

A parameter-efficient deep learning approach to predict conversion from mild cognitive impairment to Alzheimer's disease within three years

Authors: Simeon Spasov^{a*}, Luca Passamonti^b, Andrea Duggento^c, Pietro Liò^{a#}, and Nicola Toschi^{c,d#}

a

University of Cambridge, Cambridge, Department of Computer Science and Technology, William Gates Building, 15 J J Thomson Ave, Cambridge, CB3 0FD, UK (email: ses88@cam.ac.uk).

b

Department of Clinical Neurosciences, University of Cambridge, Herchel Smith Building, Forvie Site, Robinson Way, Cambridge Biomedical Campus, Cambridge CB2 0SZ Cambridge (e-mail: lp337@medschl.cam.ac.uk).

c

Department of Biomedicine and Prevention, University of Rome "Tor Vergata", Via Cracovia, 00133 Roma RM, Italy (e-mail: toschi@med.uniroma2.it).

d

A.A. Martinos Center for Biomedical Imaging, Massachusetts General Hospital and Harvard Medical School, Boston (USA) (e-mail: toschi@med.uniroma2.it).

* corresponding author

these authors contributed equally to this publication

key words: deep learning, neural networks, classification, Mild Cognitive Impairment, Alzheimer's disease, Magnetic resonance imaging, ADNI, Early diagnosis

Abstract

Some forms of mild cognitive impairment (MCI) can be the clinical precursor of severe dementia like Alzheimer's disease (AD), while other types of MCI tend to remain stable over-time and do not progress to AD pathology. To choose an effective and personalized treatment for AD, we need to identify which MCI patients are at risk of developing AD and which are not.

Here, we present a novel deep learning architecture, based on dual learning and an ad hoc layer for 3D separable convolutions, which aims at identifying those people with MCI who have a high likelihood of developing AD.

Our deep learning procedures combine structural magnetic resonance imaging (MRI), demographic, neuropsychological, and APOe4 genotyping data as input measures. The most novel characteristics of our machine learning model compared to previous ones are as follows: 1) multi-tasking, in the sense that our deep learning model jointly learns to simultaneously predict both MCI to AD conversion, and AD vs healthy classification which facilitates the relevant feature extraction for prognostication; 2) the neural network classifier employs relatively few parameters compared to other deep learning architectures (we use ~500,000 network parameters, orders of magnitude lower than other network designs) without compromising network complexity and hence significantly limits data-overfitting; 3) both structural MRI images and warp field characteristics, which quantify the amount of volumetric change compared to the common template, were used as separate input streams to extract as much information as possible from the MRI data. All the analyses were performed on a subset of the Alzheimer's Disease Neuroimaging Initiative (ADNI) database, for a total of n=785 participants (192 AD, 409 MCI, and 184 healthy controls (HC)).

We found that the most predictive combination of inputs included the structural MRI images and the demographic, neuropsychological, and APOe4 data, while the warp field metric added little predictive value. We achieved an area under the ROC curve (AUC) of 0.92 with a 10-fold cross-validated accuracy of 86%, a sensitivity of 87.5% and specificity of 85% in classifying MCI patients who developed AD in three years' time from those individuals showing stable MCI over the same time-period. To the best of our knowledge, this is the highest performance reported on a test set achieved in the literature using similar data. The same network provided an AUC of 1 and 100% accuracy, sensitivity and specificity when classifying NC from AD. We also demonstrated that our classification framework was robust to different co-registration templates and possibly irrelevant features / image sections.

Our approach is flexible and can in principle integrate other imaging modalities, such as PET, and a more diverse group of clinical data.

The convolutional framework is potentially applicable to any 3D image dataset and gives the flexibility to design a computer-aided diagnosis system targeting the prediction of any medical condition utilizing multi-modal imaging and tabular clinical data.

Introduction

More than 30 million people have a clinical diagnosis of Alzheimer's disease (AD) worldwide and this number is expected to triple by 2050 (Barnes and Yaffe, 2011), due to increased life expectancy and improvements in care (Ferri et al., 2005). AD is a form of dementia characterized by extracellular β -amyloid peptide plaque deposits and abnormal tau accumulation and phosphorylation which ultimately lead to neuronal and synaptic loss (Murphy et al. 2010). AD-related neurodegeneration follows specific patterns which arise from subcortical areas and spread to the cortical mantle (Braak and Braak et al. 1996). The classic clinical hallmark of the most common form of AD (i.e., the amnesic type) is represented by impairments in episodic memory, followed by visuo-spatial and orientation problems, and ultimately frank dementia.

Mild cognitive impairment (MCI) is a wide and heterogeneous spectrum of disorders which causes relatively less acute and noticeable memory deficit than AD. However, around 10%-15% of MCI patients convert to AD per year (Braak and Braak, 1995; Mitchell and Shiri-Feshki, 2008) within less than 5 years, although the conversion rate decreases later on. As the majority of MCI to AD conversions happen within 5 years, it is crucial to early identify the MCI subjects at risk of developing AD as soon as possible. The MCI patients who do not convert to AD tend to either remain stable, develop other forms of dementia, or even revert to a healthy state, which suggest that MCI is a conundrum of disorders which are likely to be associated with the several etio-pathogenetic mechanisms. AD-related neuropathological markers have been also observed several years before the clinical manifestation of memory symptoms (Braak and Braak, 1996; Delacourte et al., 1999; Morris et al., 1996; Serrano-Pozo et al., 2011; Mosconi et al., 2007), which suggests that AD development could be predicted before clinical onset via in vivo biomarkers (e.g. PET and MR imaging as well as blood or cerebrospinal fluid (CSF) biomarkers) (Markesbery, 2010; Baldacci et al., 2018; Hampel et al. 2018; Teipel et al., 2018). In this context, MRI imaging has garnered interest in AD diagnosis, and perhaps more importantly in prognosticating the MCI to AD conversion. Relative to CSF and PET biomarkers, MRI measures have the advantage of not using ionising radiation, of being non-invasive, less expensive and more widely available in less specialized medical environments. MRI markers also enable the possibility to gather multimodal information (e.g. structural and functional) within the same session.

In this context, there has been a growing interest in developing MRI-based computational tools to discriminate AD patients from healthy individuals, and most importantly in distinguishing between stable MCI individuals and MCI patients who go on to develop AD. To this end, different clinical data and imaging modalities have been employed with variable rates of success, including PET studies (Choi et al. 2018; Mosconi et al. 2004, Mosconi et al. 2007, Shaffer et al. 2013, Young et al. 2013), MRI studies (Filipovych et al. 2011; Moradi et al. 2015; Mosconi et al. 2007; Tong et al. 2017, Young et al. 2013), cognitive testing studies (Casanova et al. 2011; Moradi et al. 2015), and CSF biomarker studies (Davatzikos et al. 2011; Hansson et al. 2006; Riemenschneider et al. 2002; Sonnen et al. 2010). Most of the above-mentioned studies employ a classification pipeline, which relies on two independent steps. First, a dimensionality reduction method, such as ICA (Shaffer et al. 2013), L1 regularization (Moradi et al. 2015; Tong et al. 2017) or morphometry (Davatzikos et al. 2011; Fan et al. 2007), is used to reduce the raw images or volumes to a relatively small number of (possibly) highly descriptive factors. Then, these factors are fed into a multivariate pattern classification algorithm. Notably, the dimensionality reduction and classification algorithms are two separate mathematical models which involve different assumptions, hence possibly resulting in loss of relevant information in the classification process (Nguyen and Torre, 2010). Examples of studies that have used this classification pipeline to predict MCI to AD conversion using structural MRI and cognitive measures at baseline are described in Moradi et al. 2015 and Tong et al. 2017. The authors firstly perform feature selection to extract informative MRI voxels via regularized logistic regression, and subsequently use the extracted voxels, as well as the cognitive measures, to produce support vector machine (SVM)-based predictions, achieving an area under the ROC curve (AUC) between 0.9 and 0.92. In the case of Hojjati et al., 2017, who use baseline resting state fMRI data and achieve an AUC of 0.95, features are engineered by constructing a brain connectivity matrix which is treated as a graph, and the extracted graph measures are

inputted into a SVM. Also, the most frequently used classifiers, such as SVM (Moradi et al., 2015; Hojjati et al., 2017, Tong et al., 2017) and Gaussian Processes (Young et al., 2013), require the use of kernels, or data transformations, chosen from a limited user-specified set, which map the data to a new space in the hope that it will be more easily separable. However, constructing or choosing an application-specific kernel to act as a reasonable similarity measure for the task at hand is not always possible.

The use of two disjoint pipelines and the need to construct ad-hoc kernels can be surmounted by the use of a class of algorithms known as deep learning, which afford much greater representational flexibility than kernel-based methods and also automatically learn data transformations which maximize any given performance metric. Such methods have been applied to AD vs. healthy subject discrimination (Hosseini-Asl et al., 2016; Liu et al., 2015; Liu et al., 2018; Payan and Montana, 2015) and progressive MCI (pMCI) vs stable MCI (sMCI) classification (Choi et al., 2018; Lu et al., 2018 (1); Lu et al., 2018 (2)) As an example, Choi et al., 2018 and Lu et al., 2018 (1) use deep learning to achieve one of the highest pMCI/sMCI classification performances to-date despite the use of a single (albeit very informative) imaging modality based on ionizing radiation (PET) (~84% - 82% conversion rate accuracies for these studies respectively). A comparison between recent studies and methods is provided in Table 3). As is well known however, the superior representational capacity of deep learning methods relies on a high number of neural network parameters. Frequently, this gives rise to overfitting, i.e. a satisfactory training performance which however does not generalize well to unseen samples during testing or when applying the model. Although it has been demonstrated that deep learning approaches can yield impressive performance, the data-scarce nature of medical datasets is not commonly sufficient to build a useful network architecture.

The aim of this paper is therefore to develop and employ a parameter-efficient neural network architecture, based on more recent convolutional neural network layers, namely 3D separable and grouped convolutions (which were developed specifically for computer vision tasks). Additionally, we use a combination of input streams (including structural MRI as well as clinical variables comprising demographic, neuropsychological, and APOe4 genotyping data) for the joint multi-task classification of pMCI/sMCI and AD/HC, i.e. joint or dual-learning. These newer network designs have been shown to yield superior performance on other types of visual discrimination problems (Chollet et al., 2016) while maintaining the overall network parameter count low, hence efficiently battling the overfitting problem. Additionally, developed a novel feature extractor sub-network and, in order to employ these methods efficiently, we combined the Tensorflow (Abadi et al., 2016) and Keras (Chollet et al., 2015) libraries with our own 3D implementation of 3D separable convolutions which is available freely upon request.

Methods:

1. Participants and data

All data was obtained from the Alzheimer’s Disease Neuroimaging Initiative (ADNI) and comprised 435 men and 350 women aged between 55 and 91 years. The majority of subjects identified as white (>94%) and non-Hispanic (99.98%). All data we used is summarized in Table 1. Differences in median age across groups were tested using Friedman’s ANOVA and group x gender interactions were tested using Fisher’s exact test. None of these interactions resulted statistically significant ($p > 0.05$). For all participants, we employed the Magnetization Prepared Rapid Gradient-Echo (MPRAGE) T1-weighted image (structural MRI) as well as the following meta-data: demographic data (age, gender, ethnic and racial categories, education), neuropsychological cognitive assessment tests like dementia rating scales (CDRSB), the Alzheimer’s disease assessment scale (ADAS11, ADAS13), episodic memory evaluations in the Rey Auditory Verbal Learning Test (RAVLT), and APOe4 genotyping.

	No. of subjects	Age (years)	Male/Female	years in education	APOe4 expression			CDRSB	ADAS11	ADAS13	RAVLT			
					0	1	2				immediate	learning	forgetting	% forget
AD	192	75.6±7	103/81	15±2.9	57	86	41	4.4±1.6	18.8±6	29±7.3	23±7	1.7±1.8	4.4±1.9	89.4±21.2
HC	184	74.6±6	92/100	16.3±2.7	144	43	5	0.2±0.9	6±3.8	9.3±5.7	44±10.5	6±2.4	3.7±2.7	33.1±27.7
pMCI	181	73.7±7	108/73	15.9±2.8	61	90	30	2±1	13.5±4.2	21.9±5.5	27.2±6.5	2.9±2.2	4.9±2.1	78.3±27
sMCI	228	72.2±7	132/96	16±2.8	145	67	16	1.2±0.6	8.4±3.3	13.5±5.3	38.5±10	4.75±2.5	4.35±2.6	50±30

Table 1. The table summarizes the demographic, neuropsychological, cognitive assessment and APOe4 genotyping data used in this study to classify between progressive and stable MCI, and healthy and AD subjects. The data is presented in a mean±std format. The abbreviations used are APOe4 - Apolipoprotein E; CDRSB – Clinical Dementia Rating Sum of Boxes; ADAS – Alzheimer’s Disease Assessment Scale; RAVLT – Ray Auditory Verbal Learning Test.

2. Data Preprocessing

Prior to classification, all T1 weighted (T1w) images were registered to a common space. Two different T1 templates were used in order to assess the robustness of our classification methodology to structural misalignment. First, we built a custom T1 template specific to this study. To this end, we employed all T1w images, which (after N4 bias field correction) were nonlinearly co-registered to each other and averaged iteratively (i.e. the group average was recreated at the end of each iteration). The procedure was based on symmetrical diffeomorphic mapping and employed five total iterations. The second template was the Montreal Neurological T1 Template (MNI152_T1_1mm). All single-subject T1w images were nonlinearly registered to both templates.

All template creation and registration procedures were performed using the ANTs package (Avants et al., 2010, Avants et al., 2011). In detail, the high-dimensional non-linear transformation (symmetric diffeomorphic normalization transformation) model was initialized through a generic linear transformation which consisted of center of mass alignment, rigid, similarity and fully affine transformations followed by (metric: neighbourhood cross correlation, sampling: regular, gradient step size: 0.12, four multi-resolution levels, smoothing sigmas: 3, 2, 1, 0 voxels in the reference image space, shrink factors: 6, 4, 2, 1 voxels. We also used histogram matching of images before registration and data winsorisation with quantiles: 0.001, 0.999. The convergence criterion was set to be as follows: slope of the normalized energy profile over the last 10 iterations < 10-8). Co-registration of all scans required approximately 19200 hour of CPU time on a high-performance parallel computing cluster. After co-registration, all images were masked to include only brain tissue using brainmasks generated in template space using BET, part of FSL (Jenkinson et al., 2012).

After co-registration to both templates we also extracted the local Jacobian Determinant (JD) of the nonlinear part of the deformational field taking each image into template space. The JD maps were used to complement the co-registered MRI scans as an additional input stream in our model (see below). All images were masked. Additionally, in order to evaluate how much a priori knowledge about AD brain pathophysiology could improve our classification and also how much irrelevant features hamper classification performance, we defined a set of regions of interest (ROI) masks which included only brain areas known to be heavily involved in AD-related atrophy, namely parietal, temporal and frontal lobes (i.e. we performed an inclusion test). This was based on the [Hammers et al. 2003](#) atlas[©] Copyright Imperial College of Science, Technology and Medicine 2007 (www.brain-development.org).

Numerical normalization for the co-registered MRI images was performed per sample, i.e. each 3D volume was standardized to 0 mean and unit standard deviation. The reasoning behind this is that brain atrophy could be recognized as an in-sample shift in intensity for a certain area compared to other regions. The normalization applied to the clinical features, i.e. the demographic, neuropsychological, and APOe4 genotyping data, also follows the same feature scaling procedure, where the values of each separate clinical factor are normalized between [0, 1]. On the other hand, the extracted JD images were feature-scaled to have voxel values in the [0;1] range via subtracting the smallest value in the entire JD image set, and dividing by the difference between the largest and smallest values (also in the entire JD image set). This retains class-wise differences in volumetric changes created when co-registering an image to a template while rescaling the data to a global maximum and minimum.

3. Deep Learning Architecture

3.1. Architecture Overview

The network architecture is summarized in fig. 1 (b). In this paper, we developed a feature extractor sub-network (referred to as the *multi-modal feature extractor* in fig. 1 b), inspired by the parameter-efficient separable and grouped convolutional layers presented in AlexNet ([Krizhevsky et al., 2012](#)) and Xception ([Chollet, 2017](#), [Velickovic et al., 2016](#)). In detail, the layers of the feature extractor are shared between two tasks - MCI-to-AD conversion prediction and AD/HC classification, as we assume both problems share common underlying factors, i.e. the MCI subjects lie on the HC-AD continuum. This means similar data transformations are likely to be useful for prediction of different problems, and additionally this procedure increases the number of samples the extractor network is trained on hence reducing overfitting. Balancing between the tasks can be seen as imposing soft constraints on the network parameters, and if some of the factors that explain the variations in our data are shared between the two discrimination problems, overfitting is reduced further. The inputs for each task are its respective image modalities and clinical features. The feature extractor sub-network extracts 4-dimensional vectors for each of the two classification problems. These resulting latent representations are then processed by two separate fully connected layers (see fig. 1 a) with sigmoid activations and a binary cross-entropy loss applied at the output of each. The outputs of the fully connected layers are in the 0 to 1 range. The closer the activation is to 1, the more confident the model is that the input pattern corresponds to a diseased individual (i.e. AD or pMCI, depending on the classification task), and vice versa.

3.2. Mathematical formulation of Model

We will denote the input data and labels as pairs $(\mathbf{X}, \mathbf{Y}) = \{(\mathbf{x}^A_1, y^A_1), \dots, (\mathbf{x}^A_N, y^A_N), \dots, (\mathbf{x}^M_1, y^M_1), \dots, (\mathbf{x}^M_{N'}, y^M_{N'})\}$, where \mathbf{x}^A_i is the i -th observation from the Alzheimer's and healthy subset, and \mathbf{x}^M_j is the j -th observation from the pMCI vs sMCI subset. Both classification problems have corresponding class labels y^A_i and $y^M_j \in \{0, 1\}$. We refer to the empirical distributions over the AD/HC and MCI subsets as $\tilde{p}_A(\mathbf{x}, y)$ and $\tilde{p}_M(\mathbf{x}, y)$ respectively. The model log likelihoods (i.e. the conditional probabilities of the target variables, y , given the input data \mathbf{x} which we model with the neural network) for the two classification problems are given by:

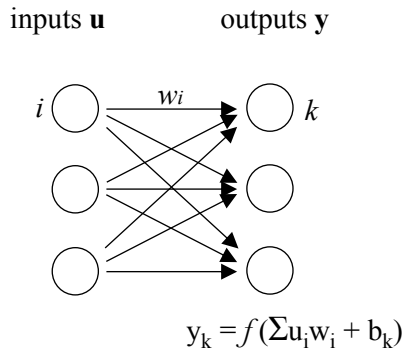


Fig. 1 (a) The figure depicts the operation of a dense or fully connected layer. The outputs y_k are formed as a non-linear transformation of the input vector \mathbf{u} . The non-linear activation works on a weighted sum of the inputs, $\sum u_i w_i$, and a bias term b_k . These layers are employed to process the clinical inputs in the Multi-modal feature extractor and to produce the output labels of our model.

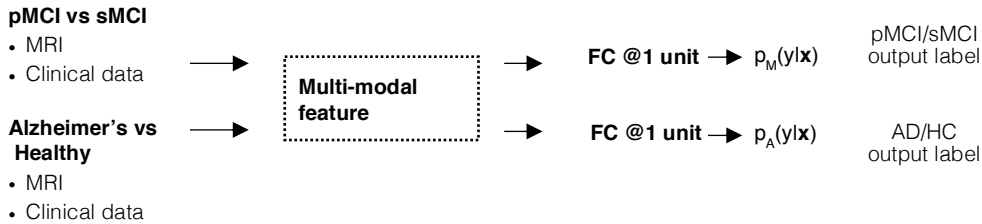


Fig. 1 (b) . The figure presents an overview of our multi-tasking neural network methodology. We have designed a sub-network (the multi-modal feature extractor) to extract 4-d feature representations from the inputs of both tasks/datasets. This sub-network (with the same θ network parameters) is applied on the data from both the pMCI/sMCI and AD vs healthy discrimination problems, as we assume the underlying factors of the conditions are similar, hence similar data transformations are likely to be useful. We then employ two fully connected layers, parametrized by ϕ and ψ , with sigmoid outputs. The sigmoid outputs approximate the conditional distribution of the labels for the two problems given the inputs ($p_A(y|\mathbf{x})$ for the AD vs healthy task and $p_M(y|\mathbf{x})$ for the pMCI vs sMCI task). We learn the network parameters such that our model outputs correspond to the true labels in the dataset by minimizing the binary cross-entropy between the observed and estimated targets. The multi-modal feature extractor is represented by a dashed-line rectangle in fig. 1 (b) and fig. 3.

$$\log p_A(y^A_i | \mathbf{x}^A_i; \theta, \phi) = f_A(y^A; \mathbf{x}^A, \theta, \phi) = -L_A \quad \log p_M(y^M_j | \mathbf{x}^M_j; \theta, \psi) = f_M(y^M; \mathbf{x}^M, \theta, \psi) = -U_M \quad (1)$$

The likelihood functions f_A and f_M are modelled as Bernoulli distributions, parametrized by neural network-based transformations of the input data as described in fig. 1 (b). The goal is to learn the network parameters such that we can approximate the *true* conditional probabilities of the labels given the inputs via the likelihood functions given by eq. 1. We use θ to denote the parameters in the feature extractor sub-network, and ϕ and ψ to denote the weights in the final fully connected layers that output the class probabilities for the Alzheimer's vs healthy and pMCI vs sMCI tasks respectively. Learning the network parameters can be represented as:

$$\operatorname{argmin} (\theta, \phi, \psi) E_{\mathbf{x}, y \sim \tilde{p}_M(\mathbf{x}, y)} [U_M] + \alpha E_{\mathbf{x}, y \sim \tilde{p}_A(\mathbf{x}, y)} [L_A] \quad (2)$$

As U_M and L_A represent negative log-likelihoods, the objective function given in eq. (2) can be viewed as minimizing the weighted sum between two binary cross-entropy terms between the observed and estimated (by our network) class probabilities. Intuitively, learning the network parameters is maximizing the probability of observing the labels in both datasets under the model, given the input cognitive, genetic and MRI biomarkers. We also introduced the α hyperparameter to control the trade-off between the two tasks during learning, and use $\alpha = 0.25$ in all experiments. Although the choice of α comes from our subjective view of the importance of one task over another and appears to lack rationale, we found that the AD/HC problem is much easier and the model quickly achieves high validation accuracy (see table 3) when $\alpha = 0.25$.

3.3. 3D Convolutions

Convolutional layers in our study work by convolving an input tensor, \mathbf{x} , with a kernel of weights \mathbf{W} , then

adding a bias term b , and finally passing the result through a non-linearity. To extract a rich set of representations we repeat this process with K different kernels (also known as channels or filters) convolving the same tensor \mathbf{x} , each resulting in a new *feature map* \mathbf{h}_k . Hence, we can write:

$$\mathbf{h}_k = f(\mathbf{W}_k * \mathbf{x} + b_k) \quad (3)$$

The feature map subscript is $k = [1, \dots, K]$. The function f can be selected from a range of differentiable non-linear transformations, such as the sigmoid $f(u) = (1 + \exp(-u))^{-1}$ and the exponential linear unit, or ELU, (Clevert et al. 2015): $f(u) = u$ if $u \geq 0$ and $f(u) = \exp(u) - 1$ if $u < 0$. We rely on the ELU transformation in our hidden layer activations and a sigmoid output for label predictions. The set of K feature maps extracted from the input \mathbf{x} defines a single layer $\ell = [1, \dots, L]$ in our convolutional neural network. Thus, the k^{th} feature map at layer ℓ is denoted as \mathbf{h}_k^ℓ . To construct a hierarchy of features we can use the outputs of layer $\ell-1$ as inputs to layer ℓ :

$$\mathbf{h}_k^\ell = f(\mathbf{W}_k^\ell * \mathbf{h}^{\ell-1} + b_k^\ell) \quad (4)$$

where \mathbf{h}^0 is \mathbf{x} . Note that in eq. (2), $\mathbf{h}^{\ell-1} = [\mathbf{h}_0^{\ell-1}, \dots, \mathbf{h}_K^{\ell-1}]$ is a 4-D tensor - a collection of the K 3D feature maps extracted at layer $\ell-1$. Consequently, \mathbf{W}_k^ℓ is also a 4-D tensor kernel of size $N^1 \times N^2 \times N^3 \times K$. This filter is multiplied element-wise during convolution with a $N^1 \times N^2 \times N^3$ patch in each of the K feature maps and then the result is summed to produce a single scalar element after adding a bias term and passing through a non-linear function. The convolutional procedure can be seen as sliding this kernel with strides in all three dimensions to produce \mathbf{h}_k^ℓ . It is important to note that the number of parameters needed to extract K^ℓ feature maps in layer ℓ from the $K^{\ell-1}$ feature maps in layer $\ell-1$ is given by:

$$(N^1 * N^2 * N^3 * K^{\ell-1} + 1) * K^\ell \quad (5)$$

where $N^1 \times N^2 \times N^3$ is the filter size used (see section 3.8 for actual values used in this paper).

3.4. Fully connected (Dense) Layers

Fully connected layers are designed to work on vectorized inputs u . The operation of the dense layer is depicted on fig. 1 (a) Each input u_i has an associated weight w_i . In order to produce an output y_k , we form the weighted sum of all inputs $\sum u_i w_i$, then add a bias term b_k , and pass the result through a differentiable non-linear function like the sigmoid or the exponential linear unit. We can repeat this procedure K times with different weight parameters to produce an output vector y , which can be used as an input to another fully connected layer. In our work we employ these dense connections to process the tabular clinical features and to produce the final output predictions (or probability scores) of our model.

3.5. Batch normalization, dropout, L2 regularization

Several standard strategies are used in our network to battle overfitting. The first one is batch normalization (Ioffe and Szegedy 2015) which normalizes a layer's outputs by subtracting their mean and dividing by the standard deviation. This whitening procedure enforces a fixed distribution of activations which has been shown to stabilize and facilitate the process of training. The technique accelerates the rate of training of deep neural nets and can act as a regularize. A second strategy is dropout by (Srivastava et al. 2014) which works by randomly dropping units and their connections during training. An intuitive explanation of its efficacy is that each unit must learn to extract useful features on its own with different sets of randomly chosen inputs. As a result, each hidden unit is more robust to random fluctuations and learns a generally useful transformation. Finally, L2 regularization penalizes weights of high absolute value, hence directly limiting the variety of functions our model can represent, i.e. its capacity.

3.6. Separable Convolutions

The separable convolutions we employ are similar to standard convolutional layers but reformulate the procedure in two steps by performing *depthwise* and then *pointwise* operations. Firstly, each input channel is spatially convolved separately, then the resulting outputs are mixed via *pointwise* convolutions with a kernel size of $1 \times 1 \times 1$. The depthwise procedure simply reformulates the convolutional operation from eq. (2) to:

$$\mathbf{h}_k^\ell = f(\mathbf{W}_k^\ell * \mathbf{h}_k^{\ell-1} + b_k^\ell) \quad (6)$$

Conceptually, the difference is that the feature map at layer ℓ , \mathbf{h}_k^ℓ , only depends on a single feature map from the previous layer $\ell-1$ for the depthwise procedure. On the other hand, standard convolutions take as an input all $K^{\ell-1}$ feature maps to produce a single output. Since in the depthwise case \mathbf{h}_k^ℓ is produced from a single map from the previous layer, the parameter count in \mathbf{W}_k^ℓ is reduced. A depthwise convolution reduces the number of parameters employed to $(N^1 * N^2 * N^3 + 1) * K^\ell$, which is $K^{\ell-1}$ times more parameter-efficient as seen when compared to eq. (5). The pointwise operation mixes all the channels together and requires $K^\ell * K^{\ell-1}$ parameters. Hence, the overall number of weights in separable convolutions is given by:

$$(N^1 * N^2 * N^3 + 1) * K^\ell + K^\ell * K^{\ell-1} \quad (7)$$

Considering the kernel sizes and number of filters in our network architecture, substituting a single convolutional layer results in ~ 20 times less parameters for the separable module. In order to achieve the above operations, we implemented an ad-hoc 3D separable convolution module as a custom Keras layer based on a TensorFlow backend.

3.7. Grouped Convolutions

The grouped layer can be viewed as a compromise between standard convolutions and the separable case. This procedure splits the previous layer's feature maps in two groups (G1 and G2) and treats them as separate

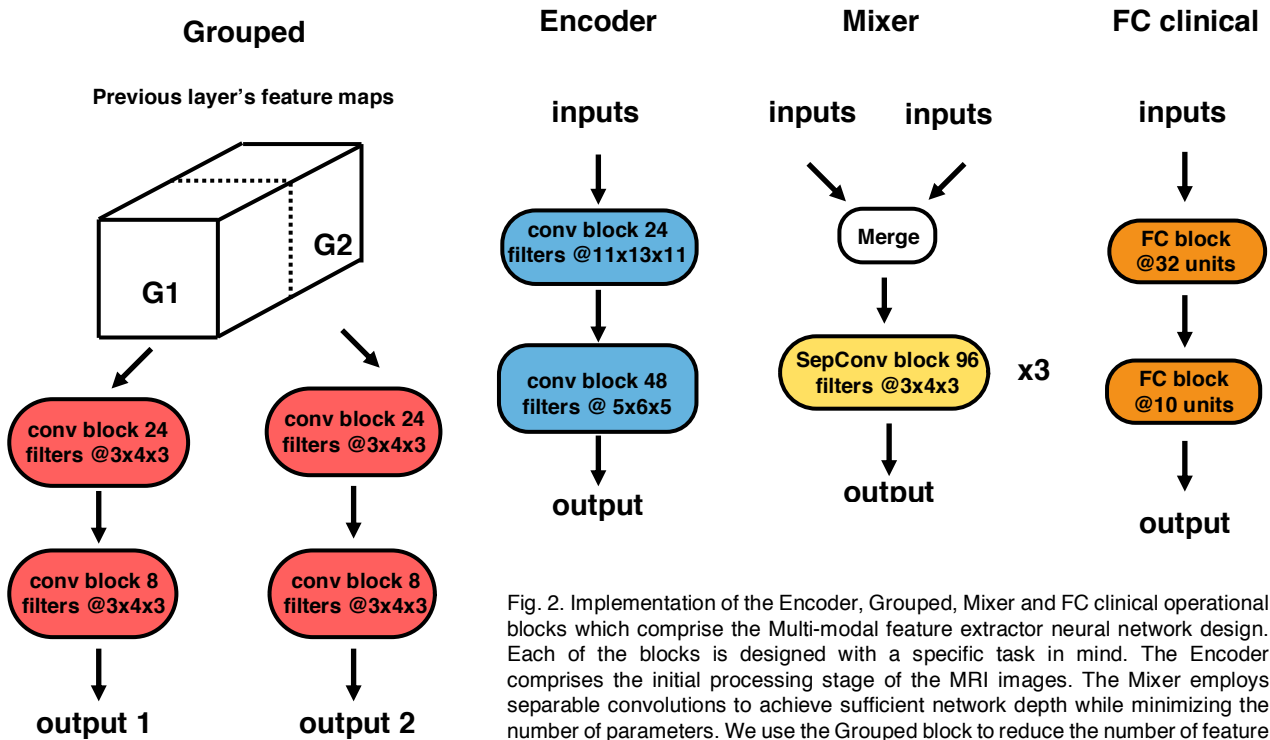


Fig. 2. Implementation of the Encoder, Grouped, Mixer and FC clinical operational blocks which comprise the Multi-modal feature extractor neural network design. Each of the blocks is designed with a specific task in mind. The Encoder comprises the initial processing stage of the MRI images. The Mixer employs separable convolutions to achieve sufficient network depth while minimizing the number of parameters. We use the Grouped block to reduce the number of feature maps in the network layers in an efficient manner. Finally, the tabular clinical data is processed via standard dense connections in FC Clinical.

when applying further transformations (fig. 2). As a result, only half of the channels are used to produce a single output feature map. The grouped layer requires twice less parameters than the standard convolutional approach, assuming the same overall number of output feature maps is generated.

3.8. Multi-modal feature extractor

Since several different sequences of layers are frequently reused, they are combined in operational blocks. Each block follows a similar pattern. For instance, convolutional blocks, used to process the 3D MRI tensors, comprise a convolutional kernel with linear activations, batch normalization and an exponential linear unit (ELU) transformation with dropout. In order to reduce the resulting spatial dimensions, max pooling is used, where only the highest value in an image patch is retained, with a window of 3 pixels and a stride of 2. Each operation is applied on the outputs of the previous one. On the other hand, the clinical features undergo a series of transformations by dense blocks. Since these blocks act on vectorized inputs, a linear dense layer is employed instead but the same regularization precautions and activations as above are applied. In a similar fashion we constructed 3D depthwise separable convolutional blocks for efficient use of model parameters. These operational blocks are then combined to form higher level blocks that implement specific tasks. Fig. 2 depicts how the convolutional, dense and depthwise separable blocks are used for this purpose. Fig. 3 shows how these higher-level blocks can be combined to extract low dimensional features from our input data, such that we can then model the conditional probabilities of the target labels for the AD/HC and pMCI/sMCI classification problems.

For example, the Encoder block is used as a first embedding stage for both the structural MRI and Jacobian Determinant inputs. Then the Mixer takes these embeddings and applies a series of three consecutive separable convolutions in order to achieve efficient network depth which has been empirically shown to yield superior performance (Chollet 2017). Since the output from the Mixer block has a relatively high number of feature maps, we use a Grouped block instead of standard convolutions as a final stage of MRI data embedding. The clinical features are compressed in the FC clinical block and concatenated with the Grouped block's outputs and finally transformed to a 4-dimensional vector. We refer to the neural network which compresses all the input biomarkers to 4-D vectors as the Multi-modal feature extractor.

The Encoder comprises a series of two convolutional blocks with 24 and 48 kernels of sizes (11 x 13 x

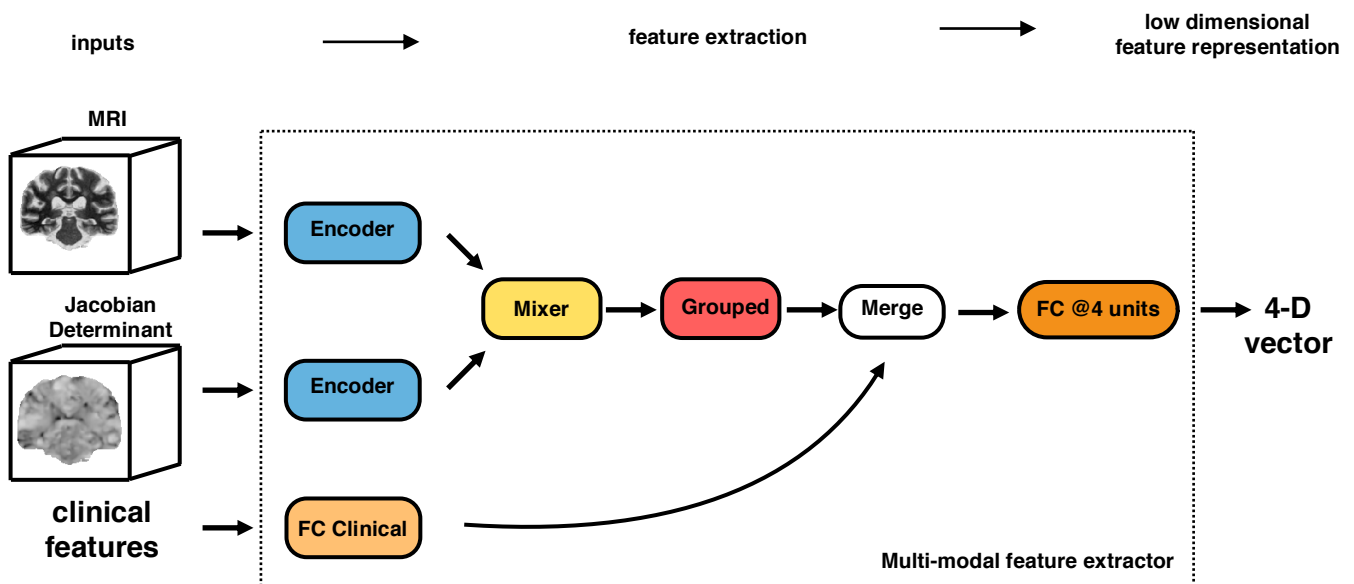


Fig. 3. The figure presents the architecture of the Multi-modal feature extractor sub-network. It is designed to take multiple 3D in order to combine the co-registered structural MRI and the Jacobian Determinant volumes. The Encoder, Mixer and Grouped blocks are designed such that the high-dimensional image inputs are efficiently embedded in a lower-dimensional vector which summarizes all useful visual information. This vector is then merged with the compressed clinical data (via FC Clinical) to output a final 4-D embedding of all input biomarkers. The operational blocks are color-coded for the ease of the reader both in fig. 2 and fig. 3. The resulting 4-D representations are passed through a dense connection with a sigmoid output (see fig. 1 a) to model the target variables for each of the two classification problems.

11) and (5 x 6 x 5) and strides of 4 and 1 in all dimensions respectively. The Mixer employs three consecutive separable convolutions with 96 kernels with depthwise spatial of size (3 x 3 x 3) and a stride of 1. The Grouped block splits the 96 feature maps from the separable block in two groups of 48 and applies two parallel streams of transformations. Each stream comprises two convolutional blocks with 24 and 8 filters of size (3 x 4 x 3) and a stride of (1 x 1 x 1). The clinical features are filtered by two consecutive dense blocks of 32 and 10 units in the FC clinical block. After concatenating the output with the flattened embeddings from the Mixer and applying a dense block with 4 units we have our final input data representation. The feature extractor sub-network has about 550 thousand parameters which is comparatively low even compared to architectures applied on 2D images.

4. Implementation

All experiments were conducted using python version 2.7.12. The neural network was built with the Keras deep learning library using TensorFlow as backend. TensorFlow, which is developed and supported by Google, is an open-source package for numerical computation with high popularity in the deep learning community. The library allows for easy deployment on multiple graphic processing units (GPUs) (CPU-based experimentation would be prohibitive because of time constraints). The Keras wrapper provides an application programming interface (API) for quicker development and has all the functionality to implement the network with the exception of 3D separable convolutions which we built as a custom layer in TensorFlow. In this paper we employed a Linux machine and two Nvidia Pascal TITAN X graphics cards with 12GB RAM each. The model was parallelized across GPUs such that the feature extractor network works on the AD vs HC and MCI-to-AD conversion problems simultaneously to speed up training. Iterating over the whole training set once, i.e. a single epoch, takes about 30 sec and prediction for a single MCI patient requires milliseconds. Since prediction would not require model parallelization or a lengthy training process, a pre-trained network is practical to be applied on a lower-end GPU (or possibly a CPU) relatively cheaply in a realistic scenario. Across all experiments certain network settings remain unchanged. These include the dropout rate - set at 0.1 for all layers and blocks; the L2 regularization penalty coefficient set at $5 \cdot 10^{-5}$ for all parameters in convolutional and fully connected layers; and the convolutional kernel weight initialization which follows the procedure described by He et al. 2015. The objective function loss is minimized using the Adam optimizer by Kingma and Ba, 2014 with an exponentially decaying learning rate:

$$lr = 0.001 * 0.3^{epoch / 10} \quad (8)$$

All other parameters are kept at their default value provided in the original Adam paper (Kingma and Ba, 2014). The network hyperparameters were picked because they led to no overfitting on the validation set during performance evaluation. A batch size of 6 samples for both the AD and MCI conversion problems is randomly sampled from the dataset until it is exhausted.

5. Performance Evaluation

For the evaluation of the classifier, we repeated the sampling strategy to divide the samples in training, validation and test set splits. Since we have 32 samples more in the MCI dataset (including both pMCI and sMCI) as compared to the AD/HC dataset, we used these 32 MCI subjects for testing purposes by randomly sampling 16 subjects from the pMCI and sMCI groups. The validation set comprised roughly 10% of the remaining dataset (36 subjects from MCI and AD/HC respectively) and was also generated by randomly picking in a balanced manner both from the progressive and stable MCI groups and from the healthy and AD patients as we were performing joint learning. Finally, the remaining 340 subjects from both the AD/HC and MCI subsets respectively (i.e. a total of 680 subjects) comprised the training set.

The model is trained for 40 epochs and the best performing model on the validation set is saved and evaluated on the test set. This procedure is then repeated 10 times with different sampling seeds so as to have different samples in the train/validation/test splits and minimize the effect of random variation. The evaluation metrics used are accuracy (ACC), sensitivity (SEN), specificity (SPE). We also perform receiver

operating characteristics (ROC) analysis and compute the AUC. The optimal operating point of the ROC curve was found via Youden's J statistic. All accuracy, sensitivity and specificity results are reported at the optimal operating point of the ROC curve. For the AD vs HC task, we report the validation results as we only defined a test set for the pMCI/sMCI classification problem. The reason is because we are interested pMCI vs sMCI classification and treat the AD/HC task as a helpful auxiliary problem. In addition, we preferred to save as much data for training as possible.

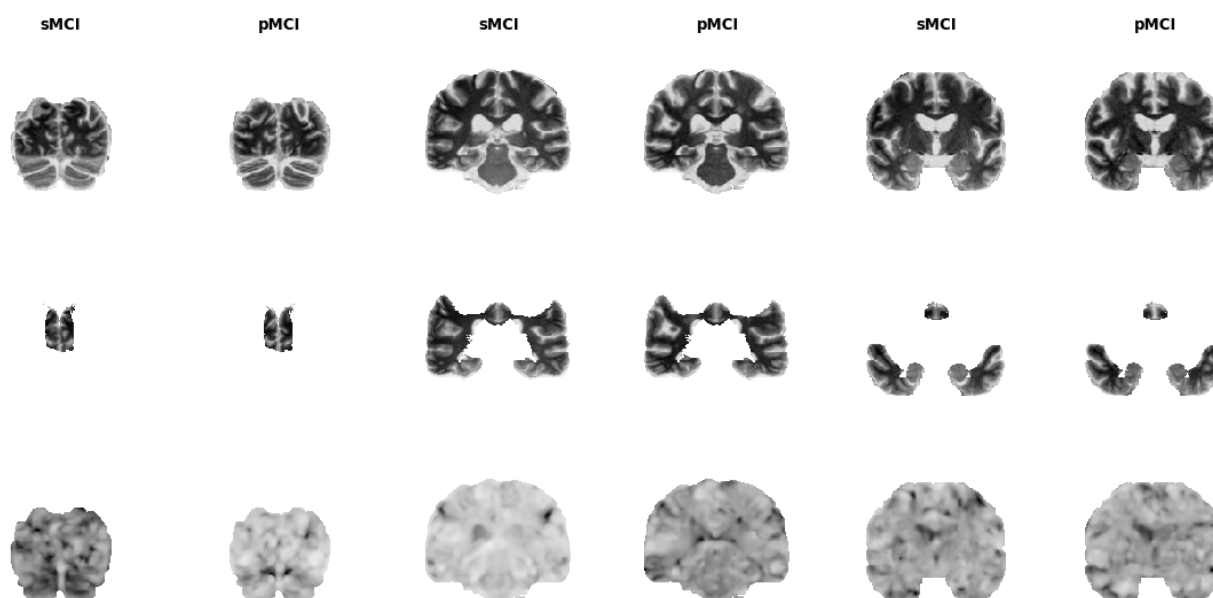


Fig. 4. Examples of the image inputs we employ in the classification framework for three different image slices. The upper row shows structural MRI images co-registered to a custom common space. The middle row displays only the brain regions we retain in the atlas-masked tests (parietal, temporal and frontal lobes). The third row shows the Jacobian Determinant images - they indicate the volumetric change a voxel in an unnormalised MRI image must undergo so as to conform to the common template.

6. Results:

Firstly, we consider the classification performance of our network on four different input biomarker combinations. Then we assess the robustness of the neural network model to MRI structural misalignment by comparing performance metrics obtained when using the custom template as opposed to the MNI152_T1_1mm template. The four input combinations are: 1) clinical features and MRI images; 2) clinical features JD images; 3) clinical features and T1w images; and 4) clinical features, JD and T1w images for a total of 7 experiments: 4 experiment in custom template space and 3 in MNI space (the input combination including the brain-atlas masked MRI was not performed in MNI space). Finally, we assess the performance of our model on the AD vs healthy task with the same input stream variants.

6.1. Multi-modal classification

Results are summarized in fig. 6 and fig. 7 and tables 2 and 3. Fig. 4 shows the retained brain regions after brain-atlas masking.

The best performance metrics are achieved by including structural MRI along with all clinical data (includes demographic, neuropsychological, and APOe4 genotyping features). The median AUC for the input combination comprising structural MRI images and clinical features is 0.92 whereas when we remove brain areas not classically associated with AD, the median AUC obtained is 0.93. Comparing these two values across folds using a Mann-Whitney U test indicated that removing brain structures unrelated to the development of AD does not hinder ($P=0.4$) discrimination in pMCI and sMCI. The median AUC when using JD images and clinical data was at 0.88 (Mann-Whitney test yielded p-value <0.041 and 0.046 when compared to the input combinations comprising structural MRI and clinical data, and atlas-masked structural MRI and clinical data results respectively). Finally, the input combination comprising all types of input

streams - T1w images, the JD data and clinical features resulted in an AUC of 0.91. Comparing this with the input variants comprising the structural MRI and clinical features, atlas-masked MRI and clinical features, or JD images and clinical features yielded p-values of 0.36, 0.38 and 0.07 respectively (Mann-Whitney-U test). These results suggest that adding structural MRI to the clinical features yields statistically significant higher performance as opposed to using only JD data as an image input stream. In addition, removing brain areas from structural MRI not classically associated with Alzheimer's disease did not show statistically different classification results compared to the experiments which retained all information. This suggests our model was not negatively impacted by the inclusion of irrelevant or only partially relevant features.

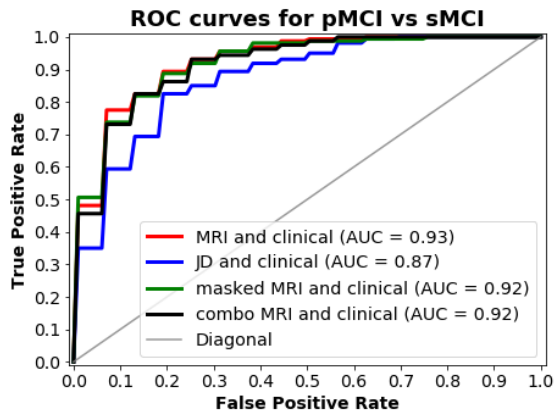


Fig. 5. ROC curves of pMCI vs sMCI classification for four input combinations: MRI images and clinical features; JD images and clinical features; Atlas-masked MRI (or just masked MRI) images and clinical features and finally a MRI; and Jacobian Determinant images and clinical features. The MRI data was co-registered to our custom template prior to performing classification.

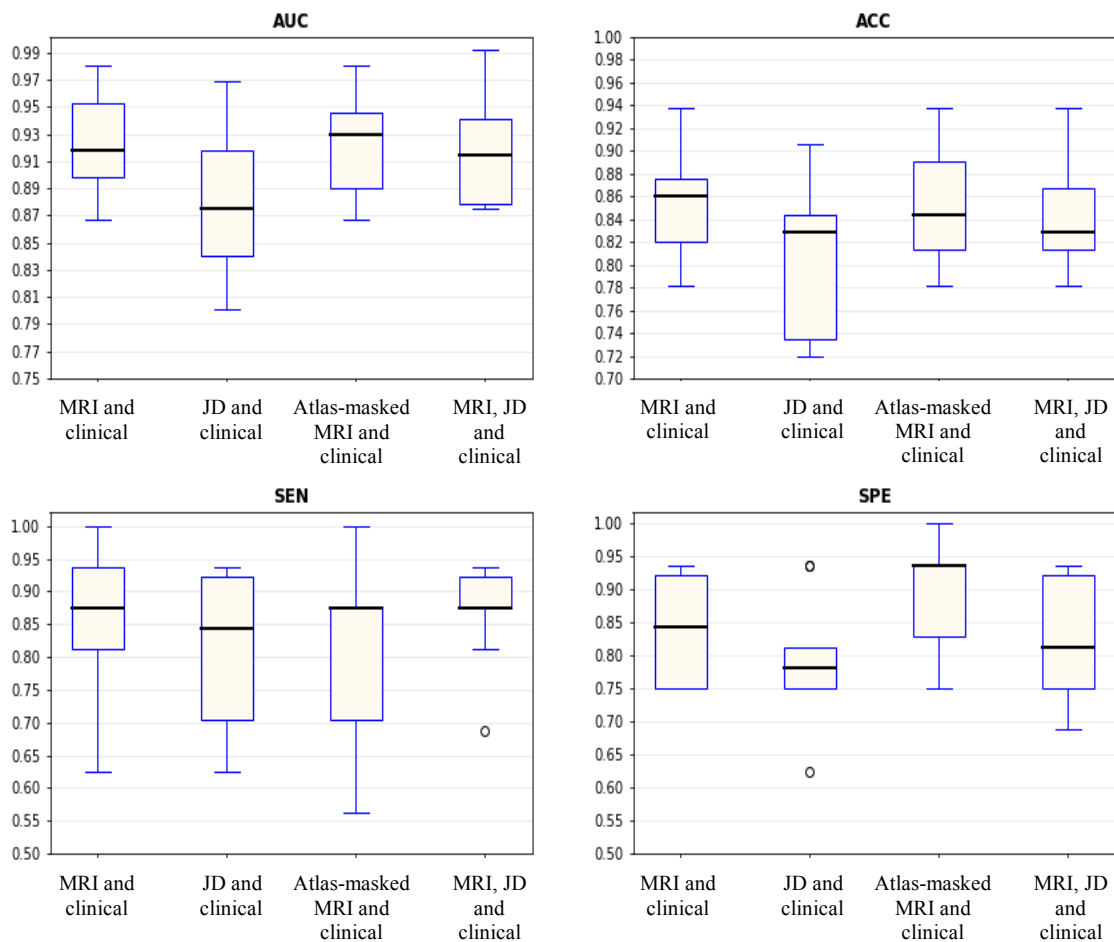


Fig. 6. Box plots for AUC, accuracy, sensitivity and specificity for pMCI vs sMCI classification based on multi-stream integration of clinical features and MRI images (co-registered to our custom template) over 10 separate test folds. The black line in each box represents the median value. The boxes encompass values between the 25th and 75th percentile whereas the tails - the top and bottom quartiles. Outliers are marked with a circle. The performance metrics correspond to the optimal operating point of each classifier.

The highest median classification accuracy we achieved was 86%, which resulted from the experiments with the structural MRI and clinical data. The atlas-masked MRI and clinical data variant yielded the second most predictive power with 84% classification accuracy, whereas the JD images and the clinical features gave 83% accuracy. Finally, employing all input features also resulted in an accuracy of 83%. Across the classification results from our four different input combinations the median sensitivity varies between 85%-87.5%, and the median specificity between 78% and 94% across the test folds.

Tables 2 and 3 summarize the performance metrics obtained from applying our deep learning methodology not only on the pMCI vs sMCI problem but also on AD vs healthy classification. Owing to the simpler nature of AD vs HC discrimination, regardless of the input streams and the co-registration template, results are close to 100% on all performance metrics.

6.2. Network robustness on non-custom template

We measured the classification performance of our deep learning network on MRI data co-registered to the Montreal Neurological Institute (MNI152) template instead of the custom normalization space. The purpose of these experiments is to assess the robustness of the methodology to possible structural misalignment in the brain areas across images as the MNI space is more "distant" from the images under study.

Results are summarized in fig. 8 and tables 2 and 3. In order to identify performance differences between using our custom and the MNI152 templates, we performed Mann Whitney U tests across folds on the obtained AUCs corresponding to the different input combination pairs (custom template vs MNI template). The obtained p-values are 0.28, 0.42 and 0.24 for the structural MRI and clinical features, Jacobian Determinants and clinical features, and the combined inputs respectively. Consequently, no statistically significant difference can be found between the performance of our classifier while operating in the two normalization spaces.

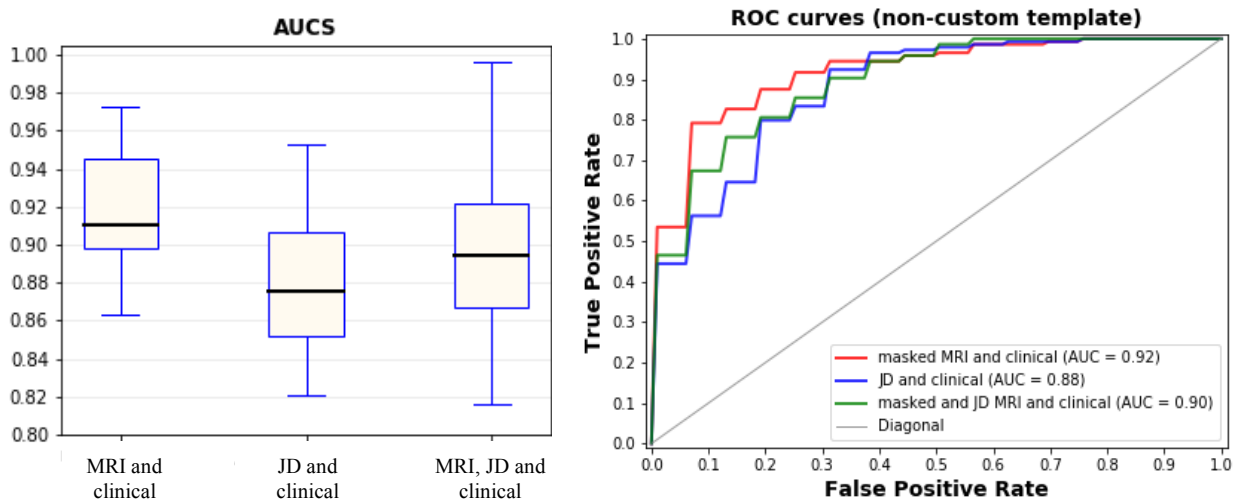


Fig. 7. Box plots for AUC, accuracy, sensitivity and specificity obtained on the pMCI vs sMCI classification task from structural MRI, Jacobian Determinant and atlas-masked structural MRI inputs (all using clinical features) over 10 separate test folds. The MRI data is co-registered to the MNI(152) template. The black line in each box represents the median value. The boxes encompass values between the 25th and 75th percentile whereas the tails - the top and bottom quartiles. Outliers are marked with a circle. The performance metrics correspond to the optimal operating point of each classifier.

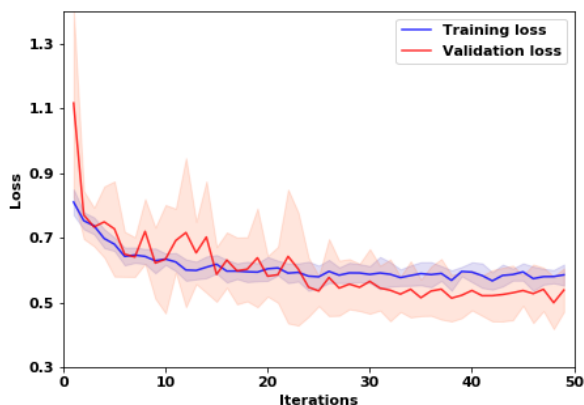


Fig. 8. A plot of the training and validation losses for our CNN architecture which utilises structural MRI and clinical features. The standard deviation of the validation loss encompasses the red area in the image, whereas the deviation of the training loss is depicted in blue.

6.3. Classification variance

Although we achieve high median performance on all metrics and on both registration templates, dispersion can be further reduced. We have plotted the standard deviation of the mean training and validation losses across the 10 test folds of the model which achieved the highest classification accuracy on fig. 8 (structure MRI and clinical features co-registered to our custom template).

One factor which contributes to the higher validation variance compared to the training loss curve is the number of samples. Since both the validation and test sets comprise an order of magnitude less subjects than the training set, we also expect them to have higher variance. Secondly, although the CNN weights were optimized using a variant of stochastic gradient descent, the hyper parameters, such as the dropout rate, the L2 regularization hyper parameter, the initial learning rate and learning rate decay were set to pre-defined values which gave good performance on only one of the validation folds. This was done as performing hyper parameter search was deemed prohibitive given the number of experiments we performed. Consequently, as the dataset is relatively small, we observed some level of overfitting or bias, depending on the specific data split employed, which indicates room for improvement on our current results. High performance metric variance is most prevalent in the sensitivity and specificity box plots since they are calculated only using either the true positives or true negatives, i.e. half the test set. Accordingly, some studies (Moradi et al. 2015, Hojjati et al. 2017, Tong et al. 2017) repeat their cross-validation loops many times (such as 100 or a 1000 times) in order to further reduce their performance variance, which was also presumed to be prohibitive for computational reasons.

7. Discussion:

Deep learning, or deep neural networks, works by extracting a hierarchy of features from the input data via flexible non-linear transformations. These new data representations are learnt such that they maximize an arbitrary performance metric, for example classification accuracy. Hence, instead of relying on expert prior knowledge, or other dimensionality reduction algorithms which might result in a non-optimal set of features, deep neural networks use the gradient in the performance metric to directly guide the feature extraction mechanism. This can result in significant improvements in classification results. Additionally, given that the feature representations are built in a multi-layered fashion (where higher level features are derived from lower level ones), articulate and information-rich images and volumes can be dealt with and incorporated easily into the classification process.

In this paper, we developed a new method with the primary goal of early identification of MCI patients with high risk of converting to Alzheimer’s disease up to three years prior to diagnosis, and the subsidiary task of Alzheimer’s patient vs. healthy control discrimination. Our approach uses a parameter-efficient deep convolutional neural network framework, inspired by grouped and separable convolutions, to extract descriptive factors from structural MRI images acquired at baseline. In this respect our work differs from previous deep learning-based methods of early AD detection in that it takes into consideration data paucity in medical datasets and introduces design precautions by reducing the number of network parameters. This in turn increases the generalization capabilities (i.e. reduces overfitting) of our model to unseen test samples, thus enabling us to achieve state-of-the-art MCI-to-AD classification performance. The structural MRI images are complemented by standard cognitive test results (CDRSB, ADAS, RAVLT), demographic information (age, gender, ethnic and racial categories, education) and APOe4 expression levels also acquired at baseline to arrive at a final score which is used to predict conversion. We chose these biomarkers in order to create a classification methodology which is as minimally invasive as possible. Hence, for example, we do not include PET imaging because of radiation exposure and CSF data owing to the potentially painful lumbar puncture which can also lead to clinical complications. Additionally, we exploited AD/HC data to limit the effects of overfitting. This was achieved by multi-task learning where the same network layers are used to extract representations from the input biomarkers for both the MCI-to-AD conversion task and the AD/HC classification problem. While previous methods employ pre-training (Payan et al. 2015; Hosseini-Asl et al. 2016; Liu et al. 2018; Liu et al. 2015) to reap similar

pMCI vs sMCI								
Input Modalities	Custom template				MNI152 template			
	AUC	ACC	SEN	SPE	AUC	ACC	SEN	SPE
MRI and clinical	0.92	86%	87.5%	84%	0.91	85%	82%	87%
Atlas-masked MRI and clinical	0.93	84%	87.5%	94%	-	-	-	-
JD and clinical	0.88	83%	84%	78%	0.88	82%	82%	81%
MRI and JD and clinical	0.91	83%	87.5%	81%	0.90	83%	77%	88%

Table 2. A comparison table between the median performance metrics on the pMCI vs sMCI classification task using our neural network model.

AD vs HC								
Input Modalities	Custom template				MNI152 template			
	AUC	ACC	SEN	SPE	AUC	ACC	SEN	SPE
MRI and clinical	1	99.5%	100%	99%	1	99.5%	100%	99%
Atlas-masked MRI and clinical	1	99.5%	100%	99%	-	-	-	-
JD and clinical	1	99%	99.5%	99%	0.99	97%	95%	99%
masked and JD and clinical	1	99.5%	100%	99%	1	99%	99%	99%

Table 3. A comparison table between the median performance metrics on the AD vs healthy classification task using our neural network model.

Author	Data	AUC	ACC	SEN	SPE	Conversion time	Validation and Testing method
Spasov et al. (this paper)	structural MRI + cognitive measures + APOe4 + demographics	0.93	86%	87.5%	85%	0-36 months	10-fold cross-validation
Hojjati et al.	rs-fMRI	0.95	91.4%	83.24%	90.1%	0-36 months	9-fold cross-validation (report on validation set)
Moradi et al.	structural MRI + cognitive measures	0.9	82%	87%	74%	0-36 months	10-fold cross-validation (report on test set)
Liu et al. (Cox)	structural MRI + FDG-PET + cognitive measures + APOe4 + demographics	0.92	84.6%	86.5%	82.4%	0-36 months	holdout
Korolev et al.	structural MRI + clinical data + plasma-proteomic data + medications	0.87	80%	83%	76%	0-36 months	10-fold cross-validation (report on test set)
Beheshti et al.	structural MRI	75.08	75%	77%	73%	0-36 months	10-fold cross-validation
Choi et al., 2018	flurodeoxyglucose and florbetapir PET	0.89	84.2%	81%	87%	0-36 months	holdout
Tong et al., 2017	structural MRI + cognitive measures	0.92	84%	88.7%	76.5%	0-36 months	10-fold cross-validation (report on
D. Lu et al. 2018 (1)	FDG-PET	-	82.5%	81.4%	83%	0-36 months	10-fold cross-validation

Table 3. A comparative table of methodologies on the pMCI vs sMCI classification task using the ADNI dataset. We provide a performance comparison table mainly for recent studies achieving classification rates close to the state-of-the-art.

benefits, this requires training the model twice, whereas dual-learning is a single-stage procedure. Our experimental procedures assess the performance of our method using two different co-registration templates (a custom one and the MNI152) as well as various input combinations of structural MRI, the local JD of the deformational field applied during MRI co-registration, as well as the clinical data. The best result we obtained was a mean AUC of 0.93 averaged across 10 different testing folds with a mean MCI-to-AD conversion prediction accuracy of 86%, sensitivity of 87.5% and specificity of 85% (see table 3). It is also important to note that, to the best of our knowledge, the only study which presents better classification results on the pMCI vs sMCI problem (Hojjati et al. 2017) does not explicitly mention the use of separate test, validation and training sets, possibly leading to double-dipping. They report their results on a validation set instead of a dedicated test set.

The main novelties of our method were 1) the use of parameter-efficient layers, such as grouped and separable convolutions (implemented as custom Keras layers for 3D inputs) to reduce the number of network parameters, hence limiting overfitting; 2) the substitution of network pre-training, which was typical in earlier deep-learning based AD classification studies (Payan et al. 2015, Hosseini-Asl et al. 2016), with multi-task learning which utilizes AD/HC data to arrive at a single-stage training approach and 3) the utilization of the JD as a complementary imaging input stream to maximize the extracted information from the structural MRI. Convolutional neural networks abstract away the manual handcrafting of useful features from medical images, such as the use of pre-defined brain regions of interest (Da et al. 2013). Intuitively, neural network-based methods should perform better as the feature extraction process is directly driven by the performance optimization procedure, however, it comes at the cost of a relatively high number of network parameters compared to the number of samples. Since there are no formal estimates of the number of training samples required for a given convolutional architecture to achieve good generalization, we are driven by the metaheuristic approach of minimizing the number of network weights and maximizing the effective number of training examples so as to

boost performance on an independent test set and consequently during clinical application. As a result, our 3D model comprises ~500,000 parameters, which is orders of magnitude lower than conventional 3D CNNs and even lower than even recent 2D CNNs. This was not done by sacrificing network depth or structural complexity but rather by inserting efficient convolutional layers. In order to facilitate the learning procedure, we hypothesized that employing an auxiliary task and minimizing the joint training objective of the MCI-to-AD conversion and AD/HC classification tasks would be an effective alternative to pre-training. In this context, AD/NC discrimination is seen as a simpler version of MCI conversion prediction, and in order to speed up training convergence we worked under the assumption that similar descriptive factors would be useful for both problems.

We also performed experiments to assess the robustness of our method to two different co-registration templates (MNI_152_1mm and a custom, study specific template), aiming to a) identify the most predictive combination of input data and assess how our model handles irrelevant features, and b) identify whether the image co-registration method had a statistically significant impact on the classification accuracies we obtain. Firstly, we found that the classification performance of our network over 10 different test folds was not statistically different between the two co-registration procedures for corresponding input combinations, with Mann-Whitney U test p-values in the 0.24-0.42 range. Secondly, we found that the most predictive input combination comprised structural MRI and clinical data. All input variants other than the JD combined with clinical data performed at a similar and statistically indistinguishable level, achieving prediction accuracies in the range 83%-86%, sensitivity at 87.5% and specificity in the range of 81%-94%. Finally, as the input variants including the brain atlas-masked MRI and non-masked MRI performed at the same level, we concluded that the inclusion of irrelevant or partially relevant features does not hinder the performance of our model.

Considering existing computer vision research, deep learning methodologies for computer-aided diagnostics would also be applicable on non-co-registered or even non-pre-processed images, however, this approach could lead to image artefacts contributing to the discriminatory performance of the algorithm, which could learn to relate center-specific (rather than disease-specific) features with disease outcomes. As with all multicentric studies, careful and unified data collection and processing is crucial to minimize this confound.

Comparing our classification metrics with recent studies indicate that only [Hojjati et al. 2017](#) et al. who use rs-fMRI outperform our results (although, as mentioned above, only reporting on a validation set comprising 4 subjects via 9-fold cross-validation). Unfortunately, at the time of writing ADNI provides limited rs-fMRI data (18 pMCI and 62 sMCI subjects) so it would be difficult to predict how their results would scale to larger populations. Additionally, using structural MRI only can significantly reduce in-patient scanner time as opposed to including a functional scan. To the best of our knowledge, the study by [Liu et al. 2017](#) is the first to produce comparable performance (at least in some metrics) to our model, at 84.6% classification accuracy vs 86% for our work. The difference is, however, that [Liu et al. 2017](#) utilize FDG-PET as an extra modality which is known to be extremely informative in AD, as well as structural MRI and all the biomarkers we have employed. [Moradi et al. 2015](#) and [Tong et al. 2017](#) both use a very similar methodology to each other and the same data (structural MRI and cognitive assessment tests) as in this paper. Their sensitivity metrics are comparable to our model at ~87%-88% but manifest lower specificity at 74% and 76% respectively, while our deep learning method achieves 85%. A possible explanation would be the inclusion of APOe4 and demographic data as well as the efficacy of the neural network. Also, as is discussed in [Moradi et al. 2015](#) the labelling of subjects varies across studies, thus hampering direct comparisons.

In summary, we developed a deep learning-based method for the early prediction of MCI-to-AD converts by combining structural MRI, neuropsychological assessment data and APOe4 expression levels obtained from the ADNI database at baseline. We achieved a very high predictive performance with an average AUC of 0.93, prediction accuracy of 86%, sensitivity of 87.5% and specificity of 85%. Our study proposes the use of more efficient neural network architectures comprising fewer parameters to limit the effects of overfitting. The convolutional framework is generic and applicable to any 3D image dataset and gives the flexibility to design a computer-aided diagnosis system targeting the prediction of any medical condition utilizing multi-modal imaging and tabular clinical data.

8. Acknowledgements

In this work we employed the database of the Alzheimer's Disease Neuroimaging Initiative (ADNI). ADNI was formed as a multicenter longitudinal study to identify imaging, clinical, genetic and biochemical biomarkers for the early detection and tracking of Alzheimer's disease (AD) and Mild Cognitive Impairment (MCI). ADNI is the result of a \$67 million partnership by the public and private sector. Financial support was obtained from the National Institute on Ageing, 13 pharmaceutical companies, and two foundations providing funding through the Foundation for the National Institutes of Health. The study can be split in three sub-initiatives - ADNI1, ADNI2 and ADNI GO. The initial phase known as ADNI1 included subjects between 55-90 years of age from approximately 50 sites from the US and Canada. ADNI2 and ADNI GO add new participants and funding to the study. The database is made available to researchers around the world and has a broad range of collaborators. The principle investigator of ADNI, who oversees all aspects, is Dr. Michael Weiner, MD, VA Medical Center and University of California - San Francisco. For up-to-date information, see www.adni-info.org. This research is supported by the Engineering and Physical Sciences Research Council [EP/ L015889/1].

References:

- M. Abadi et al., TensorFlow: A System for Large-Scale Machine Learning, Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI). Nov 2016.
- Chollet et al., Keras, 2015, available online at : <https://keras.io>, last accessed: 11.08.2018
- D.E. Barnes, K. Yaffe, The projected effect of risk factor reduction on Alzheimer's disease prevalence, *The Lancet Neurology*. 10 (2011) 819–828. doi:10.1016/s1474-4422(11)70072-2.
- C.P. Ferri, M. Prince, C. Brayne, H. Brodaty, L. Fratiglioni, M. Ganguli, K. Hall, K. Hasegawa, H. Hendrie, Y. Huang, A. Jorm, C. Mathers, P.R. Menezes, E. Rimmer, M. Scuzofca, Global prevalence of dementia: a Delphi consensus study, *The Lancet*. 366 (2005) 2112–2117. doi:10.1016/s0140-6736(05)67889-0.
- L. Mosconi, M. Brys, L. Glodzik-Sobanska, S. De Santi, H. Rusinek, M.J. de Leon, Early detection of Alzheimer's disease using neuroimaging, *Experimental Gerontology*. 42 (2007) 129–138. doi:10.1016/j.exger.2006.05.016.
- M. Paul Murphy, Harry LeVine, Alzheimer's Disease and the Amyloid- β Peptide, *JAD*. 19 (2010) 311–323. doi:10.3233/JAD-2010-1221.
- Mitchell, A.J., Shiri-Feshki, M., 2008. Temporal trends in the long term risk of progression of mild cognitive impairment: a pooled analysis. *Journal of Neurology, Neurosurgery & Psychiatry* 79, 1386–1391. <https://doi.org/10.1136/jnnp.2007.142679>
- H. Braak, E. Braak, Staging of alzheimer's disease-related neurofibrillary changes, *Neurobiology of Aging*. 16 (1995) 271–278. doi:10.1016/0197-4580(95)00021-6.
- H. Braak, E. Braak, Development of Alzheimer-related neurofibrillary changes in the neocortex inversely recapitulates cortical myelogenesis, *Acta Neuropathologica*. 92 (1996) 197–201. doi:10.1007/s004010050508.
- A. Delacourte, J.P. David, N. Sergeant, L. Buee, A. Wattez, P. Vermersch, F. Ghazali, C. Fallet-Bianco, F. Pasquier, F. Lebert, H. Petit, C. Di Menza, The biochemical pathway of neurofibrillary degeneration in aging and Alzheimer's disease, *Neurology*. 52 (1999) 1158–1158. doi:10.1212/wnl.52.6.1158.
- J.C. Morris, M. Storandt, D.W. McKeel, E.H. Rubin, J.L. Price, E.A. Grant, L. Berg, Cerebral amyloid deposition and diffuse plaques in "normal" aging: Evidence for presymptomatic and very mild Alzheimer's disease, *Neurology*. 46 (1996) 707–719. doi:10.1212/wnl.46.3.707.
- A. Serrano-Pozo, M.P. Frosch, E. Masliah, B.T. Hyman, Neuropathological Alterations in Alzheimer Disease, *Cold Spring Harbor Perspectives in Medicine*. 1 (2011) a006189–a006189. doi:10.1101/cshperspect.a006189.
- William R. Markesbery, Neuropathologic Alterations in Mild Cognitive Impairment: A Review, *JAD*. 19 (2010) 221–228. doi:10.3233/JAD-2010-1220.
- H. Hampel, N. Toschi, F. Baldacci, H. Zetterberg, K. Blennow, I. Kilimann, S.J. Teipel, E. Cavado, A. Melo dos Santos, S. Epelbaum, F. Lamari, R. Genthon, B. Dubois, R. Floris, F. Garaci, S. Lista, Alzheimer's disease biomarker-guided diagnostic workflow using the added value of six combined cerebrospinal fluid candidates: A β 1–42, total-tau, phosphorylated-tau, NFL, neurogranin, and YKL-40, *Alzheimer's & Dementia*. 14 (2018) 492–501. doi:10.1016/j.jalz.2017.11.015.
- Baldacci, F., Lista, S., O'Bryant, S.E., Ceravolo, R., Toschi, N., Hampel, H., 2018. Blood-Based Biomarker Screening with Agnostic Biological Definitions for an Accurate Diagnosis Within the Dimensional Spectrum of Neurodegenerative Diseases, in: Biomarkers for Alzheimer's Disease Drug Development. Springer New York, pp. 139–155. https://doi.org/10.1007/978-1-4939-7704-8_9
- Teipel, S.J., Cavado, E., Lista, S., Habert, M.-O., Potier, M.-C., Grothe, M.J., Epelbaum, S., Sambati, L., Gagliardi, G., Toschi, N., Greicius, M.D., Dubois, B., Hampel, H., Audrain, C., Auffret, A., Bakardjian, H., Baldacci, F., Batrancourt, B., Benakki, I., Benali, H., Bertin, H., Bertrand, A., Boukadida, L., Cacciamani, F., Causse, V., Cavado, E., Cherif Touil, S., Chiesa, P.A., Colliot, O., Dalla Barba, G., Depaulis, M., Dos Santos, A., Dubois, B., Dubois, M., Epelbaum, S., Fontaine, B., Francisque, H., Gagliardi, G., Genin, A., Genthon, R., Glasman, P., Gombert, F., Habert, M.O., Hampel, H., Hewa, H., Houot, M., Jungalee, N., Kas, A., Kilani, M., La Corte, V., Le Roy, F., Lehericy, S., Letondor, C., Levy, M., Lista, S., Lowrey, M., Ly, J., Makiese, O., Masetti, I., Mendes, A., Metzinger, C., Michon, A., Mochel, F., Nait Arab, R., Nyasse, F., Perrin, C., Poirier, F., Poisson, C., Potier, M.C., Ratovohery, S., Revillon, M., Rojkova, K., Santos-Andrade, K., Schindler, R., Servera, M.C., Seux, L., Simon, V., Skovronsky, D., Thiebaut, M., Uspenskaya, O., Vlaincu, M., 2018. Effect of Alzheimer's disease risk and protective factors on cognitive trajectories in subjective memory complainers: An INSIGHT-preAD study. *Alzheimer's & Dementia*. <https://doi.org/10.1016/j.jalz.2018.04.004>
- J.L. Shaffer, J.R. Petrella, F.C. Sheldon, K.R. Choudhury, V.D. Calhoun, R.E. Coleman, P.M. Doraiswamy, Predicting Cognitive Decline in Subjects at Risk for Alzheimer Disease by Using Combined Cerebrospinal Fluid, MR Imaging, and PET Biomarkers, *Radiology*. 266 (2013) 583–591. doi:10.1148/radiol.12120010.
- R. Filipovych, C. Davatzikos, Semi-supervised pattern classification of medical images: Application to mild cognitive impairment (MCI), *NeuroImage*. 55 (2011) 1109–1119. doi:10.1016/j.neuroimage.2010.12.066.
- L. Mosconi, D. Perani, S. Sorbi, K. Herholz, B. Nacmias, V. Holthoff, E. Salmon, J.-C. Baron, M.T.R. De Cristofaro, A. Padovani, B. Borroni, M. Franceschi, L. Bracco, A. Pupi, MCI conversion to dementia and the APOE genotype: A prediction study with FDG-PET, *Neurology*. 63 (2004) 2332–2340. doi:10.1212/01.wnl.0000147469.18313.3b.
- H. Choi, K.H. Jin, Predicting cognitive decline with deep learning of brain metabolism and amyloid imaging, *Behavioural Brain Research*. 344 (2018) 103–109. doi:10.1016/j.bbr.2018.02.017.
- J. Young, M. Modat, M.J. Cardoso, A. Mendelson, D. Cash, S. Ourselin, Accurate multimodal probabilistic prediction of conversion to Alzheimer's disease in patients with mild cognitive impairment, *NeuroImage: Clinical*. 2 (2013) 735–745. doi:10.1016/j.nicl.2013.05.004.

- E. Moradi, A. Pepe, C. Gaser, H. Huttunen, J. Tohka, Machine learning framework for early MRI-based Alzheimer's conversion prediction in MCI subjects, *NeuroImage*. 104 (2015) 398–412. doi:10.1016/j.neuroimage.2014.10.002.
- T. Tong, Q. Gao, R. Guerrero, C. Ledig, L. Chen, D. Rueckert, A.D.N. Initiative, A Novel Grading Biomarker for the Prediction of Conversion From Mild Cognitive Impairment to Alzheimer's Disease, *IEEE Transactions on Biomedical Engineering*. 64 (2017) 155–165. doi:10.1109/tbme.2016.2549363.
- R. Casanova, C.T. Whitlow, B. Wagner, J. Williamson, S.A. Shumaker, J.A. Maldjian, M.A. Espeland, High Dimensional Classification of Structural MRI Alzheimer's Disease Data Based on Large Scale Regularization, *Frontiers in Neuroinformatics*. 5 (2011). doi:10.3389/fninf.2011.00022.
- Joshua A. Sonnen, Kathleen S. Montine, Joseph F. Quinn, John C.S. Breitner, Thomas J. Montine, Cerebrospinal Fluid Biomarkers in Mild Cognitive Impairment and Dementia, *JAD*. 19 (2010) 301–309. doi:10.3233/JAD-2010-1236.
- M. Riemenschneider, N. Lautenschlager, S. Wagenpfeil, J. Diehl, A. Drzezga, A. Kurz, Cerebrospinal Fluid Tau and β -Amyloid 42 Proteins Identify Alzheimer Disease in Subjects With Mild Cognitive Impairment, *Archives of Neurology*. 59 (2002) 1729. doi:10.1001/archneur.59.11.1729.
- O. Hansson, H. Zetterberg, P. Buchhave, E. Londos, K. Blennow, L. Minthon, Association between CSF biomarkers and incipient Alzheimer's disease in patients with mild cognitive impairment: a follow-up study, *The Lancet Neurology*. 5 (2006) 228–234. doi:10.1016/s1474-4422(06)70355-6.
- C. Davatzikos, P. Bhatt, L.M. Shaw, K.N. Batmanghelich, J.Q. Trojanowski, Prediction of MCI to AD conversion, via MRI, CSF biomarkers, and pattern classification, *Neurobiology of Aging*. 32 (2011) 2322.e19-2322.e27. doi:10.1016/j.neurobiolaging.2010.05.023.
- Y. Fan, D. Shen, R.C. Gur, R.E. Gur, C. Davatzikos, COMPARE: Classification of Morphological Patterns Using Adaptive Regional Elements, *IEEE Transactions on Medical Imaging*. 26 (2007) 93–105. doi:10.1109/tmi.2006.886812.
- M.H. Nguyen, F. de la Torre, Optimal feature selection for support vector machines, *Pattern Recognition*. 43 (2010) 584–591. doi:10.1016/j.patcog.2009.09.003.
- D. Lu, K. Popuri, G.W. Ding, R. Balachandar, M.F. Beg, Multiscale deep neural network based analysis of FDG-PET images for the early diagnosis of Alzheimer's disease, *Medical Image Analysis*. 46 (2018) 26–34. doi:10.1016/j.media.2018.02.002. ——— 1
- D. Lu, K. Popuri, G.W. Ding, R. Balachandar, M.F. Beg, Multimodal and Multiscale Deep Neural Networks for the Early Diagnosis of Alzheimer's Disease using structural MR and FDG-PET images, *Scientific Reports*. 8 (2018). doi:10.1038/s41598-018-22871-z. ——— 2
- M. Liu, D. Cheng, K. Wang, Y. Wang, Multi-Modality Cascaded Convolutional Neural Networks for Alzheimer's Disease Diagnosis, *Neuroinformatics*. (2018). doi:10.1007/s12021-018-9370-4.
- K. Liu, K. Chen, L. Yao, X. Guo, Prediction of Mild Cognitive Impairment Conversion Using a Combination of Independent Component Analysis and the Cox Model, *Frontiers in Human Neuroscience*. 11 (2017). doi:10.3389/fnhum.2017.00033.
- Adrien Payan, Giovanni Montana, Predicting Alzheimer's disease: a neuroimaging study with 3D convolutional neural networks, *ICPRAM*, 2015
- E. Hosseini-Asl, R. Keynton, A. El-Baz, Alzheimer's disease diagnostics by adaptation of 3D convolutional network, in: 2016 IEEE International Conference on Image Processing (ICIP), IEEE, 2016. doi:10.1109/icip.2016.7532332.
- Avants, B.B., Yushkevich, P., Pluta, J., Minkoff, D., Korczykowski, M., Detre, J., Gee, J.C., 2010. The optimal template effect in hippocampus studies of diseased populations. *NeuroImage* 49, 2457–2466. <https://doi.org/10.1016/j.neuroimage.2009.09.062>
- Avants, B.B., Tustison, N.J., Song, G., Cook, P.A., Klein, A., Gee, J.C., 2011. A reproducible evaluation of ANTs similarity metric performance in brain image registration. *NeuroImage* 54, 2033–2044. <https://doi.org/10.1016/j.neuroimage.2010.09.025>
- Jenkinson, M., Beckmann, C.F., Behrens, T.E.J., Woolrich, M.W., Smith, S.M., 2012. FSL. *NeuroImage* 62, 782–790. <https://doi.org/10.1016/j.neuroimage.2011.09.015>
- I.O. Korolev, L.L. Symonds, A.C. Bozoki, Predicting Progression from Mild Cognitive Impairment to Alzheimer's Dementia Using Clinical, MRI, and Plasma Biomarkers via Probabilistic Pattern Classification, *PLOS ONE*. 11 (2016) e0138866. doi:10.1371/journal.pone.0138866.
- K. Liu, K. Chen, L. Yao, X. Guo, Prediction of Mild Cognitive Impairment Conversion Using a Combination of Independent Component Analysis and the Cox Model, *Frontiers in Human Neuroscience*. 11 (2017). doi:10.3389/fnhum.2017.00033.
- F. Chollet, Xception: Deep Learning with Depthwise Separable Convolutions, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, 2017. doi:10.1109/cvpr.2017.195.
- A. Krizhevsky, I. Sutskever, G.E. Hinton, ImageNet classification with deep convolutional neural networks, *Communications of the ACM*. 60 (2017) 84–90. doi:10.1145/3065386.
- Velickovic, P., Wang, D., Lane, N.D., Lio, P., 2016. X-CNN: Cross-modal convolutional neural networks for sparse datasets. 2016 IEEE Symposium Series on Computational Intelligence (SSCI). doi:10.1109/ssci.2016.7849978
- S. Ioffe and C. Szegedy, 2015. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. 2015 International Conference on Machine Learning (ICML), p. 448-456
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, Ruslan Salakhutdinov, Yoshua Bengio, Dropout: A Simple Way to Prevent Neural Networks from Overfitting, *Journal of Machine Learning Research (JMLR)*, p. 1929-1958 (2014)
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>. (accessed: 30.05.2018)

- Djork-Arné Clevert, Thomas Unterthiner, Sepp Hochreiter, Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs). 2015, CoRR abs/1511.07289
- He, K., Zhang, X., Ren, S., Sun, J., 2015. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. 2015 IEEE International Conference on Computer Vision (ICCV). doi:10.1109/iccv.2015.123
- D.Kingma and J. Ba, 2014. Adam: A Method for Stochastic Optimization. CoRR abs/1412.6980 (2014)
- S.H. Hojjati, A. Ebrahimzadeh, A. Khazaei, A. Babajani-Feremi, Predicting conversion from MCI to AD using resting-state fMRI, graph theoretical approach and SVM, *Journal of Neuroscience Methods*. 282 (2017) 69–80. doi:10.1016/j.jneumeth.2017.03.006.
- I. Beheshti, H. Demirel, H. Matsuda, Classification of Alzheimer’s disease and prediction of mild cognitive impairment-to-Alzheimer’s conversion from structural magnetic resource imaging using feature ranking and a genetic algorithm, *Computers in Biology and Medicine*. 83 (2017) 109–119. doi:10.1016/j.compbiomed.2017.02.011.
- A. Hammers, R. Allom, M.J. Koeppe, S.L. Free, R. Myers, L. Lemieux, T.N. Mitchell, D.J. Brooks, J.S. Duncan, Three-dimensional maximum probability atlas of the human brain, with particular reference to the temporal lobe, *Human Brain Mapping*. 19 (2003) 224–247. doi:10.1002/hbm.10123.
- Da, X., Toledo, J.B., Zee, J., Wolk, D.A., Xie, S.X., Ou, Y., Shacklett, A., Parmpi, P., Shaw, L., Trojanowski, J.Q., Davatzikos, C., 2014. Integration and relative value of biomarkers for prediction of MCI to AD progression: Spatial patterns of brain atrophy, cognitive scores, APOE genotype and CSF biomarkers. *NeuroImage: Clinical* 4, 164–173. doi:10.1016/j.nicl.2013.11.010