

Obtaining leaner deep neural networks for decoding brain functional connectome in a single shot



Sukrit Gupta¹, Yi Hao Chan¹, Jagath C. Rajapakse*, The Alzheimer's Disease Neuroimaging Initiative

School of Computer Science and Engineering, Nanyang Technological University, 639798, Singapore

ARTICLE INFO

Article history:

Received 30 November 2019

Revised 29 February 2020

Accepted 16 April 2020

Available online 30 January 2021

Keywords:

Alzheimer's disease

Attention deficit hyperactivity disorder

Brain decoding

Deep neural networks

Feature selection

Major depressive disorder

Mild cognitive impairment

ABSTRACT

Neuroscientific knowledge points to the presence of redundancy in the correlations of the brain's functional activity. These redundancies can be removed to mitigate the problem of overfitting when deep neural network (DNN) models are used to classify neuroimaging datasets. We propose an algorithm that removes insignificant nodes of DNNs in a layerwise manner and then adds a subset of correlated features in a single shot. When performing experiments with functional MRI datasets for classifying patients from healthy controls, we were able to obtain simpler and more generalizable DNNs. The obtained DNNs maintained a similar performance as the full network with only around 2% of the initial trainable parameters. Further, we used the trained network to identify salient brain regions and connections from functional connectome for multiple brain disorders. The identified biomarkers were found to closely correspond to previously known disease biomarkers. The proposed methods have cross-modal applications in obtaining leaner DNNs that seem to fit neuroimaging data better. The corresponding code is available at https://github.com/SCSE-Biomedical-Computing-Group/LEAN_CLIP.

© 2021 Elsevier B.V. All rights reserved.

1. Introduction

Deep neural networks (DNN) have been successfully applied to a wide variety of classification problems, ranging from image recognition to language processing [21,28]. In these applications, the choice of network architecture and the number of trainable parameters affects network performance significantly. Larger networks learn more complex and accurate mappings, which comes with increased computational complexity and possibilities of overfitting [48]. On the other hand, having insufficient number of parameters limits the network's ability to learn the correct mapping [60]. Therefore, finding optimal neural network architecture is an important problem. One possible solution involves performing elimination on a large neural network by removing redundant nodes and connections. This has been shown to produce comparable performance resulting in a smaller network but with better generalization capability [1,4].

Neural networks and machine learning techniques have increasingly been successfully applied to neuroimaging data. Non-invasive neuroimaging techniques are used to characterize functional and structural anomalies in the brain and aid in better diagnosis and treatment [13]. Functional and/or structural connectome

derived from these techniques are used as features for neural network models to classify diseased and normal subjects [7,15]. However, such studies typically involve many more input features (~ 10,000) than subject scans (~ 1,000), making neural networks prone to overfitting. Also, experimental noise introduces systematic connectivity [34], which should be ignored during the classification task. Most crucially, overfitting is exacerbated by the presence of redundancies: not all functional connectivity features are important for differentiating between scan samples from two subject groups.²

In recent works [16], the authors used feature salience scores to remove less salient features. In order to find salient features, DeepLIFT [44] was used to compute salience scores of both input features and hidden layer nodes. DeepLIFT is one of the recent attempts (besides Integrated Gradients [50], Layerwise Relevance Propagation [5] and SHAP [30]) that circumvent the issues with gradient-based approaches (such as zero gradients or discontinuities). These approaches find the contribution of nodes at each layer by propagating the contributions from the output layer. A

² Data used in preparation of this article were obtained from the ADNI database (adni.loni.usc.edu). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete list of the ADNI investigators can be found at http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf.

* Corresponding author.

¹ These authors contributed equally.

major differentiating factor of DeepLIFT is that it gives consideration to both negative and positive contributions and computes the relevance scores efficiently in a single pass. By recursively removing the least salient features, these recent works arrived at a much smaller model with comparable accuracy to the original model, reducing the problem of overfitting without compromising on classification accuracy.

However, one limitation of these previous approaches is efficiency: they involve multiple elimination and fine-tuning steps to arrive at a smaller neural network. In this paper, we propose an algorithm called Layerwise Elimination of Accessory Nodes (LEAN) that performs the elimination of accessory nodes in a single shot. Accessory nodes are defined as nodes that do not influence the classification task significantly - such nodes are assigned a low importance score by a valid decoder. LEAN uses node saliency scores derived from a valid DNN decoder to eliminate accessory nodes in the network, producing a leaner, more computationally efficient, and more generalizable DNN that is derived in a single shot.

LEAN leads to a drastic reduction in parameters, including the removal of a large number of input features. While powerful and efficient, this also leads to the loss of whole groups of correlated features. In neuroimaging datasets, such correlations exist due to the *specialization* property of the brain, whereby distinct subsystems in the brain perform specialized tasks [17]. To tackle this, we introduce another algorithm called Correlation-based eLImination of InPUts (CLIP) that identifies and retains a subset of correlated input features. When these features are combined with the remaining input features from LEAN, it leads to an improvement in the DNN performance as compared to LEAN. Combining LEAN and CLIP leads to a model that provides the best of both worlds: low number of parameters with minimal accuracy drop.

Besides the goal of getting an efficient classifier with the highest accuracy and generalization ability, a relevant and important research direction in neuroscience is to find the biological markers that are associated with a brain state (disease or cognitive task), which differentiate it from another brain state. We denote identification of brain regions and connections that are associated with a particular brain state as brain decoding. Traditional methods for brain decoding include multivariate pattern analysis [18], sparse networks based feature selection [38], latent Dirichlet allocation [40], and the use of nodal features [43]. Such methods are however based on simple or linear models. On the other hand, DNN-based approaches build deep and hierarchical models and represent key patterns underlying brain activation in its nodes and connection weights [38]. Deep learning techniques make no assumptions about application-specific a priori knowledge and therefore give consistent and unbiased decoding based on neuroimaging data. Recent applications of neural networks on neuroimaging data have derived the importance of input features in the classification task to uncover biomarkers for neurological diseases [16,19] or reveal task-related brain functional modulations [27]. Identifying disease-specific biomarkers aids in building models and classifying and predicting the progression or occurrence of unknown diseases in individuals. To do so, we use the saliency scores given by DeepLIFT to determine salient brain connections and regions for Alzheimer’s disease (AD), mild cognitive impairment (MCI), attention deficit hyperactivity disorder (ADHD), major depressive disorder (MDD), and Autism Spectrum Disorder (ASD) patients using their functional MRI scans.

In sum, we have made the following novel contributions in this work:

- We proposed an efficient layerwise node elimination approach, LEAN, that removes accessory input features and hidden layer nodes from a DNN in a single iteration. We further adapted it

for brain decoding by proposing CLIP to add in a subset of correlated features, leading to a model with minimal number of connections while maintaining the accuracy.

- The proposed approach finds salient functional connections and identifies brain regions that are responsible of classification of brain disease.
- The proposed methods were applied on multiple brain disorders for classification of patients from healthy controls and identification of disease biomarkers.

2. Methods

2.1. Feedforward DNN

Given a sample (x, d) where $x = (x_i)$ denotes the input feature vector and d denotes the sample label, we trained a feedforward neural network of L layers with the first $L - 1$ layers having rectified linear unit (ReLU) activation and a softmax layer at the end of the network. Let the weights and biases of the layer l be given by W_l and b_l , respectively. The output of layer $l \neq \{0, L\}$ is given by:

$$h_l = \text{ReLU}(W_l^T h_{l-1} + b_l) \tag{1}$$

For the input layer, $h_0 = x$. For the softmax layer, the output probabilities y of the input x belonging to class k is given by:

$$P(y = k|x) = \text{softmax}(W_L^T h_{L-1} + b_L) \tag{2}$$

where $k \in \{1, \dots, K\}$ represents the class label and the output layer weight $W_L = [w_{k,L}]$ and bias $b_L = (b_{k,L})$. To learn the parameters of the network, the cross-entropy cost is defined as:

$$J(\theta) = -E_x[\log P(y = d|x, \theta)] \tag{3}$$

where E_x is the expectation taken over all scan samples x and $\theta = \{(W_l, b_l)\}_{l=1}^L$ denotes all learnable parameters of the network.

2.2. Saliency of nodes in network layers

Let f be the neural network function mapping input x to output y . A simpler *explanation model* g is found such that g is both interpretable and an approximation of the model f . Let the number of nodes in layer l be n_l . Using an appropriate reference, let us assign to each neuron i in layer l a contribution $C_{\Delta h_{i,l} \Delta h_{k,l+1}}$ to the change in the output of neuron k in layer $l + 1$:

$$\sum_{i \leq n_l} C_{\Delta h_{i,l} \Delta h_{k,l+1}} = \Delta h_{k,l+1} \tag{4}$$

when $l (\neq L - 1, L)$ and

$$\sum_{i \leq n_{L-1}} C_{\Delta h_{i,L-1} \Delta y} = \Delta y \tag{5}$$

when $l = L - 1$, where $\Delta h_{i,l}$ is the change in the activation of the i th neuron of layer l due to the input relative to the reference. $C_{\Delta h_{i,l} \Delta h_{k,l+1}}$ is computed from the Linear, Rescale and RevealCancel rule from [44].

Given the reference input \bar{x} and the original input x , we substitute $\Delta y = f(x) - f(\bar{x})$ and $g(x) = f(x)$, giving us an equation for the model g :

$$g(x) = f(\bar{x}) + \sum_{i,l} C_{\Delta h_{i,l} \Delta y} \tag{6}$$

where the contribution of nodes in each layer l to the output y is given by $C_{\Delta h_{i,l} \Delta y}$ [44]:

$$C_{\Delta h_{i,l} \Delta y} = \Delta h_{i,l} \sum_{k,l} \frac{C_{\Delta h_{i,l} \Delta h_{k,l+1}}}{\Delta h_{i,l}} \frac{C_{\Delta h_{k,l+1} \Delta y}}{\Delta h_{k,l+1}} \tag{7}$$

Contribution of nodes in all layers $l (\neq L)$ to the output y is obtained by backpropagating contributions of layers to the output. For nodes in each layer $l (\neq L)$, we compute the *salience score vector* c_l given by:

$$c_l = \left(C_{\Delta h_{il} \Delta y} \right)_{i < n_l} \quad (8)$$

where $C_{\Delta h_{il} \Delta y}$ is derived from (7). The DeepLIFT method [44] computes the layer's salience scores based on change in the output from a reference, allowing information to propagate across the network layers even when the gradient is zero. We compute the average contributions for all layers over multiple test samples to get the final salience score for each feature.

2.3. LEAN: Layerwise elimination of accessory nodes

A decoder can identify a subset of distinguishing input features and weights that are important to classify samples. We proposed a brain decoding strategy in [16], where we obtained the salience scores for input features and hidden layer nodes of DNN classifying brain scans. In the proposed scheme, a fraction μ of nodes with the lowest salience scores c_l from layers $l (\neq L)$ were removed iteratively and the pruned model is fine-tuned at each step. Although the proposed strategy improved the classification performance, the process involved multiple iterations of elimination and fine-tuning, which were time-consuming and cumbersome. Thus, we propose here an efficient strategy to find a layerwise salience score threshold, below which nodes from the layers (input and hidden) of the DNN are removed.

For the salience score vector c_l in each layer, we determine the best fit distribution by computing the log likelihood of the model fit. We do not make any assumption about the best fit distribution and used multiple distributions (viz; power-law, log-normal, exponential and stretched exponential) to find the best fit for the salience score distribution at each layer. The power-law distribution is defined by the distribution of the salience scores c_{il} for nodes $i < n_l$ in layer l [2,10]:

$$p_{\text{power}}(c_{il}) = k c_{il}^{-\alpha}$$

where α is the scaling parameter, $k = (\alpha - 1)c_{\min}^{\alpha-1}$, and c_{\min} is the minimum degree that obeys the power-law. The distribution parameters were computed using maximum likelihood estimation [10]. We computed statistical significance for the salience scores of nodes in their respective layers and removed the nodes with p -value ≥ 0.95 . The complete overview of the layerwise elimination of nodes is given in Algorithm 1.

Algorithm 1: LEAN: Layerwise elimination of accessory nodes

Input: DNN with layers $\{l\}_{l=0}^L$

Output: Reduced DNN with layers $\{l'\}_{l'=0}^L$

Train DNN with layers $\{l\}_{l=0}^L$

for each layer l in $\{l\}_{l=0}^{L-1}$ **do**

$c_l \leftarrow C_{\Delta h_{il} \Delta y}$

$l' \leftarrow \{c_{ik} | p\text{-value}(c_{ik}) \geq 0.95\}$

2.4. CLIP: Correlation-based elimination of inputs

Besides the presence of a large number of correlated features, brain functional connectomes are also known to be modular in nature [17,47] and such modularity or clusters give rise to correlated features. Different modules correspond to different sub-systems in

the brain, performing a specific function. Examples of these modules are shown in Fig. 1(a). Modules or clusters have stronger connectivity between nodes within the module and weaker connectivity with nodes outside the module. This gives rise to the intuition that the functional connectivity for some regions of interest (ROI) - especially ROIs within the same module - are correlated with each other.

Although some modules involved in cognition and decision-making vary across subjects, others involved in functions related to perception and motor control are stable across subjects [17,32]. This gives rise to inter-subject relationships among functional connectivity features. While the extent of how high correlations is difficult to determine, correlations between pairs of inter-cluster connections would likely be lower than correlations between pairs of intra-cluster connections. This is because of the sparser and variable inter-cluster connections, and the relatively denser and more stable intra-cluster connections due to similar modularizations as shown in Fig. 1(b) and (c). For example, the dense connections within the visual module [17,32], which has a low inter-subject variability, should be highly correlated across subjects.

2.4.1. Finding clusters of correlated features

While there are other dimensionality reduction approaches such as Principal Component Analysis or Independent Component Analysis that reduce the problem of multicollinearity [22], such approaches require either prior knowledge of the number of components to be used, or extensive experimentation to arrive at the optimal number. Alternatively, a clustering-based approach presented below provides a more principled solution. Fig. 2 provides a schematic of the intuition behind finding such clusters that are spread across the distribution of salience scores.

We formulate this problem in terms of finding clusters/communities of correlated features. The similarity matrix $S = \{s_{pq}\}$ is computed where s_{pq} denotes the correlation between two training samples x^p and x^q :

$$s_{pq} = \frac{\text{Cov}(x^p, x^q)}{\sigma(x^p)\sigma(x^q)} \quad (9)$$

where Cov finds the covariance between two samples and σ finds the standard deviation across features in the sample.

We perform clustering over similarity matrix S to obtain a cluster vector (m_i) where each element m_i corresponds to the cluster label assigned to each feature i from a set $\{1, 2, \dots, M\}$ of M labels. This is done by the minimization of the normalized cut cost given by:

$$\text{cut} - \text{cost}(S, M) = \frac{1}{M} \sum_{m=1}^M \left(1 - \frac{u_m^T S u_m}{u_m^T D^k u_m} \right) \quad (10)$$

where D denotes the diagonal degree matrix of S and $U = \{u_m\}_{m=1}^M$ is a set of binary matrices representing (m_i) such that m th module is given by $u_m = (u_{mi})_{i=0}^{|\mathcal{X}|}$ where $u_{mi} = 1 (m_i = m)$ and $1(\cdot)$ denotes the identity function.

The minimization of the cut-cost in (10) is performed by multi-class spectral clustering [49]. The number of clusters is found by computing the elbow point [42] (i.e. the point of maximum curvature) of the scree plot of the eigenvalues of similarity matrix S . However, since the number of features $|\mathcal{X}|$ are often large, the similarity matrix has a large size and computing its eigendecomposition becomes expensive. Therefore, we sparsified the matrix S by thresholding correlations with values less than 0.3 (which is a recognised threshold for low correlation values [20]). Thereafter, we used the Implicitly Restarted Arnoldi Method [29] to perform eigendecomposition on the sparse matrix to obtain the eigenvalues.

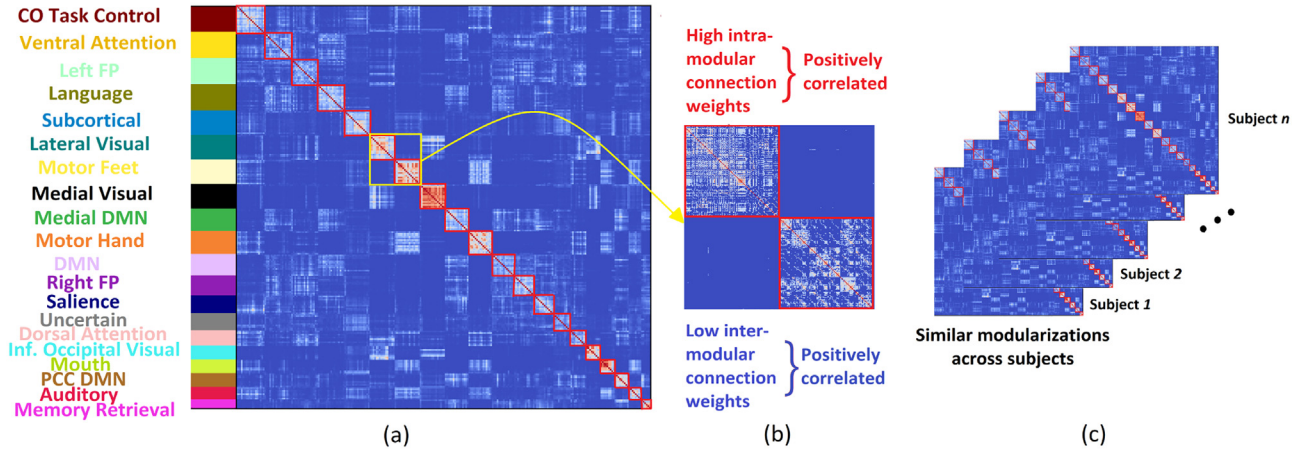


Fig. 1. Visualization of the functional connectivity matrices. (a) The matrix has its rows rearranged such that nodes in the same functional clusters are placed consecutively. The labels of different functional modules are given by the left color bar. (b) A magnified view of a subset of the correlation matrix is provided. Cells colored white have a higher weight than the cells colored blue. The white patches along the left diagonal represent intra-cluster connections, while the blue patches along the right diagonal represents inter-cluster connections. Intra-cluster connections have high and inter-cluster connections have low values across subjects and are therefore positively correlated. (c) A schematic of how some of the clusters are similar across subjects.

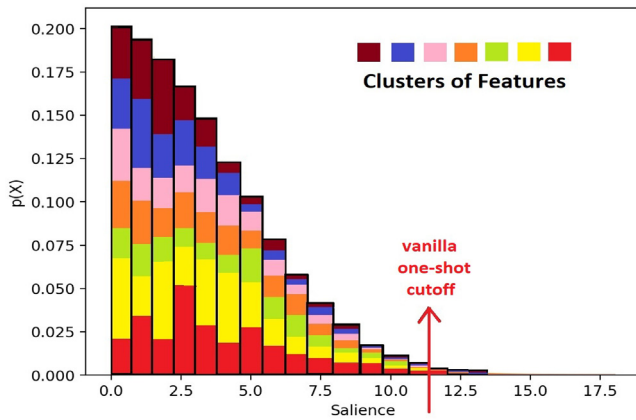


Fig. 2. Schematic of how the features are distributed in different clusters throughout the distribution of saliency scores and how saliency scores consider only features from a few clusters.

2.4.2. Selecting subsets of correlated features

For each of the computed clusters, a subset of features that have the highest correlation with features from the same cluster was selected. We do this by using the *intra-cluster degree* α_i , of feature i given by

$$\alpha_i = \sum_{m_j=m_i} S_{ij} \tag{11}$$

We select a fraction of nodes with the highest α_i values from each cluster and include them as features in addition to the ones selected from one-shot elimination pruning. The whole process for removal of subsets of correlated features is summarized in Fig. 3 and Algorithm 2.

Algorithm 2: CLIP: Correlation-based elimination of inputs

Input: Concatenated features from training data: $\{x^p\}$, k

Output: Reduced feature set l'_0

$S \leftarrow \left\{ \frac{\text{Cov}(x^p, x^q)}{\sigma(x^p)\sigma(x^q)} \right\}$

$M \leftarrow \text{elbow}(S)$

$(m_i) \leftarrow \text{minimize cut} - \text{cost}(S, M)$

$l'_0 \leftarrow l'_0 \cup \text{top } k\% \text{ features from each module in } (m_i)$

2.5. Combining LEAN and CLIP

We argue that the presence of correlated features leads to a sudden drop in accuracy with LEAN because while performing recursive elimination, we retain a fraction of features from clusters

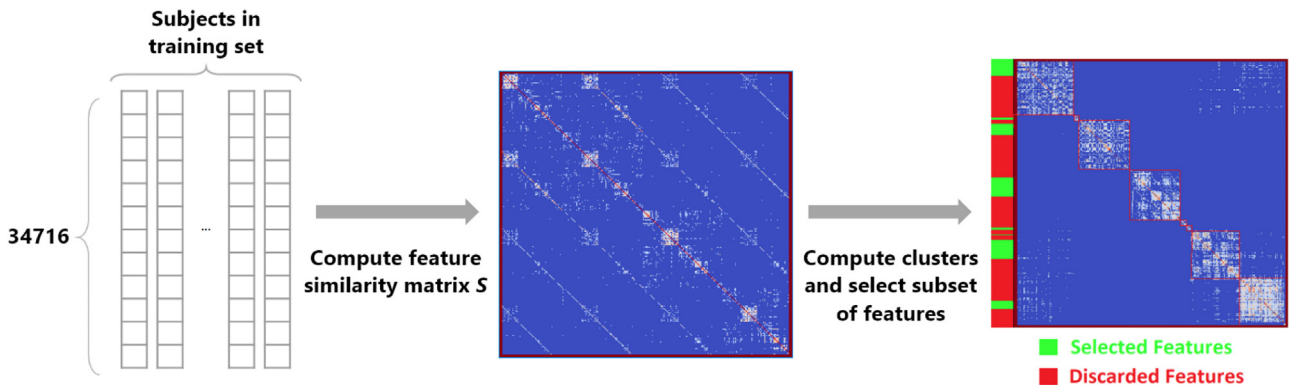


Fig. 3. Overview of the feature selection process. The vector of features representing each subject is concatenated and the similarity matrix S is computed. Thereafter, clustering is performed on S and subsets of nodes with the highest intra-cluster degree for each cluster is retained while the rest are discarded. The matrix on the right has been rearranged such that nodes in the same cluster are placed consecutively.

Table 1
Performances of the best architectures of different models on the neuroimaging datasets.

Disorder	SVM	CNN	FFN
ASD	68.0% ± 3.4%	70.6% ± 2.9%	70.9% ± 4.2%
AD	71.3% ± 6.5%	78.1% ± 5.3%	77.4% ± 4.6%
ADHD	54.8% ± 5.5%	60.5% ± 4.5%	61.1% ± 4.2%
MCI	64.6% ± 4.3%	66.5% ± 3.8%	67.5% ± 4.1%
MDD	73.3% ± 3.9%	77.1% ± 3.6%	78.3% ± 3.2%

Table 2
Comparison of functional connectivity based classifiers. Params represent percentages of the number of parameters left in the reduced model, relative to the model that uses the full feature set.

Disorder	All Features Accuracy	LEAN combined		CLIP + LEAN (Inputs Only)		CLIP + LEAN combined	
		Acc	Params	Acc	Params	Acc	Params
ASD	70.9%	68.1%	0.5%	69.3%	9.5%	68.9%	1.1%
AD	77.4%	75.8%	0.7%	78.1%	7.7%	78.1%	2.0%
ADHD	61.1%	59.0%	2.1%	59.3%	9.9%	59.4%	5.7%
MCI	67.5%	64.5%	0.6%	65.3%	6.9%	66.2%	2.0%
MDD	78.3%	75.2%	1.2%	76.4%	10.9%	76.0%	2.6%

of correlated features at different thresholds. However, with LEAN, entire clusters of correlated features are lost at once. This problem has been studied previously for various linear and non-linear classifiers [62] by using different feature selection methods [54]. For instance, in ordinary least square regression, the existence of multicollinearity leads to a larger standard error [14] and affects the interpretation of feature salience - a feature that would have been deemed as important is no longer significant when another correlated feature is present. We hypothesize that LEAN discards too many input features at once and using the additional subsets of correlated features (identified by Algorithm 2) will improve the model accuracy. Thus, while hidden layers nodes are eliminated only using LEAN (Algorithm 1), the input features are first eliminated using LEAN and then some are added back from the subsets of correlated features via CLIP (Algorithm 2).

2.6. Decoding the brain functional connectome associated with brain disease

Brain decoding is defined as the identification of brain activity patterns that emerge from activations as well as interactions among specific brain regions and connections, which distinguish one brain state from another. Decoding the DNN trained on connectome features translates to identifying salient features of the DNN, which correspond to biomarkers (i.e. key brain connections and ROI) associated with the brain state. We demonstrated this approach in [16] for the first time by using feed-forward DNN. We successfully adopted DeepLIFT in decoding brain functional connectivity in [16], which efficiently computes saliency scores for input features in a single pass and then recursively eliminate irrelevant features. Such an approach only focuses on input features and does not optimize the DNN architecture. In this paper, we improve upon our previous work by proposing a combination of LEAN and CLIP for not only to decode the input feature but also finds leaner DNN model for efficient classification without the loss of accuracy. By using LEAN and CLIP on resting-state fMRI brain scans gathered in brain disease, we achieve state-of-art accuracies for disease classification with leaner DNN models and salient input features. The decoded input features correspond to brain connections that are associated with brain disease.

3. Results

3.1. Datasets

We downloaded resting-state functional MRI scans for AD and MCI from the Alzheimer's Disease Neuroimaging Initiative (ADNI);

for ADHD from the International Neuroimaging Datasharing Initiative (INDI) [6]; for MDD from the data provided by the Creativity and Affective Neuroscience Lab in the Brain Imaging Center of Southwest University; and for ASD from the Autism Brain Imaging Data Exchange (ABIDE). The details of the acquisition protocols, subjects, and preprocessing pipelines are attached in the [supplementary materials](#).

3.2. Features for classification

We used the Power atlas [35] to obtain functionally diverse 264 ROIs spanning the entire cerebral cortex. Average time series were computed for voxels within a spherical radius of 2.5 mm surrounding each ROI, and the functional connectivity matrix for each scan was generated by computing the Pearson correlation coefficient between time-series for each pair of ROI. Since the matrices are symmetrical, we consider only the upper triangular connectivity matrix and flatten it into an input vector for the network. This resulted in an input vector with 34,716 elements.

3.3. Classification for full feature set

Both the encoder and decoder were implemented in Python using the Keras, Tensorflow and DeepLIFT libraries. For datasets with more than one scan per subject, we ensured that all the subject scans were either in the training or test set. For all datasets, a batch size of 8 was used along with a learning rate of 10^{-4} .

Besides the feedforward neural networks (FFN), convolutional neural network (CNN) architectures [8,31] and support vector machines (SVM) were implemented with different parameters. For CNNs, we gave the connectivity matrix as an input, and varied the number of filters and the number of layers. The number of filters in each layer was varied based on the number of weights in the corresponding FFN such that the number of trainable parameters were same in both architectures. The number of trainable weights were changed from 1.7×10^4 to 3.5×10^7 . For SVMs, we varied the parameters $(C \in \{0.001, 0.01, 0.1, 1, 10\}, \gamma \in \{0.001, 0.01, 0.1, 1\})$ and tried different kernels such as linear, polynomial, RBF and sigmoid. For the FFN, the number of hidden layers and the number of neurons in each hidden layer were varied. Parameters of the final architecture for the full feature set were obtained using grid search from accuracies obtained from stratified 5-fold cross validation. We added dropout of 0.1 to the hidden layers and imposed early stopping to prevent overfitting. The details of the different configurations for the DNNs are given in the [supplementary materials](#).

Table 3
The differences between test and train losses computed using cross-entropy.

Disorder	All Features	LEAN	CLIP + LEAN
ASD	0.50	0.28	0.20
AD	0.34	0.14	0.13
ADHD	0.29	0.06	0.01
MCI	0.13	0.002	0.02
MDD	0.61	0.31	0.36

Table 4
Accuracies of models produced by keeping X% most important features from SVM and logistic regression models. C + L = CLIP + LEAN.

Disorder	FFN	C + L		1%	5%	10%	100%
ASD	70.9%	68.9%	SVM	53.0%	59.0%	67.7%	68.0%
			Logreg	63.1%	67.0%	66.6%	68.1%
AD	77.4%	78.1%	SVM	66.8%	70.2%	71.6%	71.3%
			Logreg	70.9%	72.7%	72.2%	71.3%
ADHD	61.1%	59.4%	SVM	53.2%	51.6%	52.0%	53.1%
			Logreg	52.5%	51.9%	52.1%	53.2%
MCI	67.5%	66.2%	SVM	64.6%	62.9%	60.1%	58.7%
			Logreg	56.1%	57.9%	58.4%	58.7%
MDD	78.3%	76.0%	SVM	69.4%	72.5%	72.9%	73.3%
			Logreg	70.5%	72.6%	72.9%	73.3%

Table 5
Accuracies of baseline models before and after CLIP was applied.

Disorder	SVM		Logistic Regression	
	Before	After	Before	After
ASD	68.0% ± 3.3%	62.0% ± 2.9%	67.9% ± 3.1%	61.9% ± 3.1%
AD	71.3% ± 6.5%	67.6% ± 6.4%	71.2% ± 6.6%	70.7% ± 7.0%
ADHD	54.8% ± 5.5%	55.9% ± 4.8%	53.2% ± 5.7%	51.6% ± 5.4%
MCI	64.5% ± 4.3%	64.6% ± 4.3%	58.7% ± 5.4%	56.3% ± 5.5%
MDD	73.3% ± 3.9%	71.9% ± 5.0%	73.3% ± 4.3%	68.9% ± 3.7%

Table 6
Average time taken per epoch (in seconds) for each variant of the proposed algorithms. mil = million.

Disorder	Dataset Size	Initial Params	All Features	LEAN combined	CLIP + LEAN combined
ASD	823	1.74 mil	0.71	0.53	0.52
AD	299	1.74 mil	0.36	0.24	0.24
ADHD	396	1.74 mil	0.47	0.34	0.37
MCI	554	0.69 mil	0.53	0.42	0.41
MDD	457	34.78 mil	2.09	0.42	0.54

We observed that the highest accuracies were still given by FFN (except in case of CN vs AD, where the difference was insignificant). The average accuracies along with standard deviations for SVM, CNN and FFN models are reported in Table 1 and the corresponding architectures are reported in Table S5 in the supplementary materials. These accuracies are obtained from repeating the experiments using 10 different seeds. For each seed, 5-fold cross-validation was performed.

3.4. Classification with feature subsets from CLIP and LEAN

We tried three different elimination strategies: ‘LEAN’, ‘CLIP + LEAN (Inputs Only)’ and ‘LEAN + CLIP’. ‘LEAN’ (described in Algorithm 1) involves keeping only statistically significant features at each layer; ‘CLIP + LEAN (Inputs Only)’ includes subsets of correlated features (Algorithm 2) in addition to DeepLIFT features, but no elimination is performed for the hidden layers; and ‘CLIP + LEAN’ goes even further to prune the hidden layers, as described in Section 2.5. For CLIP, we generated the similarity matrix (described in Section 2.4.1) to derive a subset of correlated

features that are combined with features obtained from LEAN. We did this for each of the seeds and folds and retained 5% of the top features for each cluster. Table 2 summarises the changes in classification accuracy and remaining number of trainable parameters (compared to the original model) for different approaches.

As seen in Table 2, we observe that relative to the full feature set, LEAN is able to reduce the model parameters drastically to around 0.5% to 2.1% of the original number of parameters. Intuitively, keeping the most important features should lead to a model with minimal drop in accuracy. However, we found that the drop in accuracy is large (1.6% in ADHD to 3.1% in MDD) and thus LEAN cannot be directly used. One can argue that a different p-value threshold can be used to allow more features to be included but it would take several iterations to arrive at the threshold that gives the optimal model. Furthermore, such a tuned threshold is unlikely to generalize to other situations.

We reduced the drop in accuracy by identifying subsets of correlated features in ‘CLIP + LEAN (Inputs Only)’. Evidently, our approach leads to a smaller drop in accuracy as compared to the LEAN approach (and even lead to an increase in accuracy for CN/

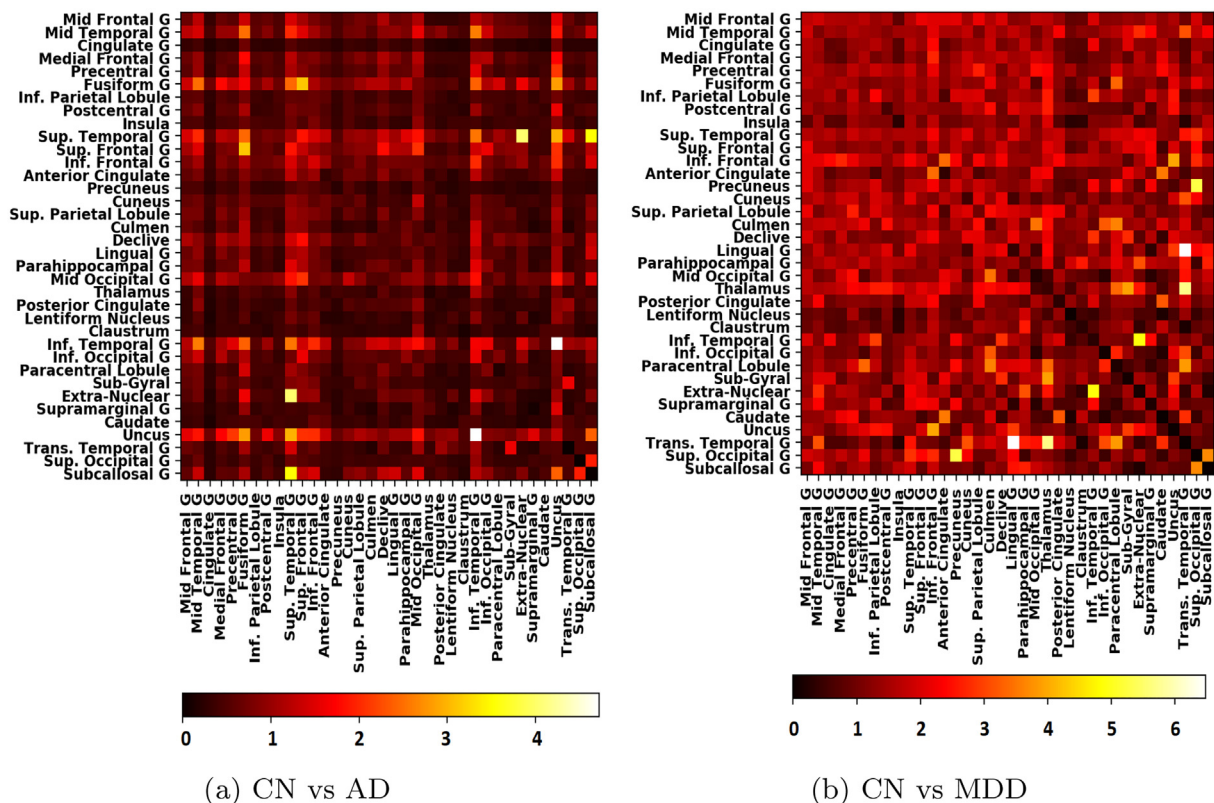


Fig. 4. The saliency scores of functional connections between brain regions derived while classifying normal and diseased participants for (a) Alzheimer’s Disease; and (b) Major Depressive Disorder. Mid: Middle, Inf: Inferior, Sup: Superior, G: Gyrus.

AD). Although we obtain similar classification performance when we only prune the inputs (and not the hidden layers), the number of remaining parameters (7% to 11% of the original) are much larger than LEAN. By extending the node elimination process to the hidden layers in ‘CLIP + LEAN’, we further reduced the number of parameters involved to levels quite close to the LEAN approach, with similar accuracies (and with consistently higher accuracies than LEAN).

Comparing the results of ‘CLIP + LEAN (Inputs Only)’ and ‘CLIP + LEAN’, there is no significant trend in terms of change in accuracy. This shows that we can safely remove hidden nodes to obtain a learner model. Thus, our subsequent analysis will be focused on LEAN and CLIP + LEAN.

3.4.1. Reduction in overfitting

Overfitting happens during model training when the test loss increases after reaching a minimum [60]. If a model overfits, it fails to generalise and remembers only the training data. To evaluate overfitting, we use the difference between the test set loss and train set loss as evaluated by cross-entropy when the best model was found. Because of the small size of the datasets, the difference of train and test losses renders a stable measure of overfitting [33].

The model was first trained with early stopping such that the best model was chosen by looking at the test accuracies of each of the epochs. The best model was chosen at the epoch where test accuracy was highest. Table 3 shows the extent of overfitting for each of the datasets with cross-entropy being the metric used to compute the loss. The reported scores were calculated by averaging over 5 folds. As shown in the table, there is a clear reduction of overfitting across all datasets for both LEAN and CLIP + LEAN, after these algorithms were used to generate the learner model.

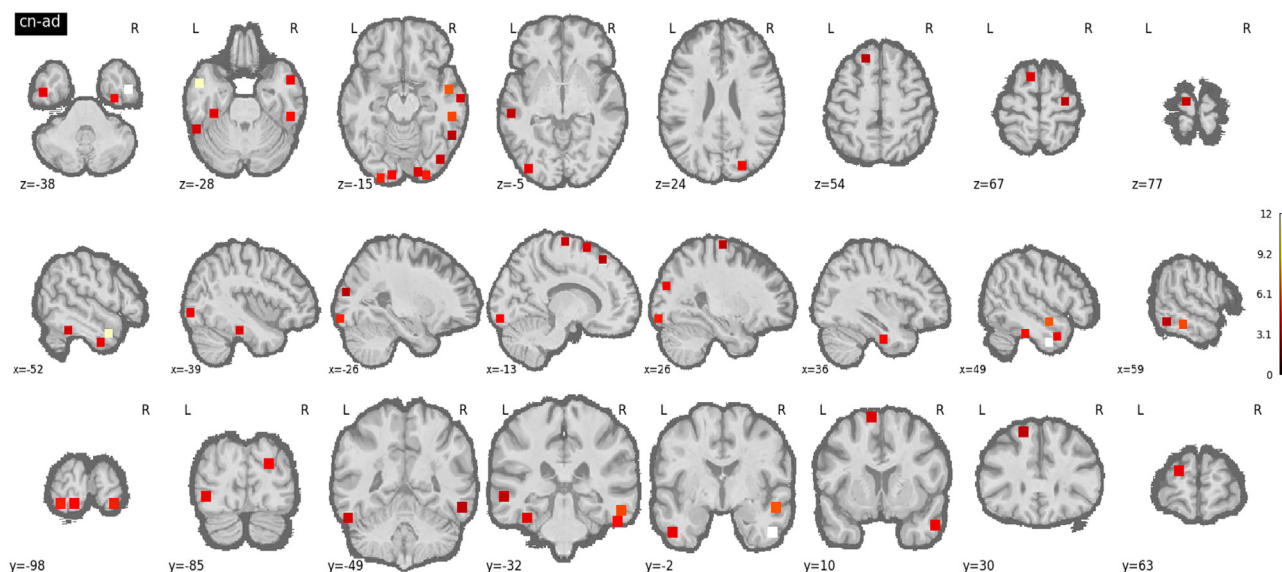
3.4.2. Comparison with other redundancy removal methods

We compared performances of CLIP + LEAN with those of SVM and logistic regression, including their own feature importance methods. Experiments were conducted, keeping 1%, 5% and 10% of features as this is the same range of remaining features obtained from our proposed algorithms. For SVM, only the linear kernel is able to provide importance scores and for this comparison, the best model with a linear kernel was chosen for each dataset. As seen in Table 4, as the number of features remaining increases, model accuracy for logistic regression generally increases but such a trend is less pronounced for SVM. Ultimately, CLIP + LEAN still generally outperforms both the SVM and logistic regression models.

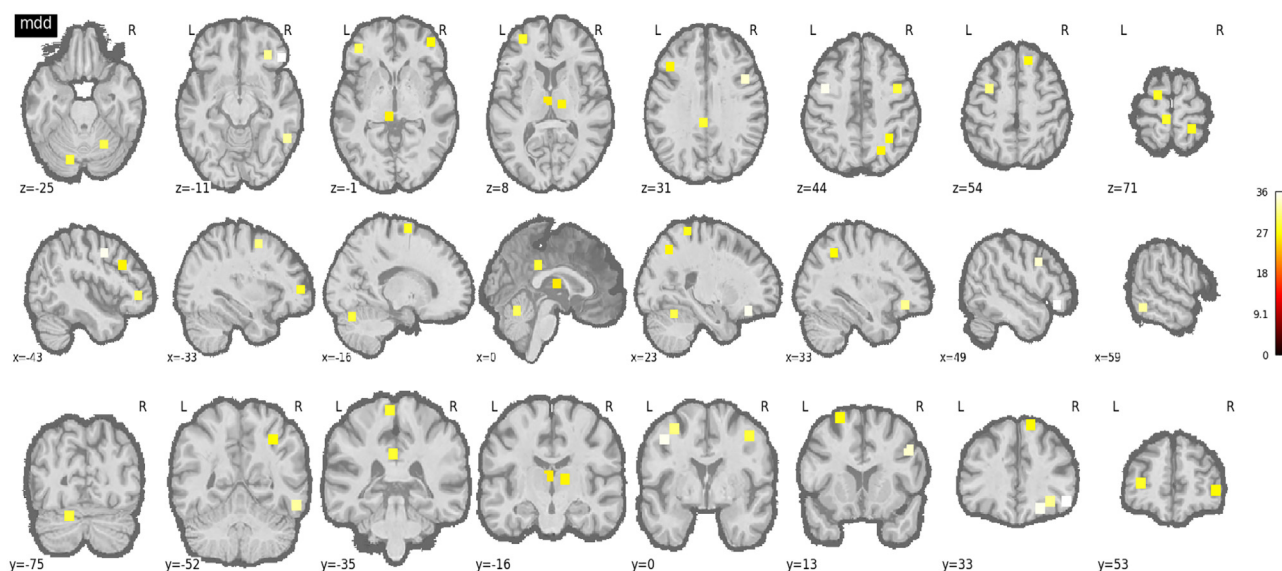
Additionally, since traditional models like SVM and logistic regression assume that the used features are uncorrelated, we used CLIP to retain a subset of correlated features and compared the performance of the resulting models. As seen in Table 5, the classification accuracy usually suffers when a feature subset is used. However, the accuracy is significantly lower ($p - value < 10^{-3}$) than the corresponding FFN in all the cases (except for MCI).

3.4.3. Comparison of algorithm runtime

In order to evaluate time efficiency of different algorithms, average time taken per epoch to train the neural network were computed. As shown in Table 6, all variants of our proposed redundancy reduction algorithms lead to a reduction of time taken - the increase being greater especially if the original ‘All Features’ model used was large like in the case of MDD. Here, we also specified the time taken for the one time operation to compute the feature clusters in CLIP. On average (over datasets and seeds), pre-processing took 135 min where the eigenvalue computation of the similarity matrix S is the source of the bottleneck. This can be significantly reduced by lowering the number of eigenvalues to be computed - in our experiments, we defined it



(a) CN vs AD



(b) CN vs MDD

Fig. 5. The axial, sagittal and coronal views of top 10% salient ROI differentiating normal and diseased patients for (a) Alzheimer’s Disease, and (b) Major Depressive Disorder.

to be 1000. Also, importance score computation took, on average, a range of approximately 20s for AD to approximately 2 min for MDD (which starts off from a large model). All experiments were performed using 4 x NVIDIA Tesla P100 16 GB on a Linux server with 72 cores. Comparing the proposed algorithms with each other, there are slight variations in the timings but the differences are largely insignificant (except for the case of MDD) and both variants take approximately the same amount of time to run. We also concluded that the pruned models, LEAN and CLIP + LEAN, take a significantly shorter amount of time to run as compared to the full model ($p - value < 10^{-4}$). When a relatively small model is used, both variants take approximately the same amount of time to run.

3.5. Potential disease biomarkers

Thus far, we have seen how salience scores help to reduce the number of parameters involved in DNN, leading to a leaner model

with minimal reduction in accuracy (or at times, even leading to an increase in model accuracy). We also use the salience scores to identify important input features that enable identification key brain ROI (or potential biomarkers) associated with different diseases. The input fed into our DNN was a 34,716-dimensional vector and each element in the vector represents a correlation score between a pair of ROI (i.e., strength of a functional connection). We analyse the importance of the functional connections by using the salience scores and the importance of ROI by computing the sum of salience scores of the functional connections incident on the ROI. Since the Power atlas [35] does not provide anatomical labels, the Crossley atlas [12] (which has anatomical labels) was used to map ROI from the Power atlas on the basis of Euclidean distance. Then, ROI with the top 10% highest scores were visualised using Nilearn.

From Fig. 4b, we found that the salient connections for MDD lie between the transverse temporal gyrus and lingual gyrus; the

transverse temporal gyrus and thalamus; superior occipital gyrus and precuneus, and between extra-nuclear and inferior temporal gyrus. For AD, we found that the connections between the superior temporal and subcallosal; superior temporal and extra nuclear, and between the uncus and inferior temporal regions were salient (Fig. 4a). For MCI, however, we found that only the connections between the uncus, subcallosal and uncus inferior temporal were salient (supplementary figure S2(b)). For ADHD, the salient connections were found between the paracentral and superior temporal; paracentral and inferior frontal; inferior temporal and superior occipital; supramarginal and sub-gyral; and supramarginal and extra nuclear regions (supplementary figure S2(a)). For ASD, salient connections were found between regions in the frontal inferior orbital to the putamen, frontal superior orbital, and frontal inferior operculum; the putamen and temporal pole superior; and cerebellum and the vermis (supplementary figure S2(c)).

For all the diseases, ROI in the inferior temporal gyrus and paracentral lobule were common among the salient regions. The other most distinctive ROI for MDD were located in the subcallosal gyrus, supplementary motor area, the medial frontal gyrus, the parahippocampal gyrus and the posterior cingulate (Fig. 5b); for ADHD in the postcentral gyrus, inferior frontal gyrus, transverse temporal gyrus, superior occipital gyrus, the insula and the medial frontal gyrus (supplementary figure S3(a)). For both AD and MCI, the salient regions were found to be in the cerebellum, temporal pole middle, frontal superior medial, parahippocampal and fusiform gyrus (Fig. 5a and supplementary figure S3(b), respectively). For ASD, salient regions were found in the inferior temporal lobe, the rectus, heschl's gyrus, cerebellum, and the paracentral lobule (supplementary figure S3(c)).

4. Discussion

4.1. DNN models for full feature set

As seen in Table 1, the feedforward DNN using the entire feature set is generally able to achieve better classification accuracies than CNNs and SVMs. The rather high standard deviation is attributable to the size of the datasets [57], which typically do not contain more than a few hundred scans. The low accuracy for ADHD is attributable to the mismatch between the age of the subjects and the age of the subjects used to derive the Power atlas [35]. Also, more subtle connectivity differences with respect to CN subjects makes it harder to classify MCI as compared to AD.

We found that CNNs and DNNs perform consistently better than SVMs. However, the difference in performance between CNNs and feedforward DNNs are subtle: except for CN-AD, feedforward DNNs performed better than CNNs but the difference is rarely more than 1%. From these results, CNNs - even the customised ones - do not seem to be getting any additional significant information for functional connectivity. Although CNNs have desirable features such as parameter sharing, they only capture information within their local receptive field (usually a small square-shaped subset, or a row or column along a matrix). Such information is not sufficient to capture the functional connectivity present in brains as these receptive fields do not consider global connectivity patterns. On the other hand, the fact that CNNs have fewer parameters than feedforward DNNs shows the possibility of redundancies in feedforward DNNs which makes it important to eliminate accessory nodes from feedforward DNNs. With LEAN, we are able to get a leaner model while methodically removing less salient features and nodes (instead of being constricted by the spatial limitations of local receptive fields).

4.2. LEAN + CLIP: high accuracy with few parameters

From the results in Table 2, we observe that LEAN has the largest accuracy drop but the smallest number of remaining parameters;

'CLIP + LEAN (Inputs Only)' results in the small drop in accuracy but retains the largest number of parameters. However, 'CLIP + LEAN' gives the best of both worlds - the number of weights retained is close to levels from LEAN, but the accuracy drop is similar to that of 'CLIP + LEAN (Inputs Only)'. An implication of our results is the presence of correlated features in functional connectomes, which can be exploited to reduce input feature set in functional/structural connectomes. In our case, we observed that LEAN led to drastic reduction in the input feature set, thereby removing sets of correlated features entirely. Our results show that adding a subset of the correlated features in 'CLIP + LEAN' (and also in 'CLIP + LEAN (Inputs Only)') led to an improvement in the performance of the classifier.

Crucially, our results show that there is a large amount of redundancy in neural network models. Our experiments with quantification of overfitting and time taken per epoch for the proposed model show that the proposed model not only leads to a reduction in overfitting but is also faster. The latter was expected since there is a significant reduction in the number of trainable network parameters. Feedforward networks have been used by researchers in the past [19,25,27] to perform classification on functional connectivity data. However, these works often use the full set of features and overfitting is a key limitation in such approaches. In our work, we have proposed a way to reduce the effects of overfitting and showed how even as it leads to a drastic reduction of parameters, the drop in accuracy is minimal. Thus, just a simple feedforward DNN is sufficient to capture crucial patterns in the data to classify healthy and diseased brains. Adding more nodes in the hidden layers does not improve model accuracy much (and in some cases such as CN-AD classification, doing so even leads to poorer generalization).

Although we have only tested our approaches on neuroimaging datasets, they are also widely applicable to other datasets. LEAN (Algorithm 1) is applicable in settings where the number of features greatly outnumber the size of the dataset, or in cases where the dataset is small and will benefit from using fewer features so as to increase generalizability of the model. CLIP (Algorithm 2) is applicable to any datasets where there are features that have strong correlations with each other.

4.3. Identified biomarkers are clinically relevant

For our algorithm to perform well, it is important that the features classified as salient by our decoder have clinical relevance to the respective condition. To verify this, we compare the results of our decoder with results from previous studies finding disease biomarkers.

For AD/MCI, previous studies have consistently pointed out to changes in the hippocampal and medial temporal lobe [9,11,52]. Accordingly, we found alterations in the inferior temporal, parahippocampal [58], fusiform [9], and paracentral lobule [52]. We observed that the most salient connections were from the uncus for both AD and MCI, an anterior extremity of the parahippocampal gyrus known for disruption during both AD [59] and MCI [65].

Multiple functional networks such as the fronto-parietal task-control network, involved in attention and emotion regulation; the default mode network, for internal medication; the dorsal attention network, for directing external attention; and the salience network, which helps in emotion processing or monitoring salient events have been implicated in MDD [23,46]. Deficits in cognitive control can be traced to the anomalies in the fronto-parietal task-control network whereas too much internal rumination (and lesser engagement with the external world) may be reflected in the aberrant connectivity of the default mode network with other goal-directed networks. The amygdala [26,36] (part of the subcallosal region), regions in the inferior temporal [56], the medial frontal lobe [56], posterior cingulate

cortex [63] and supplementary motor area [64] have been implicated in MDD. In [61], disruption between the connections between regions in the thalamus and the transverse temporal gyrus were reported, which our consistent with our results.

Multiple studies have reported anomalies in the default-mode and dorsal attention networks in children and adults suffering with ADHD suggesting altered functional connectivity with attention and cognitive processing [45,51,53,55]. The postcentral gyrus and paracentral lobule areas involved in motor functioning [45,51], medial and inferior frontal gyrus [39,45], superior occipital gyrus [55], and the insula [53] have been found to be different for children/young adolescents having ADHD, which are consistent with our results.

Likewise, the frontal orbital regions and Heschl's gyrus are involved in sensory integration, speech processing and decision-making [24,37], the putamen responsible for focusing attention [41], the cerebellar (and vermis) involved in motor regulation [3] have been previously identified as biomarkers for ASD.

5. Conclusion

In summary, we have proposed two algorithms - LEAN and CLIP - to reduce overfitting in DNNs making them more generalizable than models that use the entire feature set. Our approach leads us to an optimal neural network architecture in a single shot that is more efficient than previous methods that relied on recursive removal of features and nodes. Furthermore, via CLIP, the approach is customised to reduce the effects of correlated features that are present in neuroimaging datasets. Our experiments show that using both LEAN and CLIP took into account both redundancy and correlation in input features and gave the best balance between drop in accuracy and reduction in trainable weights. We successfully applied the proposed algorithms on 3 datasets (and 4 different neurological disorders), showing its application in brain decoding and biomarker identification. The proposed approach has applications to both structural and functional neuroimaging datasets and to investigate brain disease.

CRedit authorship contribution statement

Sukrit Gupta: Conceptualization, Methodology, Investigation, Writing - original draft. **Yi Hao Chan:** Conceptualization, Methodology, Investigation, Writing - original draft. **Jagath C. Rajapakse:** Conceptualization, Methodology, Investigation, Supervision, Writing - review & editing.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgement

This work was partially supported by AcRF Tier 1 grant RG 149/17 of Ministry of Education, Singapore. Data collection and sharing for this project was funded by the Alzheimer's Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and DOD ADNI (Department of Defense award number W81XWH-12-2-0012), the International Neuroimaging Datasharing Initiative (INDI) and the Creativity and Affective Neuroscience Lab in the Brain Imaging Center of Southwest University. We would like to thank Dr. Jiang Qiu and Dr. Dongtao Wei from the Creativity and Affective Neuroscience Lab, Brain Imaging Center of Southwest University for providing us the dataset for MDD.

Appendix A. Supplementary data

Supplementary data associated with this article can be found, in the online version, at <https://doi.org/10.1016/j.neucom.2020.04.152>.

References

- [1] Y.S. Abu-Mostafa, M. Magdon-Ismail, H.-T. Lin, *Learning from data*, Vol. 4, AMLBook New York, NY, USA, 2012.
- [2] R. Albert, A.-L. Barabási, Statistical mechanics of complex networks, *Rev. Modern Phys.* 74 (1) (2002) 47.
- [3] G. Allen, E. Courchesne, Differential effects of developmental cerebellar abnormality on cognitive and motor functions in the cerebellum: an fMRI study of autism, *Am. J. Psychiatry* 160 (2) (2003) 262–273.
- [4] S. Anwar, K. Hwang, W. Sung, Structured pruning of deep convolutional neural networks, *ACM J. Emerging Technol. Computing Syst. (JETC)* 13 (3) (2017) 32.
- [5] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, *PLoS one* 10 (7) (2015) e0130140.
- [6] P. Bellec, C. Chu, F. Chouinard-Decorte, Y. Benhajali, D.S. Margulies, R.C. Craddock, The neuro bureau ADHD-200 preprocessed repository, *Neuroimage* 144 (2017) 275–286.
- [7] R.M. Birn, The role of physiological noise in resting-state functional connectivity, *Neuroimage* 62 (2) (2012) 864–870.
- [8] Brown, C. J., Kawahara, J., Hamarneh, G., 2018. Connectome priors in deep neural networks to predict autism. In: 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018). IEEE, pp. 110–113..
- [9] D. Chan, N.C. Fox, R.I. Scallan, W.R. Crum, J.L. Whitwell, G. Leschziner, A.M. Rossor, J.M. Stevens, L. Cipolletti, M.N. Rossor, Patterns of temporal lobe atrophy in semantic dementia and Alzheimer's disease, *Ann. Neurol.* 49 (4) (2001) 433–442.
- [10] A. Clauset, C.R. Shalizi, M.E. Newman, Power-law distributions in empirical data, *SIAM Rev.* 51 (4) (2009) 661–703.
- [11] A. Convit, J. De Asis, M. De Leon, C. Tarshish, S. De Santi, H. Rusinek, Atrophy of the medial occipitotemporal, inferior, and middle temporal gyri in non-demented elderly predict decline to Alzheimer's disease, *Neurobiol. Aging* 21 (1) (2000) 19–26.
- [12] N.A. Crossley, A. Mechelli, P.E. Vértes, T.T. Winton-Brown, A.X. Patel, C.E. Ginestet, P. McGuire, E.T. Bullmore, Cognitive relevance of the community structure of the human brain functional coactivation network, *Proc. Nat. Acad. Sci.* 110 (28) (2013) 11583–11588.
- [13] G. Deco, M.L. Kringsbach, Great expectations: using whole-brain computational connectomics for understanding neuropsychiatric disorders, *Neuron* 84 (5) (2014) 892–905.
- [14] C.F. Dormann, J. Elith, S. Bacher, C. Buchmann, G. Carl, G. Carré, J.R.G. Marquéz, B. Gruber, B. Lafourcade, P.J. Leitão, et al., Collinearity: a review of methods to deal with it and a simulation study evaluating their performance, *Ecography* 36 (1) (2013) 27–46.
- [15] Y. Du, Z. Fu, V.D. Calhoun, Classification and prediction of brain disorders using functional connectivity: promising but challenging, *Front. Neurosci.* 12 (2018).
- [16] S. Gupta, Y.H. Chan, J.C. Rajapakse, Decoding brain functional connectivity implicated in AD and MCI, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2019, pp. 781–789.
- [17] S. Gupta, J.C. Rajapakse, Iterative consensus spectral clustering improves detection of subject and group level brain functional modules. *Scientific Reports (Under Review)*, 2019..
- [18] M.N. Hebart, C.I. Baker, Deconstructing multivariate decoding for the study of brain function, *Neuroimage* 180 (2018) 4–18.
- [19] A.S. Heinsfeld, A.R. Franco, R.C. Craddock, A. Buchweitz, F. Meneguzzi, Identification of autism spectrum disorder using deep learning and the ABIDE dataset, *NeuroImage: Clinical* 17 (2018) 16–23.
- [20] D.E. Hinkle, W. Wiersma, S.G. Jurs, Applied statistics for the behavioral sciences, Vol. 663, Houghton Mifflin College Division, 2003.
- [21] G. Hinton, L. Deng, D. Yu, G. Dahl, A.-R. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, B. Kingsbury, et al., Deep neural networks for acoustic modeling in speech recognition, *IEEE Signal Processing Magazine* 29 (2012).
- [22] A. Hyvärinen, E. Oja, Independent component analysis: algorithms and applications, *Neural Networks* 13 (4–5) (2000) 411–430.
- [23] R.H. Kaiser, J.R. Andrews-Hanna, J.M. Spielberg, S.L. Warren, B.P. Sutton, G.A. Miller, W. Heller, M.T. Banich, Distracted and down: neural mechanisms of affective interference in subclinical depression, *Social Cognitive Affective Neurosci.* 10 (5) (2014) 654–663.
- [24] R.K. Kana, T.A. Keller, V.L. Cherkassky, N.J. Minshew, M.A. Just, Atypical frontal-posterior synchronization of theory of mind regions in autism during mental state attribution, *Social Neurosci.* 4 (2) (2009) 135–152.
- [25] J. Kim, V.D. Calhoun, E. Shim, J.-H. Lee, Deep neural network with weight sparsity control and pre-training extracts hierarchical features and enhances classification performance: Evidence from whole-brain resting-state functional connectivity patterns of schizophrenia, *Neuroimage* 124 (2016) 127–146.
- [26] L. Kong, K. Chen, Y. Tang, F. Wu, N. Driesen, F. Womer, G. Fan, L. Ren, W. Jiang, Y. Cao, et al., Functional connectivity between the amygdala and prefrontal

- cortex in medication-naïve individuals with major depressive disorder, *J. Psychiatry Neurosci.* JPN 38 (6) (2013) 417.
- [27] Koyamada, S., Shikachu, Y., Nakae, K., Koyama, M., Ishii, S., 2015. Deep learning of fMRI big data: a novel approach to subject-transfer decoding. arXiv preprint arXiv:1502.00093..
- [28] Krizhevsky, A., Sutskever, I., Hinton, G. E., 2012. Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems*. pp. 1097–1105..
- [29] R.B. Lehoucq, D.C. Sorensen, C. Yang, ARPACK users' guide: solution of large-scale eigenvalue problems with implicitly restarted Arnoldi methods. Vol. 6. Siam, 1998..
- [30] S.M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: *Advances in Neural Information Processing Systems*, 2017, pp. 4765–4774.
- [31] R.J. Meszlényi, K. Buza, Z. Vidnyánszky, Resting state fMRI functional connectivity-based classification using a convolutional neural network architecture, *Front. Neuroinformatics* 11 (2017) 61.
- [32] S. Mueller, D. Wang, M.D. Fox, B.T. Yeo, J.E.A. Sepulcre, Individual variability in functional connectivity architecture of the human brain, *Neuron* 77 (3) (2013) 586–595.
- [33] Poggio, T., Kawaguchi, K., Liao, Q., Miranda, B., Rosasco, L., Boix, X., Hidary, J., Mhaskar, H., 2017. Theory of deep learning iii: explaining the non-overfitting puzzle. arXiv preprint arXiv:1801.00173..
- [34] J.D. Power, K.A. Barnes, A.Z. Snyder, B.L. Schlaggar, S.E. Petersen, Spurious but systematic correlations in functional connectivity mri networks arise from subject motion, *Neuroimage* 59 (3) (2012) 2142–2154.
- [35] J.D. Power, A.L. Cohen, S.M. Nelson, G.S. Wig, et al., Functional network organization of the human brain, *Neuron* 72 (4) (2011) 665–678.
- [36] R. Ramasubbu, N. Konduru, F. Cortese, S. Bray, I. Gaxiola, B. Goodyear, Reduced intrinsic connectivity of amygdala in adults with major depressive disorder, *Front. Psychiatry* 5 (2014) 17.
- [37] A. Rausch, W. Zhang, K.V. Haak, M. Mennes, E.J. Hermans, E. van Oort, G. van Wingen, C.F. Beckmann, J.K. Buitelaar, W.B. Groen, Altered functional connectivity of the amygdaloid input nuclei in adolescents and young adults with autism spectrum disorder: a resting state fmri study, *Molecular Autism* 7 (1) (2016) 13.
- [38] M.J. Rosa, L. Portugal, T. Hahn, A.J. Fallgatter, M.I. Garrido, J. Shawe-Taylor, J. Mourao-Miranda, Sparse network-based models for patient classification using fMRI, *Neuroimage* 105 (2015) 493–506.
- [39] K. Rubia, R. Halari, A. Christakou, E. Taylor, Impulsiveness as a timing disturbance: neurocognitive abnormalities in attention-deficit hyperactivity disorder during temporal processes and normalization with methylphenidate, *Phil. Trans. R. Soc. B: Biolog. Sci.* 364 (1525) (2009) 1919–1931.
- [40] T.N. Rubin, O. Koyejo, K.J. Gorgolewski, M.N. Jones, R.A. Poldrack, T. Yarkoni, Decoding brain activity using a large-scale probabilistic functional-anatomical atlas of human cognition, *PLoS Comput. Biol.* 13 (10) (2017) e1005649.
- [41] W. Sato, Y. Kubota, T. Kochiyama, S. Uono, S. Yoshimura, R. Sawada, M. Sakiyama, M. Toichi, Increased putamen volume in adults with autism spectrum disorder, *Front. Human Neurosci.* 8 (2014) 957.
- [42] V. Satopaa, J. Albrecht, D. Irwin, B. Raghavan, Finding a kneedle in a haystack: Detecting knee points in system behavior, in: *2011 31st international conference on distributed computing systems workshops, IEEE, 2011*, pp. 166–171.
- [43] B. Sen, S.-H. Chu, K.K. Parhi, Ranking regions, edges and classifying tasks in functional brain graphs by sub-graph entropy, *Sci. Rep.* 9 (1) (2019) 7628.
- [44] A. Shrikumar, P. Greenside, A. Kundaje, Learning important features through propagating activation differences, in: *Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR.org, 2017*, pp. 3145–3153..
- [45] A.B. Smith, E. Taylor, M. Brammer, R. Halari, K. Rubia, Reduced activation in right lateral prefrontal cortex and anterior cingulate gyrus in medication-naïve adolescents with attention deficit hyperactivity disorder during time discrimination, *J. Child Psychol. Psychiatry* 49 (9) (2008) 977–985.
- [46] H.R. Snyder, Major depressive disorder is associated with broad impairments on neuropsychological measures of executive function: a meta-analysis and review, *Psychol. Bull.* 139 (1) (2013) 81.
- [47] O. Sporns, R.F. Betzel, Modular brain networks, *Annu. Rev. Psychol.* 67 (2016) 613–640.
- [48] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, Dropout: a simple way to prevent neural networks from overfitting, *J. Mach. Learn. Res.* 15 (1) (2014) 1929–1958.
- [49] X.Y. Stella, J. Shi, Multiclass spectral clustering, in: *Proc. of International Conference on Computer Vision, IEEE, 2003*, pp. 313–319.
- [50] M. Sundararajan, A. Taly, Q. Yan, Axiomatic attribution for deep networks. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR.org, 2017*, pp. 3319–3328..
- [51] S.J. Suskauer, D.J. Simmonds, B.S. Caffo, M.B. Denckla, J.J. Pekar, S.H. Mostofsky, fmri of intrasubject variability in adhd: anomalous premotor activity with prefrontal compensation, *J. Am. Acad. Child Adolescent Psychiatry* 47 (10) (2008) 1141–1150.
- [52] S.N. Thiyagesh, T.F. Farrow, R.W. Parks, H. Accosta-Mesa, C. Young, I.D. Wilkinson, M.D. Hunter, P.W. Woodruff, The neural basis of visuospatial perception in Alzheimer's disease and healthy elderly comparison subjects: an fMRI study, *Psychiatry Research: Neuroimaging* 172 (2) (2009) 109–116.
- [53] L. Tian, T. Jiang, Y. Wang, Y. Zang, Y. He, M. Liang, M. Sui, Q. Cao, S. Hu, M. Peng, et al., Altered resting-state functional connectivity patterns of anterior cingulate cortex in adolescents with attention deficit hyperactivity disorder, *Neuroscience Letters* 400 (1–2) (2006) 39–43.
- [54] L. Tološi, T. Lengauer, Classification with correlated features: unreliability of feature ranking and solutions, *Bioinformatics* 27 (14) (2011) 1986–1994.
- [55] D. Tomasi, N.D. Volkow, Abnormal functional connectivity in children with attention-deficit/hyperactivity disorder, *Biological Psychiatry* 71 (5) (2012) 443–450.
- [56] J.D. Townsend, N.K. Eberhart, S.Y. Bookheimer, N.I. Eisenberger, L.C. Foland-Ross, I.A. Cook, C.A. Sugar, L.L. Altschuler, fmri activation in the amygdala and the orbitofrontal cortex in unmedicated subjects with major depressive disorder, *Psychiatry Research: Neuroimaging* 183 (3) (2010) 209–217.
- [57] G. Varoquaux, Cross-validation failure: small sample sizes lead to large error bars, *Neuroimage* 180 (2018) 68–77.
- [58] L. Wang, Y. Zang, Y. He, M. Liang, X. Zhang, L. Tian, T. Wu, T. Jiang, K. Li, Changes in hippocampal connectivity in the early stages of Alzheimer's disease: evidence from resting state fMRI, *Neuroimage* 31 (2) (2006) 496–504.
- [59] Z. Wang, M. Zhang, Y. Han, H. Song, R. Guo, K. Li, Differentially disrupted functional connectivity of the subregions of the amygdala in alzheimer's disease, *J. X-ray Sci. Technol.* 24 (2) (2016) 329–342.
- [60] A. Weigend, On overfitting and the effective number of hidden units. In: *Proceedings of the 1993 connectionist models summer school*. Vol. 1. 1994. pp. 335–342..
- [61] S.-W. Xue, D. Wang, Z. Tan, Y. Wang, Z. Lian, Y. Sun, X. Hu, X. Wang, X. Zhou, Disrupted brain entropy and functional connectivity patterns of thalamic subregions in major depressive disorder, *Neuropsychiatric Disease Treatment* 15 (2019) 2629.
- [62] K. Yan, D. Zhang, Feature selection and analysis on correlated gas sensor data with recursive feature elimination, *Sensors Actuators B: Chemical* 212 (2015) 353–363.
- [63] Z. Yao, L. Wang, Q. Lu, H. Liu, G. Teng, Regional homogeneity in depression and its relationship with separate depressive symptom clusters: a resting-state fMRI study, *J. Affective Disorders* 115 (3) (2009) 430–438.
- [64] B. Zhang, M. Li, W. Qin, L.R. Demenescu, C.D. Metzger, B. Bogerts, C. Yu, M. Walter, Altered functional connectivity density in major depressive disorder at rest, *Eur. Arch. Psychiatry Clinical Neuroscience* 266 (3) (2016) 239–248.
- [65] Zhao, Z., Lu, J., Jia, X., Chao, W., Han, Y., Jia, J., Li, K., 2014. Selective changes of resting-state brain oscillations in aMCI: an fMRI study using ALFF. *BioMed Research International* 2014..



Dr. Sukrit Gupta obtained a Bachelor in Engineering (Computer Science) from Punjab Engineering College, India and a PhD in Computer Science from NTU, Singapore. He is currently a Scientist at the Deep Learning for Healthcare Division at the Institute of Infocomm Research, A*STAR, Singapore. He works in the areas of deep learning, complex network analysis and neuroimaging. His work has mainly been focused on the detection of brain functional modules and using brain topological features as disease biomarkers, thus obtaining more efficient deep neural networks with better diagnostic accuracy.



Mr. Yi Hao Chan obtained a Bachelor in Engineering in Computer Science from NTU, Singapore. He is presently a PhD Student in Computer Science at NTU, Singapore. He works in the areas of deep learning with functional MRI and Diffusion Tensor Imaging data. His work so far has mainly focused on detecting functional and structural disease biomarkers to obtain smaller and more efficient disease classification models.



Jagath C. Rajapakse is a Professor of Computer Engineering at NTU, Singapore. His research interests are in deep learning, brain imaging, and biological networks, and has published over 300 peer reviewed articles in these areas. He was a Visiting Professor to Massachusetts Institute of Technology (MIT), USA, and Visiting Scientist to the Max-Planck Institute of Cognitive and Brain Sciences, Germany and the National Institute of Health, USA. He is a Fellow of IEEE and has served as Associate Editor for IEEE Transactions on Medical Imaging, IEEE Transactions on Neural Networks and Learning Systems, and IEEE Transactions on Computational Biology and Bioinformatics.