science (LNCS 8695)*, 345–360 (Springer, 2014).
59. Hariharan, B., Arbeláez, P., Girshick, R. & Malik, J. Simultaneous detection and segmentation. In *Computer Vision–ECCV 2014, Lecture Notes in Computer Science (LNCS 8695)*, 297–312 (Springer, 2014).
60. Long, J., Shelhamer, E. & Darrell, T. Fully convolutional networks for semantic segmentation. *arXiv preprint arXiv:1411.4038* (2014).
61. Ronneberger, O., Fischer, P. & Brox, T. U-net: Convolutional networks for biomedical image segmentation. *arXiv preprint arXiv:1505.04597* (2015).
62. Litjens, G. *et al.* A survey on deep learning in medical image analysis. *arXiv preprint arXiv:1702.05747* (2017).
63. Greenspan, H., van Ginneken, B. & Summers, R. M. Guest editorial deep learning in medical imaging: Overview and future promise of an exciting new technique. *IEEE Transactions on Medical Imaging* **35**, 1153–1159 (2016).
64. Kleesiek, J. *et al.* Deep mri brain extraction: a 3d convolutional neural network for skull stripping. *NeuroImage* **129**, 460–469 (2016).
65. Zhang, W. *et al.* Deep convolutional neural networks for multi-modality iso-intense infant brain image segmentation. *NeuroImage* **108**, 214–224 (2015).
66. Moeskops, P. *et al.* Automatic segmentation of mr brain images with a convolutional neural network. *IEEE transactions on medical imaging* **35**, 1252–1261 (2016).
67. Milletari, F. *et al.* Hough-cnn: Deep learning for segmentation of deep brain regions in mri and ultrasound. *arXiv preprint arXiv:1601.07014* (2016).

68. Chen, H., Dou, Q., Yu, L. & Heng, P.-A. Voxresnet: Deep voxelwise residual networks for volumetric brain segmentation. *arXiv preprint arXiv:1608.05895* (2016).
69. Nie, D., Wang, L., Gao, Y. & Sken, D. Fully convolutional networks for multi-modality isointense infant brain image segmentation. In *Biomedical Imaging (ISBI), 2016 IEEE 13th International Symposium on*, 1342–1345 (IEEE, 2016).
70. Shakeri, M. *et al.* Sub-cortical brain structure segmentation using f-cnn's. *arXiv preprint arXiv:1602.02130* (2016).
71. Pereira, S., Pinto, A., Alves, V. & Silva, C. A. Brain tumor segmentation using convolutional neural networks in mri images. *IEEE transactions on medical imaging* **35**, 1240–1251 (2016).
72. Havaei, M. *et al.* Brain tumor segmentation with deep neural networks. *Medical Image Analysis* (2016).
73. Havaei, M., Guizard, N., Chapados, N. & Bengio, Y. Hemis: Hetero-modal image segmentation. *arXiv preprint arXiv:1607.05194* (2016).
74. Zhao, L. & Jia, K. Multiscale cnns for brain tumor segmentation and diagnosis. *Computational and mathematical methods in medicine* **2016** (2016).
75. Ghafoorian, M. *et al.* Deep multi-scale location-aware 3d convolutional neural networks for automated detection of lacunes of presumed vascular origin. *NeuroImage: Clinical* **14**, 391–399 (2017).
76. Dou, Q. *et al.* Automatic cerebral microbleeds detection from mr images via independent subspace analysis based hierarchical features. In *2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, 7933–7936 (IEEE, 2015).
77. Dou, Q. *et al.* Automatic detection of cerebral microbleeds from mr images via 3d convolutional neural networks. *IEEE transactions on medical imaging* **35**, 1182–1195 (2016).
78. Kamber, M., Shinghal, R., Collins, D. L., Francis, G. S. & Evans, A. C. Model-based 3-d segmentation of multiple sclerosis lesions in magnetic resonance brain images. *Medical Imaging, IEEE Transactions on* **14**, 442–453 (1995).
79. Hervé, D., Mangin, J.-F., Molko, N., Bousser, M.-G. & Chabriat, H. Shape and volume of lacunar infarcts a 3d mri study in cerebral autosomal dominant arteriopathy with subcortical infarcts and leukoencephalopathy. *Stroke* **36**, 2384–2388 (2005).
80. Jenkinson, M. & Smith, S. A global optimisation method for robust affine registration of brain images. *Medical Image Analysis* **5**, 143–156 (2001).
81. Mazziotta, J. *et al.* A four-dimensional probabilistic atlas of the human brain. *Journal of the American Medical Informatics Association* **8**, 401–430 (2001).
82. Jenkinson, M., Beckmann, C. F., Behrens, T. E., Woolrich, M. W. & Smith, S. M. Fsl. *Neuroimage* **62**, 782–790 (2012).
83. Smith, S. M. Fast robust automated brain extraction. *Human Brain Mapping* **17**, 143–155 (2002).
84. Zhang, Y., Brady, M. & Smith, S. Segmentation of brain mr images through a hidden markov random field model and the expectation-maximization algorithm. *Medical Imaging, IEEE Transactions on* **20**, 45–57 (2001).
85. Pastor-Pellicer, J., Zamora-Martinez, F., España-Boquera, S. & Castro-Bleda, M. J. F-measure as the error function to train neural networks. In *Advances in Computational Intelligence*, Lecture Notes in Computer Science (LNCS 7902), 376–384 (Springer, 2013).
86. Scherer, D., Müller, A. & Behnke, S. Evaluation of pooling operations in convolutional architectures for object recognition. In *Artificial Neural Networks—ICANN 2010*, Lecture Notes in Computer Science (LNCS 6354), 92–101 (Springer, 2010).
87. Bottou, L. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, 177–186 (Springer, 2010).
88. Dauphin, Y. N., de Vries, H., Chung, J. & Bengio, Y. Rmsprop and equilibrated adaptive learning rates for non-convex optimization. *arXiv preprint arXiv:1502.04390* (2015).
89. Maas, A. L., Hannun, A. Y. & Ng, A. Y. Rectifier nonlinearities improve neural network acoustic models. In *Proc. ICML*, vol. 30 (2013).
90. Bottou, L. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade*, Lecture Notes in Computer Science (LNCS 7700), 421–436 (Springer, 2012).
91. Glorot, X. & Bengio, Y. Understanding the difficulty of training deep feedforward neural networks. In *Aistats*, vol. 9, 249–256 (2010).
92. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* **15**, 1929–1958 (2014).
93. Kooi, T. *et al.* Large scale deep learning for computer aided detection of mammographic lesions. *Medical Image Analysis* **35**, 303–312 (2017).

Acknowledgements

This work was supported by a VIDI innovational grant from the Netherlands Organisation for Scientific Research (NWO, grant 016.126.351). The authors also would like to acknowledge Lucas J.B. van Oudheusden and Koen Vijverberg for their contributions to this study.

Author Contributions

M.G. performed the experiments and wrote the manuscript. B.P., E.M., N.K., T.H., and B.v.G. designed the study and the experiments and supervised the project. I.v.U. and F.E.d.L. provided the dataset and the manual annotations. G.L. and C.S. provided comments and feedback on the study and the results. All the authors reviewed the manuscript.

Additional Information

Competing Interests: The authors declare that they have no competing interests.

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2017