

# Shallow vs deep learning architectures for white matter lesion segmentation in the early stages of multiple sclerosis

Francesco La Rosa<sup>1,3</sup>, Mário João Fartaria<sup>1,2,4</sup>, Tobias Kober<sup>1,2,4</sup>, Jonas Richiardi<sup>2,4</sup>, Cristina Granziera<sup>5,6</sup>, Jean-Philippe Thiran<sup>1,4</sup>, and Meritxell Bach Cuadra<sup>1,3,4</sup>

<sup>1</sup> LTS5, Ecole Polytechnique Fédérale de Lausanne, Switzerland

<sup>2</sup> Siemens Healthcare AG, Lausanne, Switzerland

<sup>3</sup> Medical Image Analysis Laboratory, CIBM, University of Lausanne, Switzerland

<sup>4</sup> Radiology Department, Lausanne University Hospital, Switzerland

<sup>5</sup> Translational Imaging in Neurology Basel, Department of Medicine and Biomedical Engineering, University Hospital Basel and University of Basel, Basel, Switzerland

<sup>6</sup> Neurologic Clinic and Policlinic, Departments of Medicine, Clinical Research and Biomedical Engineering, University Hospital Basel and University of Basel, Basel, Switzerland

**Abstract.** In this work, we present a comparison of a shallow and a deep learning architecture for the automated segmentation of white matter lesions in MR images of multiple sclerosis patients. In particular, we train and test both methods on early stage disease patients, to verify their performance in challenging conditions, more similar to a clinical setting than what is typically provided in multiple sclerosis segmentation challenges. Furthermore, we evaluate a prototype naive combination of the two methods, which refines the final segmentation. All methods were trained on 32 patients, and the evaluation was performed on a pure test set of 73 cases. Results show low lesion-wise false positives (30%) for the deep learning architecture, whereas the shallow architecture yields the best Dice coefficient (63%) and volume difference (19%). Combining both shallow and deep architectures further improves the lesion-wise metrics (69% and 26% lesion-wise true and false positive rate, respectively).

## 1 Introduction

Multiple Sclerosis (MS) is a demyelinating disease that affects the central nervous system. Demyelination results in focal lesions that appear with higher frequency in the white matter (WM). Magnetic resonance imaging (MRI) is a fundamental tool for MS diagnosis and monitoring of disease evolution as well as response to therapy. Currently, expert’s manual annotations are considered the clinical gold standard for MS lesion identification. However, as this task is time-consuming and prone to inter and intra-observer variations, many automated methods for MS lesion detection and segmentation have been proposed in the literature [1]. In this context, supervised techniques that learn and train from manually annotated

examples have been proven to be the most successful in detection of MS WM lesions [2, 3, 4]. In the last years, deep learning architectures have achieved remarkable successes and have recently proven good performance in MS lesion segmentation as well [5, 6, 7].

In order to compare automated lesion segmentation methods, several computational imaging challenges have been proposed at international conferences [2, 3, 4], providing very valuable benchmark datasets for validation. However, these evaluation scenarios are based on patients with relatively high lesion load, and reported results are often computed on scans exhibiting relative large lesion sizes. Thus, the performance of deep supervised techniques on early stages of MS and at small lesion sizes remains to be proven.

In this work, we aim at comparing shallow with novel deep learning architectures using data from early stages of the disease in challenging conditions, i.e. exploring minimum lesion sizes as given by neuroradiological conventions [8] and even pushing the limit below.

To this end, we have selected two recently published MS segmentation methods. First, we have applied a supervised k-NN method combined with partial volume (PV) modeling [9, 10], specifically developed on subjects with a low disease burden and small lesions. Second, we have used a recently and publicly available deep learning approach based on a cascade of two 3D patch-wise convolutional neural networks (CNNs) [5]. At the time of the writing of this work, this CNNs method achieved the best result on the MICCAI2008 and MSSEG2016 challenges [2, 4] and competitive performance on other clinical datasets. Furthermore, we explore a straightforward prototype combination of these two methods. Both methods and their combination are trained on the same clinical dataset and validated on a pure test set. The results are analyzed considering different minimum lesion volume and total lesion load, as these are important evidences for early stage disease patients with low disabilities.

## 2 Methodology

### 2.1 Datasets

The study was approved by the Ethics Committee of our institution, and all patients gave written informed consent prior to participation. The training dataset was composed of 32 patients, 18 female / 14 male, mean age  $34 \pm 10$  years, with Expanded Disability Status Scale (EDSS) scores ranged from 1 to 2 (mean  $1.6 \pm 0.3$ ). Mean lesion volume is  $0.11 \pm 0.40$  ml (range 0.001-7.03 ml). Mean lesion load per case was  $6.0 \pm 7.2$  ml (range 0.3-37.2 ml). MRI acquisitions were performed on a 3T MRI scanner (Magnetom Trio, Siemens Healthcare, Erlangen, Germany). Both 3D MPRAGE and 3D FLAIR were acquired with a resolution of  $1 \times 1 \times 1.2$  mm<sup>3</sup>.

The test dataset was made up of 73 patients, 50 females and 23 males (mean age  $38 \pm 10$  years). EDSS scores ranged from 1 to 7.5 (mean  $2.6 \pm 1.5$ ). Mean lesion volume was  $0.25 \pm 3.29$  ml (range 0.002-159.827 ml). Mean lesion load

per case was  $14.3 \pm 27.9$  ml (range 0.2-162.9 ml). Both 3D MPRAGE and 3D FLAIR were acquired at  $1 \times 1 \times 1$  mm<sup>3</sup> but with different Siemens scanners: 5 subjects at 1.5T with MAGNETOM Aera, and the other patients at 3T with either Prisma\_fit, TrioTim, or Skyra systems.

**Manual segmentation:** In the training set, MS lesions were detected by consensus by one radiologist and one neurologist, with respectively 6 and 11 years of experience. The lesion volumes were then delineated in each image by a trained technician. Testing set lesions segmentation was performed by the Medical Image Analysis Center-MIAC [11] based on a standardized semi-automated method and further experts quality check, which has been extensively applied to phase II and III clinical trials.

## 2.2 Pre-processing

The same pre-processing steps were applied to the training and testing datasets. First, the two image contrasts were rigidly registered to the same space (MPRAGE) using the ELASTIX C++ library [12]. Second, all cases were skull-stripped using BET [13] and bias-corrected using N4 [14, 15].

## 2.3 LeMan-PV

LeMan-PV is a Bayesian PV estimation (PVE) algorithm, where spatial constraints for GM and lesions are included to drive the segmentation [10]. The spatial constraint for GM is an atlas-based probability map, and spatial constraints for lesions are derived from a kNN-supervised-based approach [16]. LeMan-PV has proven its good performance, and improvements as compared to state-of-the-art methods, in a leave-one-out experiments with MS patients with low lesion loads and small lesions. As in [9], initial mean tissue intensities and hyperparameters (symmetric penalty matrix  $A$ , and amount of spatial smoothness  $\beta$ ) were set and a patient with relatively high lesion load chosen as a reference to train the PV estimation algorithm. Specifically,  $A$  coefficients were  $a_1 = 11.25$ ,  $a_4 = 14.33$ ,  $a_5 = 0.47$ ,  $a_6 = 12.21$ ,  $a_7 = 1.33$ ,  $a_8 = 16.93$ , and  $\beta = 0.5$ . Patient mean intensities were set beforehand by histogram matching with the same reference patient used for hyperparameter setting [17].

## 2.4 CNNs

A novel MS segmentation method based on a cascade of two 3D patch-wise CNNs has recently been proposed [5]. The two networks have the same architecture and number of parameters, but don't share the same weights. Added to the above pre-processing steps, additional intensity normalization was performed, applying a histogram matching technique [17]. Afterwards, the first CNN receives as input patches of size 11x11x11 from different MRI modalities, centered around a voxel of interest. Only voxels with a FLAIR intensity over a threshold optimized in the validation phase are considered. Lesion candidates from the first CNN are

then given as input to the second one, which mainly has the task of reducing the false positives. In order to overcome the problem of data imbalance, before each CNN the negative class is undersampled, and the same number of positive and negative patches are obtained. Binary output masks are computed by linearly thresholding the probabilistic lesion masks given as output by the second network.

Table 1: Network architecture.  $c$  indicates the number of MRI modalities.

Layer	Type	Output size	Feature maps
0	Input	$c \times 11 \times 11 \times 11$	-
1	Convolutional	$32 \times 11 \times 11 \times 11$	32
2	Convolutional	$32 \times 11 \times 11 \times 11$	32
3	Max-pooling	$32 \times 5 \times 5 \times 5$	-
4	Convolutional	$64 \times 5 \times 5 \times 5$	64
5	Convolutional	$64 \times 5 \times 5 \times 5$	64
6	Max-pooling	$64 \times 2 \times 2 \times 2$	-
7	FC	256	256
8	Softmax	2	2

We have applied the same architecture [5] publicly available at [18] (see Table 1). Each convolutional layer is followed by a ReLU activation function and a batch normalization regularization. Dropout ( $p=0.5$ ) is applied before the first fully-connected layer. The networks were trained with the adaptive learning rate method (ADADELTA) [19], a batch size of 128, and early stopping as in the original paper. From the training dataset 7 cases were kept for validation, leaving 25 cases for training. With these the binarisation threshold was optimized considering equally the dice coefficient and the lesion false positive rate. In the original work [5] the CNNs were trained with 20 to 35 cases. Therefore, having a comparable number of patients for training, we hypothesize that this method should not perform worse in our study.

## 2.5 Combination of LeMan-PV with CNNs

It has been shown that for segmentation tasks, CNNs can benefit from prior probability maps fed in as an additional input channel [20, 21]. Moreover, combining different classifiers has also been a successful technique for improving the final results in supervised learning in several works [22, 23, 24]. Here, we propose a naive prototype combination (PV-CNNs) of both approaches described above. The concentration lesion maps generated by LeMan-PV are included as an additional input channel of the first CNN during training and testing. In this way, additional prior information on lesions was given to the network with the aim of improving the final segmentation.

### 3 Results

We compared the results of LeMan-PV, CNNs, and PV-CNNs strategies (see Figure 1). In line with three MS lesion segmentation challenges [2,3,4], we computed the following evaluation metrics: overlap Dice coefficient (Dice), lesion-wise false positive (LFPR) and lesion-wise true positive (LTPR) rates, voxel-wise true positives (TP), and volume difference (VD), according to [3,25]. Rather than a leave-one-out analysis [5,10], we present our results on a pure testing set of 73 patients cases acquired with different scanners. These two factors allow us to evaluate the generalization of the proposed methods in a setting close to the clinical scenario (shown in Table 3).

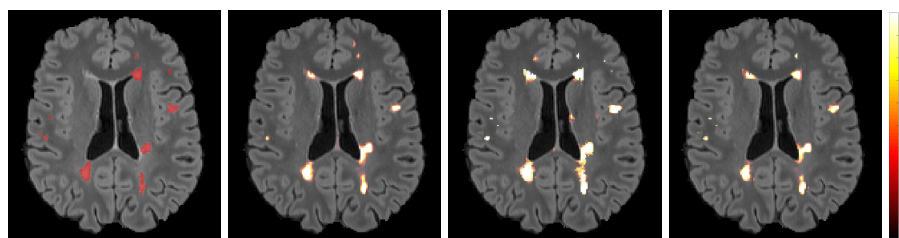


Fig. 1: Segmentation results (lesion probability (CNNs, PV-CNNs) and concentration (LeMan-PV)), from left to right: ground truth, LeMan-PV, CNNs, PV-CNNs. Reduction of LFPR is observed in PV-CNNs.

Quantitative evaluation at different lesion sizes (5, 10, 15 mm<sup>3</sup>) is given by ROC curves in Figure 2. Both LeManPV and PV-CNNs performed better at bigger minimum lesion size. However, CNNs did not show this behavior in our cohort, presenting similar ROC curves for all minimum lesion sizes.

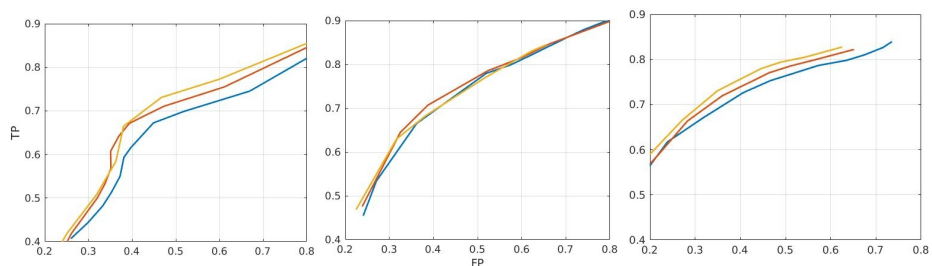


Fig. 2: ROC curves for different minimum lesion size: 5 (blue), 10 (orange), 15 (yellow) mm<sup>3</sup>. From left to right: LeMan-PV, CNNs, PV-CNNs.

As in the original studies [5,9], in what follows, a minimum lesion size of 5 mm<sup>3</sup> is considered. Median values for the whole test dataset are reported in

Table 2. LeManPV achieved the best Dice coefficient and volume difference. However, in terms of LFPR and LTPR, the CNNs performed better. The combination of the two methods outperformed them singularly in these lesion-wise metrics.

Table 2: Median values of the evaluation metrics for each method considered.

Method	Dice	LFPR	LTPR	TP	VD
LeMan-PV	<b>0.63</b>	0.37	0.57	<b>0.66</b>	<b>0.19</b>
CNNs	0.57	0.30	0.66	0.56	0.26
PV-CNNs	0.60	<b>0.26</b>	<b>0.69</b>	<b>0.66</b>	0.40

Segmentation results by Dice coefficient, TP, LFPR, and LTPR are given in the boxplots of Figure 3. Results are split in groups of patients according to their total lesion volume (TLV). In agreement with [16], we considered a low ( $TLV < 5\text{ ml}$ ), moderate ( $5\text{ ml} \leq TLV \leq 15\text{ ml}$ ) and high ( $TLV > 15\text{ ml}$ ) total lesion burden. Statistically significant differences between the methods are computed with Wilcoxon signed-rank test ( $p < 0.05$  uncorrected). Interestingly, PV-CNNs achieved the best Dice coefficient for low and medium TLV, but its performance drops for high lesion load. We hypothesize that the lower number of cases in this category (only 15 patients) downgrades the classification results for CNNs weakness to statistics. Overall, besides the presence of some outliers, LeMan-PV and CNNs showed a similar behavior at low and medium lesion loads. Regarding the TP, there are not significant differences between the three TLV. On the other hand, the LFPR decreases for all methods as the TLV increases. This represents an understandable behavior, as higher lesion load cases are expected to be better segmented. Curiously, and opposite to TP, the LTPR follows a similar trend. However, as stated above, the low number of patients at highest lesion load prevents us from drawing conclusions.

Volume differences are given (top row, for low and medium TLV patients only, bottom row: all dataset) by Bland-Altman plots (Figure 4). Slightly better results were obtained when combining both architectures, with a mean volume difference of  $-133.21\text{ ml}$ . However, a different behavior is shown when including the high TLV patients, with an increase of the mean volume difference to  $3250$ ,  $6410$  and  $7483\text{ ml}$  for LeMan-PV, CNNs and PV-CNNs respectively.

Finally, the effect of the scanner type is briefly investigated. Table 3 shows the mean Dice coefficient for the four different scanner types used to acquire the testing cases. For all segmentation methods, the highest Dice coefficient is achieved for the cases acquired with the TrioTim scanner. However, in this work the number of cases for each scanner is highly unbalanced. Therefore, further studies, with enlarged datasets, will be needed to quantify accuracy versus scanner type.

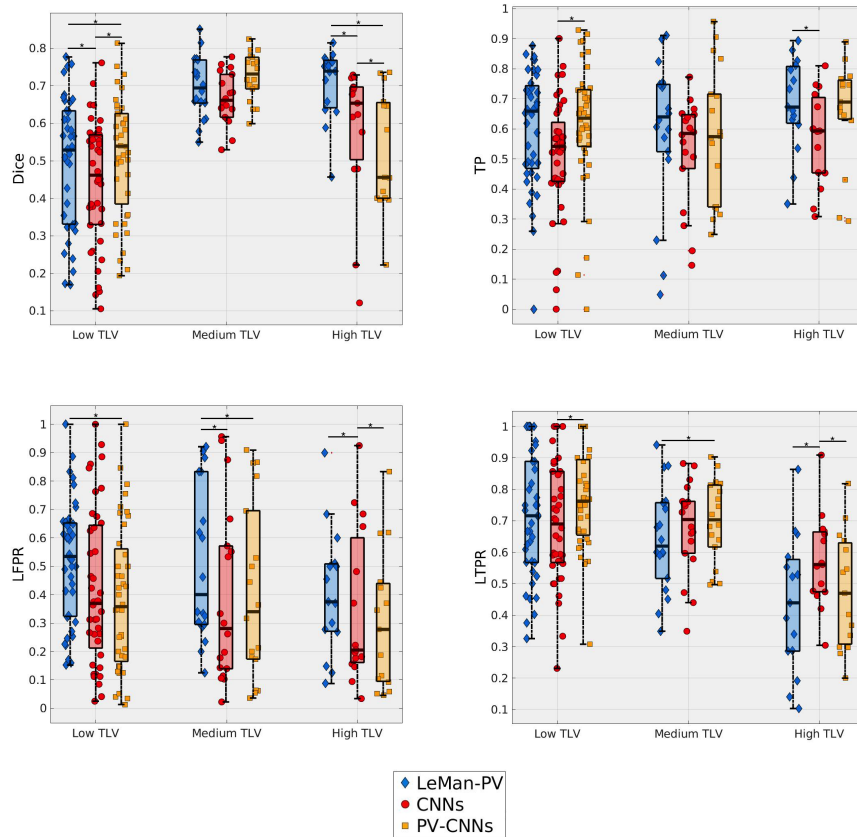


Fig. 3: Box plots of the Dice coefficient, TP, LFPR, and LTPR for the three methods considering the different TLV of the testing cases ( $p < 0.05$  is indicated by \*).

## 4 Conclusion

In this work, we presented the comparison of two of the most recent automated methods for WM lesions segmentation published in literature. In particular, we have tested a Bayesian partial volume estimation algorithm (LeMan-PV) [9] and a novel deep learning architecture based on a cascade of CNNs [5]. Both methods were tested on a pure test dataset composed of 73 cases, mainly belonging to early stage disease patients. The CNNs achieved the lowest LFPR of 30%. This confirms, as claimed in the original paper [5], that they are an effective method for reducing false positives. However, LeMan-PV showed the best segmentation results with the highest Dice coefficient (63%) and smallest volume difference (19%), indicating that PV might be still an asset for good delineation. Further analysis indicates a slight dependence of LeMan-PV performance on the

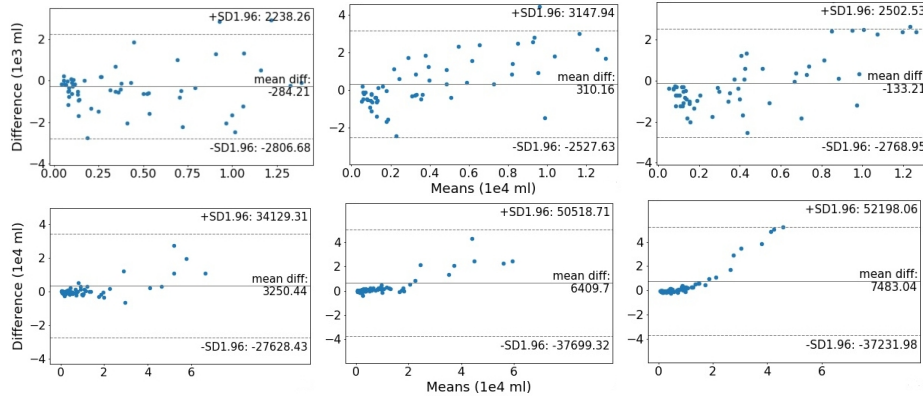


Fig. 4: Bland-Altman plots of low and medium TLV cases (top row) and of the whole dataset (bottom row) showing the volume differences of the three methods analyzed. From left to right: LeMan-PV, CNNs, PV-CNNs.

Table 3: Mean Dice coefficient of the testing cases for the different scanners they were acquired with.

Scanner	N. cases	Dice (range)		
		LeMan-PV	CNNs	PV-CNNs
Aera	5	0.59 (0.47-0.64)	0.47 (0.12-0.75)	0.52 (0.29-0.63)
TrioTim	6	0.63 (0.50-0.81)	0.61 (0.44-0.75)	0.65 (0.51-0.80)
Prisma.fit	11	0.54 (0.31-0.74)	0.48 (0.26-0.64)	0.53 (0.34-0.72)
Skyra	51	0.59 (0.16-0.84)	0.53 (0.11-0.78)	0.56 (0.19-0.80)

minimum lesion size considered, whereas the CNNs didn't show this behavior. Furthermore, a combination of the two methods (PV-CNNs) was implemented. Providing the CNNs with the probability maps of the LeMan-PV improved the LFPR (26%) and LTPR (69%) but did not perform well in terms of VD. Those results confirm that the hybrid of the two methods is also effective for WM lesion segmentation of early stages disease cases. However, further improvements are needed to increase the segmentation accuracy of low lesion burden cases, in which these automated methods achieved the worst performance (median Dice around 0.5). These cases are indeed of great importance for detecting MS lesions in the early stages of the disease. Future work will include experimenting with advanced combinations of these methods, training and testing on different datasets, and verifying if the results depend on the scanner used.

## Acknowledgements

The work is supported by the Centre d'Imagerie BioMédicale (CIBM) of the University of Lausanne (UNIL), the Swiss Federal Institute of Technology Lausanne



(EPFL), the University of Geneva (UniGe), the Centre Hospitalier Universitaire Vaudois (CHUV), the Hôpitaux Universitaires de Genève (HUG), and the Leenaards and Jeantet Foundations. This project is also supported by the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie project TRABIT (agreement No 765148). CG is supported by the Swiss National Science Foundation grant SNSF Professorship PP00P3-176984.

## References

1. Garcia-Lorenzo, Daniel, et al. "Review of automatic segmentation methods of multiple sclerosis white matter lesions on conventional magnetic resonance imaging." *Medical image analysis* 17.1 (2013).
2. Styner, Martin, et al. "3D segmentation in the clinic: A grand challenge II: MS lesion segmentation." *Midas Journal* 2008 (2008).
3. Carass, Aaron, et al. "Longitudinal multiple sclerosis lesion segmentation: Resource and challenge". *NeuroImage* 148 (2017).
4. Commowick, Olivier, Frdric Cervenansky, and Roxana Ameli. "MSSEG Challenge Proceedings: Multiple Sclerosis Lesions Segmentation Challenge Using a Data Management and Processing Infrastructure." *MICCAI* 2016.
5. Valverde, Sergi, et al. "Improving automated multiple sclerosis lesion segmentation with a cascaded 3D convolutional neural network approach." *NeuroImage* 155 (2017): 159-168.
6. Brosch, Tom, et al. "Deep 3D convolutional encoder networks with shortcuts for multiscale feature integration applied to multiple sclerosis lesion segmentation." *IEEE Trans.Med.Img.* 35.5 (2016).
7. Roy, Snehashis, et al. "Multiple Sclerosis Lesion Segmentation from Brain MRI via Fully Convolutional Neural Networks". *arXiv:1803.09172* (2018).
8. Grahl, S., et al. "Defining a minimal meaningful lesion size in multiple sclerosis." *Multiple Sclerosis Journal* 23 (2017).
9. Fartaria, Mário João, et al. "Partial volume-aware assessment of multiple sclerosis lesions." *NeuroImage: Clinical* 18 (2018): 245-253.
10. Fartaria, Mário João, et al. "Segmentation of Cortical and Subcortical Multiple Sclerosis Lesions Based on Constrained Partial Volume Modeling." *MICCAI* 2017.
11. <https://miac.swiss/en/>
12. Klein, Stefan, et al. "Elastix: a toolbox for intensity-based medical image registration.", *IEEE Trans.Med.Img.* 29.1 (2010).
13. Smith, Stephen M. "Fast robust automated brain extraction." *Human brain mapping* 17.3 (2002): 143-155.
14. Tustison, Nicholas J., et al. "N4ITK: improved N3 bias correction." *IEEE transactions on medical imaging* 29.6 (2010).
15. Kikinis R, Pieper SD, Vosburgh K (2014) *3D Slicer: a platform for subject-specific image analysis, visualization, and clinical support*. ISBN: 978-1-4614-7656-6.
16. Fartaria, Mário João, et al. "Automated detection of white matter and cortical lesions in early stages of multiple sclerosis." *J. Mag. Res. Imag.* 43.6 (2016).
17. Nyul, Lszl G., Jayaram K. Udupa, and Xuan Zhang. "New variants of a method of MRI scale standardization." *IEEE Trans.Med.Img.* 19.2 (2000).
18. <https://github.com/sergivalverde/cnn-ms-lesion-segmentation>

19. Zeiler, Matthew D. "ADADELTA: an adaptive learning rate method". arXiv:1212.5701 (2012).
20. Luo, Kunming, et al. "A CNN-based segmentation model for segmenting foreground by a probability map." Intelligent Signal Processing and Communication Systems (ISPACS), IEEE ISBI 2017.
21. Zotti, Clement, et al. "GridNet with automatic shape prior registration for automatic MRI cardiac segmentation". arXiv:1705.08943 (2017).
22. Kotsiantis, Sotiris B., Ioannis D. Zaharakis, and Panayiotis E. Pintelas. "Machine learning: a review of classification and combining techniques." Artificial Intelligence Review 26.3 (2006).
23. Fartaria, Mário João, et al. "An ensemble of 3D convolutional neural networks for central vein detection in white matter lesions." MIDL 2018 Abstract Submission.
24. Kamnitsas, Konstantinos, et al. "Ensembles of Multiple Models and Architectures for Robust Brain Tumour Segmentation." arXiv:1711.01468 (2017).
25. E. Geremia, O. Clatz, B.H. Menze, E. Konukoglu, A. Criminisi, N. Ayache, Spatial decision forests for MS lesion segmentation in multi-channel magnetic resonance images, 57, NeuroImage.