

BrainGNN: Interpretable Brain Graph Neural Network for fMRI Analysis*

Xiaoxiao Li^{1†}, Yuan Zhou⁴, Siyuan Gao^{1‡}, Nicha Dvornek^{1,4‡}, Muhan Zhang^{6‡}, Juntang Zhuang¹, Shi Gu⁵, Dustin Scheinost⁴, Lawrence Staib^{1,3,4}, Pamela Ventola², and James Duncan^{1,3,4}

¹ Biomedical Engineering, Yale University, New Haven, CT, 06511, USA

² Child Study Center, Yale School of Medicine, New Haven, CT, 06511, USA

³ Electrical Engineering, Yale University, New Haven, CT, 06511, USA

⁴ Radiology & Biomedical Imaging, Yale School of Medicine, New Haven, CT, 06511, USA

⁵ Department of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu, China

⁶ Facebook AI

Abstract. Understanding how certain brain regions relate to a specific neurological disorder or cognitive stimuli has been an important area of neuroimaging research. We propose BrainGNN, a graph neural network (GNN) framework to analyze functional magnetic resonance images (fMRI) and discover neurological biomarkers. In contrast to feedforward neural networks (FNN) and convolutional neural networks (CNN) in traditional functional connectivity-based fMRI analysis methods, we construct weighted graphs from fMRI and apply a GNN to fMRI brain graphs. Considering the special property of brain graphs, we design novel brain ROI-aware graph convolutional layers (Ra-GNN) that leverages the topological and functional information of fMRI. Motivated by the need for transparency in medical image analysis, our BrainGNN contains ROI-selection pooling layers (R-pool) that highlight salient ROIs (nodes in the graph), so that we can infer which ROIs are important for prediction. Furthermore, we propose regularization terms - unit loss, topK pooling (TPK) loss and group-level consistency (GLC) loss - on pooling results to encourage reasonable ROI-selection and provide flexibility to preserve either individual- or group-level patterns. We apply the BrainGNN framework on two independent fMRI datasets: Autism Spectral Disorder (ASD) fMRI dataset and Human Connectome Project (HCP) 900 Subject Release. We investigate different choices of the hyperparameters and show that BrainGNN outperforms the alternative FNN, CNN and GNN-based fMRI image analysis methods in terms of classification accuracy. The obtained community clustering and salient ROI detection results show high correspondence with the previous neuroimaging-derived evidence of biomarkers for ASD and specific task states decoded in task-fMRI.

*In submission

†Corresponding Author: Xiaoxiao Li, xiaoxiao.li@aya.yale.edu

‡Equal contribution

Keywords: Graph Neural Network, Neuroimaging, Interpretability

1 Introduction

2 The brain is an exceptionally complex system and understanding its functional
3 organization is the goal of modern neuroscience. Using fMRI, large strides in
4 understanding this organization have been made by modeling the brain as a
5 graph—a mathematical construct describing the connections or interactions (i.e.
6 edges) between different discrete objects (i.e. nodes). To create these graphs,
7 nodes are defined as brain regions of interest (ROIs) and edges are defined as the
8 functional connectivity between those ROIs, computed as the pairwise correlations
9 of functional magnetic resonance imaging (fMRI) time series, as illustrated
10 in Fig. 1. Traditional graph-based analyses for fMRI have focused on using graph
11 theoretical metrics to summarize the functional connectivity for each node into
12 a single number [46,23]. However, these methods do not consider higher-order
13 interactions between ROIs, as these interactions cannot be preserved in a single
14 number. Additionally, due to the high dimensionality of fMRI data, usually
15 ROIs are clustered into highly connected communities to reduce dimensionality.
16 Then, features are extracted from these smaller communities for further analysis
17 [32,12]. For these two-stage methods, if the results from the first stage are not
18 reliable, significant errors can be induced in the second stage.

19 The past few years have seen the growing prevalence of the use of graph
20 neural networks (GNN) for end-to-end graph learning applications. GNNs are
21 the state-of-the-art deep learning methods for most graph-structured data anal-
22 ysis problems. They combine node features, edge features, and graph structure
23 by using a neural network to embed node information and pass information
24 through edges in the graph. As such, they can be viewed as a generalization of
25 the traditional convolutional neural networks (CNN) for images. Due to their
26 high performance and interpretability, GNNs have been a widely applied graph
27 analysis method. [26,25,50,28,51]. Most existing GNNs are built on graphs that
28 do not have correspondence between the nodes of different instances, such as
29 social networks and protein networks, limiting interpretability. These methods
30 – including the current GNN methods for fMRI analysis – use the same kernel
31 over different nodes, which implicitly assumes brain graphs are translation
32 invariant. However, nodes in the same brain graph have distinct locations and
33 unique identities. Thus, applying the same kernel over all nodes is problematic.
34 In addition, few GNN studies have explored both individual-level and group-level
35 explanations, which are critical in neuroimaging research.

36 In this work, we propose a graph neural network-based framework for map-
37 ping regional and cross-regional functional activation patterns for classification
38 tasks, such as classifying neurodisorder patients versus healthy control subjects
39 and performing cognitive task decoding. Our framework jointly learns ROI cluster-
40 ing and the downstream whole-brain fMRI analysis. This not only reduces pre-
41 conceived errors, but also learns particular clustering patterns associated with
42 the downstream tasks. Specifically, from estimated model parameters, we can

43 retrieve ROI clustering patterns. Also, our GNN design facilitates model inter-
44 pretability by regulating intermediate outputs with a novel loss term, which
45 provides the flexibility to choose between individual-level and group-level expla-
46 nations.

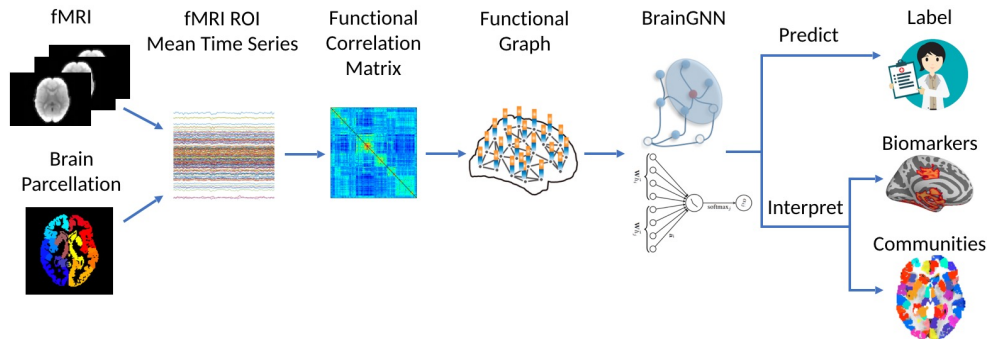


Fig. 1: The overview of the pipeline. fMRI images are parcellated by an atlas and transferred to graphs. Then, the graphs are sent to our proposed BrainGNN, which gives the prediction of specific tasks. Jointly, BrainGNN selects salient brain regions that are informative to the prediction task and clusters brain regions into prediction-related communities.

47 2 Methods and Materials

48 2.1 Preliminaries

49 Notation and Problem Definition

50 First we parcellate the brain into N regions of interest (ROIs) based on its T1
51 structural MRI. We define ROIs as graph nodes $\mathcal{V} = \{v_1, \dots, v_N\}$ and the nodes
52 are preordered. As brain ROIs can be aligned by brain parcellation atlases based
53 on their location in the structure space, we define the brain graphs as ordered
54 aligned graphs. We define an undirected weighted graph as $G = (\mathcal{V}, \mathcal{E})$, where
55 \mathcal{E} is the edge set, i.e., a collection of (v_i, v_j) linking vertices from v_i to v_j . In
56 our setting, G has an associated node feature set $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_N\}$, where \mathbf{h}_i
57 is the feature vector associated with node v_i . For every edge connecting two nodes,
58 $(v_i, v_j) \in \mathcal{E}$, we have its strength $e_{ij} \in \mathbb{R}$ and $e_{ij} > 0$. We also define $e_{ij} = 0$
59 for $(v_i, v_j) \notin \mathcal{E}$ and therefore the adjacency matrix $E = [e_{ij}] \in \mathbb{R}^{N \times N}$ is well
60 defined.

61 Architecture Overview

62 Classification on graphs is achieved by first embedding node features into a low-
63 dimensional space, then grouping nodes and summarizing them. The summarized

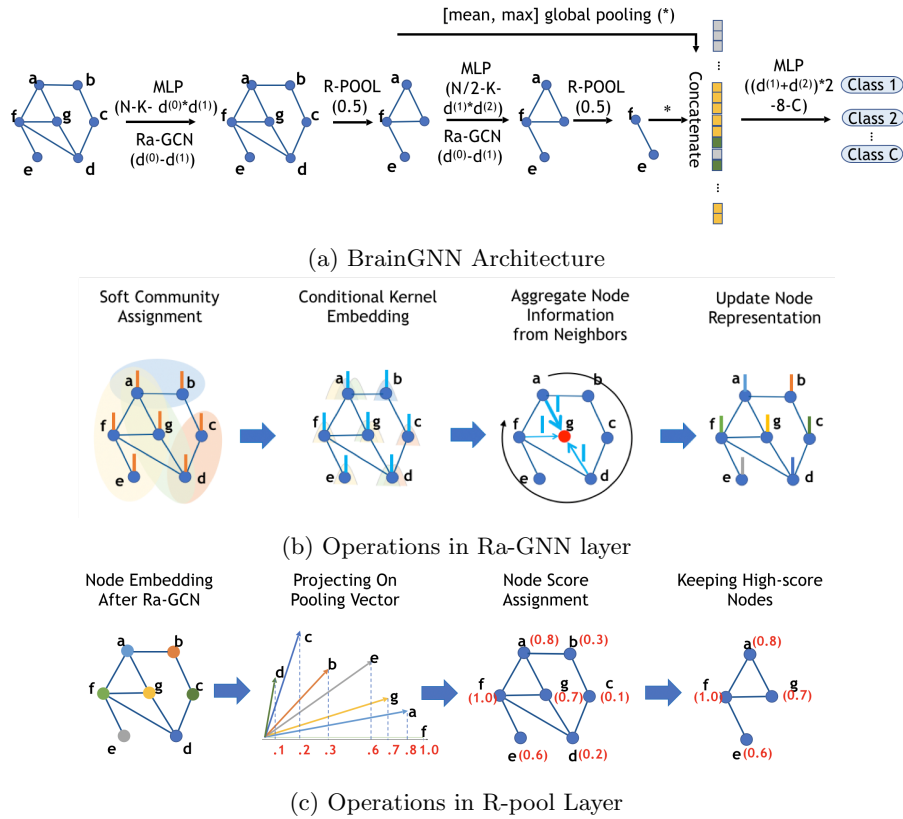


Fig. 2: (a) introduces the BrainGNN architecture that we propose in this work. BrainGNN is composed of Ra-GNN and R-pool blocks. It takes graphs as inputs and outputs graph-level predictions. (b) shows how the Ra-GNN layer embeds node features. First, nodes are softly assigned to communities based on their membership scores to the communities. Each community is associated with a different basis vector. Each node is embedded by the particular basis vectors based on the communities that it belongs to. Then, by aggregating a node's own embedding and its neighbors' embedding, the updated representation is assigned to each node on the graph. (c) shows how R-pool selects nodes to keep. First, all the nodes' representations are projected to a learnable vector. The nodes with large projected values are retained with their corresponding connections.

64 vector is then fed into a classifier, such as a multilayer perceptron (MLP), poten-
 65 tially in an end-to-end fashion. Our proposed network architecture is illustrated
 66 in Figure 2a. It is formed by three different types of layers: graph convolutional
 67 layers, node pooling layers and a readout layer. Generally speaking, GNNs induc-
 68 tively learn a node representation by recursively transforming and aggregating
 69 the feature vectors of its neighboring nodes.

70 A **graph convolutional layer** is used to probe the graph structure by using
71 edge features, which contain important information about graphs. For example,
72 the weights of the edges in brain fMRI graphs can represent the relationship
73 between different ROIs.

74 Following [39], we define $\mathbf{h}_i^{(l)} \in \mathbb{R}^{d^{(l)}}$ as the features for the i^{th} node in the
75 l^{th} layer, where $d^{(l)}$ is the dimension of the l^{th} layer features. The propagation
76 model for the forward-pass update of node representation is calculated as:

$$\mathbf{h}_i^{(l)} = \sigma \left(W_0^{(l-1)} \mathbf{h}_i^{(l-1)} + \sum_{j \in \mathcal{N}(i)} \phi \left(W_1^{(l-1)} \mathbf{h}_j^{(l-1)}, e_{ij} \right) \right), \quad (1)$$

77 where $\mathcal{N}(i)$ denotes the set of indices of neighboring nodes of node v_i and e_{ij}
78 denotes the features associated with the edge from v_i to v_j , W_0, W_1 denote the
79 model's parameters to be learned, and ϕ is any linear/nonlinear function that
80 can be applied on the neighboring nodes' feature embedding. σ is the activation
81 function.

82 A **node pooling** layer is used to reduce the size of the graph, either by
83 grouping the nodes together or pruning the original graph G to a subgraph G_s
84 by keeping some important nodes only. We will focus on the pruning method,
85 as it is more interpretable and can help detect biomarkers.

86 A **readout** layer is used to summarize the node feature vectors $\{\mathbf{h}_i^{(l)}\}$ into a
87 single vector \mathbf{z} which is finally fed into a classifier for graph classification.

88 2.2 Proposed Approach

89 In this section, we provide insights and highlight the innovative design aspects
90 of our proposed BrainGNN architecture.

91 ROI-aware Graph Convolutional Layer

92 **Overview** We propose an ROI-aware graph convolutional neural network (Ra-
93 GNN) with two insights. First, when computing the node embedding, we allow
94 Ra-GNN to learn different convolutional kernels conditioned on the ROI (geo-
95 metric information of the brain), instead of using the same kernel W on all the
96 nodes as it is shown in Eq. (1). Second, we include edge weights for message
97 filtering, as the magnitude of edge weights presents the connection strength be-
98 tween two ROIs. We assume more closely connected ROIs have higher impact.

99 **Design** We begin by assuming the graphs have additional regional information
100 and the nodes of the same region from different graphs have similar properties.
101 We propose to encode the regional information to the embedding kernel function
102 for the nodes. Given node i 's regional information \mathbf{r}_i , such as the node's coor-
103 dinates in a mesh graph, we propose to learn the vectorized embedding kernel
104 $\text{vec}(W_i^{(l)})$ based on \mathbf{r}_i on the l^{th} Ra-GNN:

$$\text{vec}(W_i^{(l)}) = f_{MLP^{(l)}}(\mathbf{r}_i) = \Theta_2^{(l)} \text{relu}(\Theta_1^{(l)} \mathbf{r}_i) + \mathbf{b}^{(l)}, \quad (2)$$

105 where the MLP network with parameters $\{\Theta_1^{(l)}, \Theta_2^{(l)}\}$ maps \mathbf{r}_i to a $d^{(l)} \cdot d^{(l-1)}$
 106 dimensional vector then reshapes the output to a $d^{(l)} \times d^{(l-1)}$ matrix $W_i^{(l)}$.

107 Given a brain parcellated into N ROIs, we order the ROIs in the same manner
 108 for all the brain graphs. Therefore, the nodes in the graphs of different subjects
 109 are aligned. However, the convolutional embedding should be independent of the
 110 ordering methods. Given an ROI ordering for all the graphs, we use one-hot en-
 111 coding to represent the ROI's location information, instead of using coordinates,
 112 because the nodes in the brain are aligned well. Specifically, for node v_i , its ROI
 113 representation \mathbf{r}_i is a N -dimensional vector with 1 in the i^{th} entry and 0 for the
 114 other entries. Assume that $\Theta_1^{(l)} = [\boldsymbol{\alpha}_1^{(l)}, \dots, \boldsymbol{\alpha}_{N^{(l)}}^{(l)}]$, where $N^{(l)}$ is the number of
 115 ROIs left on the l^{th} layer, $\boldsymbol{\alpha}_i^{(l)} = [\alpha_{i1}^{(l)}, \dots, \alpha_{iK^{(l)}}^{(l)}]^T \in \mathbb{R}^{K^{(l)}}$, $\forall i \in \{1, \dots, N^{(l)}\}$,
 116 where $K^{(l)}$ can be seen as the number of clustered communities for the $N^{(l)}$
 117 ROIs. Assume $\Theta_2^{(l)} = [\boldsymbol{\beta}_1^{(l)}, \dots, \boldsymbol{\beta}_{K^{(l)}}^{(l)}]$ with $\boldsymbol{\beta}_j^{(l)} \in \mathbb{R}^{d^{(l)} \cdot d^{(l-1)}}$, $\forall j \in \{1, \dots, K^{(l)}\}$.
 118 Then Eq. (2) can be rewritten as

$$\text{vec}(W_i^{(l)}) = \sum_{j=1}^{K^{(l)}} (\alpha_{ij}^{(l)})^+ \boldsymbol{\beta}_j^{(l)} + \mathbf{b}^{(l)}. \quad (3)$$

119 We can view $\{\boldsymbol{\beta}_j^{(l)} : j = 1, \dots, K^{(l)}\}$ as a basis and $(\alpha_{ij}^{(l)})^+$ as the coordinates.
 120 From another perspective, $(\alpha_{ij}^{(l)})^+$ can be seen as the non-negative assign-
 121 ment score of ROI i to community j . If we train different embedding kernels
 122 for different ROIs on l^{th} Ra-GNN, the total parameters to be learned will be
 123 $N^{(l)} d^{(l)} d^{(l-1)}$. Usually we have $K^{(l)} \ll N^{(l)}$. By Eq. (3), we can reduce the
 124 number of learnable parameters to $K^{(l)} d^{(l)} d^{(l-1)} + N^{(l)} K^{(l)}$ parameters, while
 125 still assigning a separate embedding kernel for each ROI. The ROIs in the same
 126 community will be embedded by the similar kernel so that nodes in different
 127 communities are embedded in different ways.

128 As the graph convolution operations in [17], the node features will be multi-
 129 plied by the edge weights, so that neighbors connected with stronger edges have
 130 a larger influence. The GNN layer using ROI-aware kernels and edge weights for
 131 filtering can be written as:

$$\mathbf{h}_i^{(l)} = W_i^{(l-1)} \mathbf{h}_i^{(l-1)} + \sum_{j \in \mathcal{N}(i)} \tilde{e}_{ij} W_j^{(l-1)} \mathbf{h}_j^{(l-1)}, \quad (4)$$

132 To avoid increasing the scale of output features, the edge features need to be
 133 normalized, as in GAT [43] and GNN [27]. Due to the aggregation mechanism,
 134 we normalize the weights by $\tilde{e}_{ij} = e_{ij} / \sum_{j \in \mathcal{N}(i)} e_{ij}$.

135 ROI-topK Pooling Layer

136 **Overview** To perform graph-level classification, a layer for dimensionality re-
 137 duction is needed since the number of nodes and the feature dimension per node
 138 are both large. Recent findings have shown that some ROIs are more indicative
 139 of predicting neurological disorders than the others [22,2], suggesting that they

140 should be kept in the dimensionality reduction step. Therefore the node (ROI)
 141 pooling layer (R-pool) is designed to keep the most indicative ROIs, thereby
 142 reducing dimensionality and removing *noisy* nodes.

143 **Design** To make sure that down-sampling layers behave idiomatically with re-
 144 spect to different graph sizes and structures, we adopt the approach in [6,15] for
 145 reducing graph nodes. The choice of which nodes to drop is determined based
 146 on projecting the node attributes onto a learnable vector $\mathbf{w}^{(l-1)} \in \mathbb{R}^{d^{(l-1)}}$. The
 147 nodes receiving lower scores will experience less feature retention. We denote
 148 $H^{(l-1)} = [\mathbf{h}_1^{(l-1)}, \dots, \mathbf{h}_{N^{(l-1)}}^{(l-1)}]^T$, where $N^{(l-1)}$ is the number of nodes at the
 149 $(l-1)^{th}$ layer. Fully written out, the operation of this pooling layer (computing
 150 a pooled graph, $(\mathcal{V}^{(l)}, \mathcal{E}^{(l)})$, from an input graph, $(\mathcal{V}^{(l-1)}, \mathcal{E}^{(l-1)})$), is expressed
 151 as follows:

$$\begin{aligned} \mathbf{s}^{(l-1)} &= H^{(l-1)} \mathbf{w}^{(l-1)} / \|\mathbf{w}^{(l-1)}\| \\ \tilde{\mathbf{s}}^{(l-1)} &= (\mathbf{s}^{(l-1)} - \mu(\mathbf{s}^{(l-1)})) / \sigma(\mathbf{s}^{(l-1)}) \\ \mathbf{i} &= \text{top}k(\tilde{\mathbf{s}}^{(l-1)}, k) \\ H^{(l)} &= (H^{(l-1)} \odot \text{sigmoid}(\tilde{\mathbf{s}}^{(l-1)}))_{\mathbf{i},:} \\ E^{(l)} &= E_{\mathbf{i},\mathbf{i}}^{(l-1)}. \end{aligned} \quad (5)$$

152 Here $\|\cdot\|$ is the L_2 norm, μ and σ take the input vector and output the mean
 153 and standard deviation of its elements. The notation $\text{top}k$ finds the indices
 154 corresponding to the largest k elements in score vector \mathbf{s} , \odot is (broadcasted)
 155 element-wise multiplication, and $(\cdot)_{\mathbf{i},\mathbf{j}}$ is an indexing operation which takes el-
 156 ements at row indices specified by \mathbf{i} and column indices specified by \mathbf{j} (colon
 157 denotes all indices). The pooling operation retains sparsity by requiring only a
 158 projection, a point-wise multiplication and a slicing into the original features
 159 and adjacency matrix. Different from [6], we induce the constraint $\|\mathbf{w}^{(l)}\|_2 = 1$
 160 implemented by adding an additional regularization loss $\sum_{l=1}^L (\|\mathbf{w}^{(l)}\|_2 - 1)^2$ to
 161 avoid identifiability issues. In addition, we added element-wise score normaliza-
 162 tion $\tilde{\mathbf{s}}^{(l)} = (\mathbf{s}^{(l)} - \mu(\mathbf{s}^{(l)})) / \sigma(\mathbf{s}^{(l)})$, which is important for calculating the GLC
 163 loss and TPK loss (introduced in Section 2.3).

164 Readout Layer

165 Lastly, we seek a “flattening” operation to preserve information about the input
 166 graph in a fixed-size representation. Concretely, to summarize the output graph
 167 of the l th conv-pool block, $(\mathcal{V}^{(l)}, \mathcal{E}^{(l)})$, we use

$$\mathbf{z}^{(l)} = \text{mean } \mathcal{H}^{(l)} \parallel \max \mathcal{H}^{(l)}, \quad (6)$$

168 where $\mathcal{H}^{(l)} = \{\mathbf{h}_i^{(l)} : i = 1, \dots, N^{(l)}\}$, mean and max operate elementwisely,
 169 and \parallel denotes concatenation. To retain information of a graph in a vector, we
 170 concatenate both mean and max summarization for a more informative graph-
 171 level representation. The final summary vector is obtained as the concatenation
 172 of all those summaries (i.e. $\mathbf{z} = \mathbf{z}^{(1)} \parallel \mathbf{z}^{(2)} \parallel \dots \parallel \mathbf{z}^{(L)}$) and it is submitted to a
 173 MLP for obtaining final predictions.

174 **2.3 Putting Layers Together and Loss Functions**

175 All in all, the architecture (as shown in Fig. 2a) consists of two kinds of layers.
 176 The input is the weighted graph with its node attributes constructed from fMRI.
 177 We form a two-layer GNN block starting with ROI-aware node embedding by
 178 the proposed Ra-GNN layer in Section 2.2, followed by the proposed R-pool
 179 layer in Section 2.2. The whole network sequentially concatenates these GNN
 180 blocks, and readout layers are added after each GNN block. The final summary
 181 vector concatenates all those summaries, and an MLP is applied after that to
 182 give final predictions. Now we describe the loss function for the neural network.

183 The classification loss is the cross entropy loss:

$$L_{ce} = -\frac{1}{M} \sum_{m=1}^M \sum_{c=1}^C y_{m,c} \log(\hat{y}_{m,c}), \quad (7)$$

184 where M is the number of instances, C is the number of classes, y_{mc} is the
 185 ground truth label and \hat{y}_{mc} is the model output.

186 We add several loss terms to regulate the learning process and control the
 187 interpretability. First, as we mentioned in Section 2.2, to avoid the problem of
 188 identifiability, we propose unit loss:

$$L_{unit}^{(l)} = (\|\mathbf{w}^{(l)}\|_2 - 1)^2. \quad (8)$$

189 Note that $\tilde{\mathbf{s}}^{(l)}$ in Eq. (5) is computed from the input $H^{(l)}$. Therefore, for different
 190 inputs $H^{(l)}$, the selected entries of $\tilde{\mathbf{s}}^{(l)}$ can be very different. For our application,
 191 we want to find the common patterns/biomarkers for a certain neuro-prediction
 192 task. Thus, we add regularization to force the $\tilde{\mathbf{s}}^{(l)}$ vectors to be similar for differ-
 193 ent input instances after the first pooling layer and call it group-level consistency
 194 (GLC). We do not constrain GLC for the second pooling layer because the nodes
 195 after the first pooling layer in different graphs might be different.

196 In each training batch, suppose there are M instances, which can be parti-
 197 tioned into C subsets based on the class labels, $\mathcal{I}_c = \{m : m = 1, \dots, M, y_{m,c} =$
 198 $1\}$, for $c = 1, \dots, C$. And $y_{m,c} = 1$ indicates the m^{th} instance belonging to
 199 class c . We form the scoring matrix for the instances belonging to class c as
 200 $S_c^{(1)} = [\tilde{\mathbf{s}}_i^{(1)} : i \in \mathcal{I}_c]^T \in \mathbb{R}^{M_c \times N}$, where $M_c = |\mathcal{I}_c|$. The GLC loss can be
 201 expressed as:

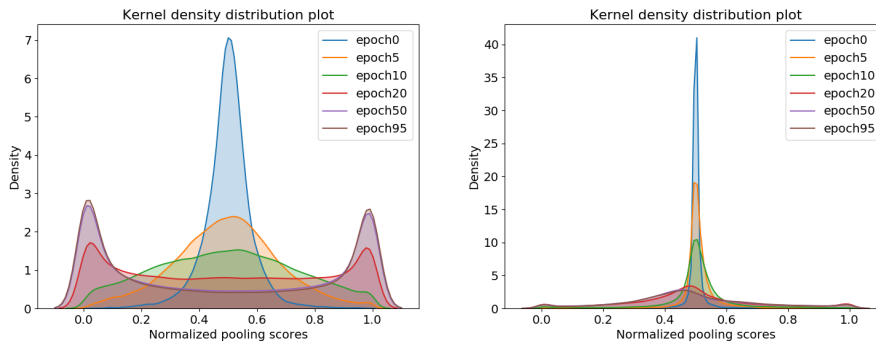
$$\begin{aligned} L_{GLC} &= \sum_{c=1}^C \sum_{i,j \in \mathcal{I}_c} \|\tilde{\mathbf{s}}_i^{(1)} - \tilde{\mathbf{s}}_j^{(1)}\|^2 \\ &= \sum_{c=1}^C \{2\text{Tr}((S_c^{(1)})^T D_c S_c^{(1)}) - 2\text{Tr}((S_c^{(1)})^T W_c S_c^{(1)})\} \\ &= 2 \sum_{c=1}^C \text{Tr}((S_c^{(1)})^T L_c S_c^{(1)}), \end{aligned} \quad (9)$$

202 where W_c is a $M_c \times M_c$ matrix with all 1s, D_c is a $M_c \times M_c$ diagonal matrix with
 203 M_c as diagonal elements, and $L_c = D_c - W_c$ is a symmetric positive semidefinite
 204 matrix [45].

205 In addition, we hope the top k selected indicative ROIs should have signifi-
 206 cantly different scores than those of the unselected nodes. Ideally, the scores for
 207 the selected nodes should be close to 1 and the scores for the unselected nodes
 208 should be close to 0. To achieve this, we rank $\text{sigmoid}(\hat{s}_m^{(l)})$ for the m th instance
 209 in a descending order, denote it as $\hat{s}_m^{(l)} = [\hat{s}_{m,1}^{(l)}, \dots, \hat{s}_{m,N^{(l)}}^{(l)}]$, and apply a con-
 210 straint to all the M training instances to make the values of $\hat{s}_m^{(l)}$ more dispersed.
 211 We define TPK loss using binary cross-entropy as:

$$L_{TPK}^{(l)} = -\frac{1}{M} \sum_{m=1}^M \frac{1}{N^{(l)}} \left(\sum_{i=1}^k \log(\hat{s}_{m,i}^{(l)}) + \sum_{i=1}^{N^{(l)}-k} \log(1 - \hat{s}_{m,i+k}^{(l)}) \right), \quad (10)$$

212 We show the kernel density estimate plots of normalized node pooling scores
 213 (indication of the importance of the nodes) changing over the training epoch
 214 in Fig. 3 when $k = \frac{1}{2}N^{(l)}$. It is clear to see that the pooling scores are more
 215 dispersed over time, Hence the top 50% selected nodes have significantly higher
 importance scores than the unselected ones.



(a) The change of the distribution of node pooling scores \hat{s} of the 1st R-pool layer over 100 training epochs. (b) The change of the distribution of node pooling scores \hat{s} of the 2nd R-pool layer over 100 training epochs.

Fig. 3: Effect of TopK pooling (TPK) loss. With the TPK loss regularization, the node pooling scores of the selected nodes and those of the unselected nodes become significantly separate.

218 Finally, the final loss function is formed as:

$$L_{total} = L_{ce} + \sum_{l=1}^L \lambda_1^{(l)} L_{unit}^{(l)} + \sum_{l=1}^L \lambda_2^{(l)} L_{TPK}^{(l)} + \lambda_3 L_{GLC}, \quad (11)$$

219 where λ 's are tunable hyper-parameters, l indicates the l^{th} GNN block and L
220 is the total number of GNN blocks. GLC loss is only calculated based on the
221 pooling scores of the first pooling layer.

222 2.4 Interpretation from BrainGNN

223 **Community Detection from Convolutional Layers** The important contri-
224 bution of our proposed ROI-aware convolutional layer is the implied community
225 clustering patterns in the graph. Discovering brain community patterns is crit-
226 ical to understanding co-activation and interaction in the brain. Revisiting Eq.
227 (3), α_{ij}^+ provides the membership of ROI i to community j . The community
228 assignment is soft and overlaid. It is similar to tensor decomposition-based com-
229 munity detection methods, such as PARAFAC [7], that decompose the tensor
230 to discover overlapping functional brain networks. Parameter α_i^+ can be seen
231 as the loading vector in PARAFAC that presents the membership of each node
232 to a certain community. Hence, we consider region i belongs to community j
233 if $\alpha_{ij} > \mu(\alpha_i^+) + \sigma(\alpha_i^+)$ [29]. This gives us a collection of community indices
234 indicating region membership $\{i_j \subset \{1, \dots, N\} : j = 1, \dots, K\}$.

235 **Biomarker Detection from Pooling Layers** Without the added TPK loss
236 (Eq. (10)), the significance of the nodes left after pooling cannot be guaranteed.
237 With TPK loss, pooling scores are more dispersed over time, hence the selected
238 nodes have significantly higher importance scores than the unselected ones. The
239 strength of the GLC loss controls the tradeoff between individual-level interpre-
240 tation and group-level interpretation. On the one hand, for precision medicine,
241 individual-level biomarkers are desired for planning targeted treatment. On the
242 other hand, group-level biomarkers are essential for understanding the common
243 characteristic patterns associated with the disease. We can tune the coefficient
244 λ_3 to control different levels of interpretation. Large λ_3 encourages selecting
245 similar nodes, while small λ_3 allows various node selection results for different
246 instances.

247 2.5 Datasets

248 Two independent datasets, the Biopoint Autism Study Dataset (Biopoint) [44]
249 and the Human Connectome Project (HCP) 900 Subject Release [42], are used
250 in this work. For the Biopoint dataset, the aim is to classify Autism Spectrum
251 Disorder (ASD) and Healthy Control (HC). For the HCP dataset, the aim is
252 to classify 7 task states - gambling, language, motor, relational, social, working
253 memory (WM), emotion.

254 **Biopoint Dataset** The Biopoint Autism Study Dataset [44] contains 72 ASD
255 children and 43 age-matched ($p > 0.124$) and IQ-matched ($p > 0.122$) neurotyp-
256 ical healthy controls (HCs). For the fMRI scans, subjects perform the "biopoint"
257 task, viewing point-light animations of coherent and scrambled biological motion
258 in a block design [22] (24s per block).

259 The fMRI data are preprocessed using FSL as follows: 1) motion correction
260 using MCFLIRT, 2) interleaved slice timing correction, 3) BET brain extraction,
261 4) spatial smoothing (FWHM=5mm), and 5) high-pass temporal filtering. The
262 functional and anatomical data are registered to the MNI152 standard brain
263 atlas [44] using FreeSurfer. The first few frames are discarded, resulting in 146
264 frames for each fMRI sequence.

265 The Desikan-Killiany [11] atlas is used to parcellate brain images into 84
266 ROIs. The mean time series for each node is extracted from a random 1/3 of
267 voxels in the ROI (given an atlas) by bootstrapping. We augment the data 30
268 times, resulting in 2160 ASD graphs and 1290 HC graphs separately. Edges are
269 defined by thresholding (top 10% positive) partial correlations to achieve sparse
270 connections. For node attributes, we concatenate seven handcrafted features: the
271 degree of node, the mean and standard deviation of the task-fMRI time series,
272 General Linear Model (GLM) coefficients, and Pearson correlation coefficient to
273 node 1 – 84. Pearson correlation and partial correlation are different measures
274 of fMRI connectivity. We aggregate them by using one to build edge connec-
275 tions and the other to build node features. For the GLM coefficients, they are
276 the coefficients of the biological motion matrix, the coefficient of the scramble
277 motion matrix, and the coefficients of the previous two matrices' derivatives in
278 the "biopoint task". Hence, node feature $\mathbf{h}_i^{(0)} \in \mathbb{R}^{(7+84)}$.

279 **HCP Dataset** For this dataset, we restrict our analyses to those individuals
280 who participated in all nine fMRI conditions (seven tasks, two rests) with full
281 length of scan, whose mean frame-to-frame displacement is less than 0.1 mm
282 and whose maximum frame-to-frame displacement is less than 0.15 mm (n=506;
283 237 males; ages 22–37). This conservative threshold for exclusion due to motion
284 is used to mitigate the substantial effects of motion on functional connectivity;
285 only left-right (LR) phase encoding run data are considered.

286 fMRI data were processed with standard methods (see [14] for more details)
287 and parcellated into 268 nodes using a whole-brain, functional atlas defined in a
288 separate sample (see [18] for more details). Task functional connectivity was cal-
289 culated based on the raw task time series: the mean time series of each node pair
290 were used to calculate the Pearson correlation and partial correlation. Matrices
291 were generated for LR phase encoding runs in the HCP data, and these matrices
292 were averaged for each condition, thus generating one 268×268 Pearson correla-
293 tion connectivity matrix and partial correlation connectivity matrix per individ-
294 ual per task condition. We define a weighted undirected graph with 268 nodes per
295 individual per task condition resulting in 3542 graphs in total. The same graph
296 construction method as for the Biopoint data was used: nodes represent parcel-
297 lated brain regions, and edges are constructed by thresholding (top 10% positive)

298 partial correlation. For node attributes, we concatenate three handcrafted fea-
299 tures: degree of node, mean and standard deviation of task-fMRI time series,
300 and Pearson correlation coefficient to node 1 – 268, as GLM parameters are not
301 useful for task state classification. Hence, node feature $\mathbf{h}_i^{(0)} \in \mathbb{R}^{(3+268)}$.

302 2.6 Implementation Details

303 We trained and tested the algorithm on Pytorch in the Python environment
304 using a NVIDIA Geforce GTX 1080Ti with 11GB GPU memory. The model
305 architecture was implemented with 2 conv layers and 2 pooling layers as shown
306 in Fig. 2a, with parameter $N = 84, K^{(1)} = K^{(2)} = 8, d^{(0)} = 91, d^{(1)} = 16, d^{(2)} =$
307 $16, C = 2$ for the Biopoint dataset and $N = 268, K^{(1)} = K^{(2)} = 8, d^{(0)} =$
308 $271, d^{(1)} = 32, d^{(2)} = 32, C = 7$ for HCP dataset. The pooling ratio is 0.5.
309 $\lambda_1^{(1)}$ and $\lambda_1^{(2)}$ were fixed to 1. The motivation of $K = 8$ comes from the eight
310 functional networks defined by Finn et al. [14].

311 We will discuss the variation of $\lambda_2^{(1)}, \lambda_2^{(2)}$ and λ_3 in Section 3.1. We randomly
312 split the data into five folds based on subjects, which means that the graphs
313 from a single subject can only appear in either the training or testing dataset.
314 Four folds were used as training data, and the left-out fold was used for testing.
315 Adam was used as the optimizer. We trained BrainGNN for 100 iterations with
316 an initial learning rate of 0.001 and annealed to half every 20 epochs. Each
317 batch contained 400 graphs for Biopoint data and 100 graphs for HCP data.
318 The weight decay parameter was 0.005.

319 3 Results

320 3.1 Ablation Study and Hyperparameter Discussion

321 Ablation studies were performed to investigate the ROI-aware graph convolu-
322 tional mechanism. We compared our proposed Ra-GNN layer with the strategy
323 of directly learning embedding kernels W . We denoted the alternative strat-
324 egy as 'GNN.' We tuned the coefficients $(\lambda_2^{(1)} - \lambda_2^{(2)} - \lambda_3)$ in the loss function
325 in Eq. (11). $\lambda_2^{(1)}$ and $\lambda_2^{(2)}$ encouraged more separable node importance scores
326 for selected and unselected nodes after pooling. λ_3 controlled the similarity of
327 the nodes selected by different instances, which could control the level of in-
328 terpretability between individual-level and group-level. Small λ_3 would result
329 in variant individual-specific patterns, while large λ_3 would force the model to
330 learn common group-level patterns. As task classification on HCP could achieve
331 consistently high accuracy over the parameter variations, we only showed the
332 results on the Biopoint dataset in Fig. 4 to better examine the effect of model
333 variations.

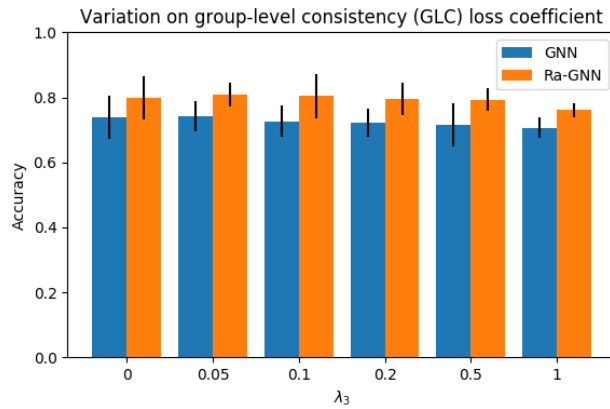
334 First, we investigated the effects of λ_3 on the accuracy, as a suitable range
335 of λ_3 should be determined in order to not sacrifice model accuracy. In Fig. 4a,
336 $\lambda_2^{(1)}$ and $\lambda_2^{(2)}$ were fixed to 0. We noticed that the results were stable to the
337 variation of λ_3 in the range 0 - 0.5. When $\lambda_3 = 1$, the accuracy dropped. The

338 accuracy reached the peak when $\lambda_3 = 0.1$. As the other deep learning models
339 behaved, BrainGNN was overparameterized. Without regularization ($\lambda_3 = 0$),
340 the model was easier to overfit to the training set, while larger regularization
341 on consistency might result in underfitting on the training set. Next, we fixed
342 $\lambda_2^{(1)} = \lambda_2^{(2)} = 0.1$ and varied λ_3 again. As the results presented in Fig. 4b, the
343 accuracy dropped if we increased λ_3 after 0.2, which followed the same trend
344 in Fig. 4a. However, the accuracy under the setting of $\lambda_3 = 0$ was better than
345 that in Fig. 4a. Probably the λ_2 terms worked as regularization and mitigated
346 the overfitting issue. Then, we fixed $\lambda_3 = 0.1$ and varied $\lambda_2^{(1)}$ and $\lambda_2^{(2)}$ from
347 0 – 0.5. As the results shown in Fig. 4c, when we increased $\lambda_2^{(1)}$ and $\lambda_2^{(2)}$ to 0.2,
348 the accuracy slightly dropped, while the accuracy sharply dropped when they
349 were increased to 0.5. For the following baseline comparison experiments, we set
350 $\lambda_2^{(1)} - \lambda_2^{(2)} - \lambda_3$ to be 0.1 – 0.1 – 0.1. As the results shown in Fig. 4, Ra-GNN
351 overall outperformed the GNN strategy under all the parameter settings. The
352 reason could be better node embedding from multiple embedding kernels in Ra-
353 GNN, as the traditional GNN strategies treated ROIs (nodes) identically and
354 used the same kernel for all the ROIs. Hence, we claim that Ra-GNN can better
355 characterize the heterogeneous representations of brain ROIs.

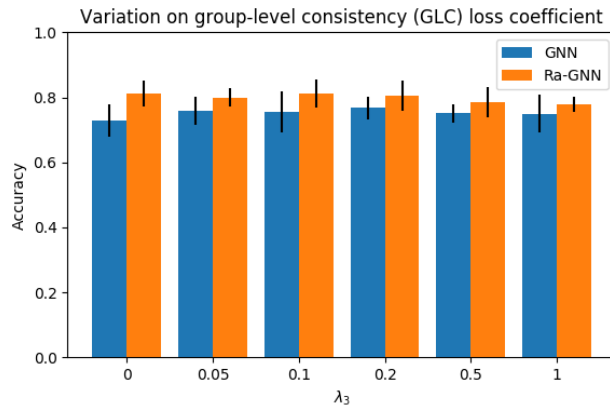
356 3.2 Comparison with Baseline Methods

357 First, we compared our method with traditional machine learning (ML) methods
358 for fMRI analysis, which took vectorized correlation matrices as inputs. The ML
359 baseline methods included Random Forest (1000 trees), SVM (RBF kernel), and
360 MLP (2 layers with 20 hidden nodes). Second, we compared our method with
361 the state-of-the-art deep learning (DL) methods, including BrainNetCNN [24],
362 and other GNN methods: 1) replace Ra-GNN layer with the graph convolutional
363 layers in Li et al. [28], 2) GraphSAGE [19] and 3) GAT (1 head) [43]. It is worth
364 noting that GraphSAGE [19] and GAT [43] did not take edge weights in the ag-
365 gregation step of the graph convolutional operation. The inputs of BrainNetCNN
366 were correlation matrices. We used the parameter settings indicated in the orig-
367 inal paper [24]. The inputs of the alternative GNN methods were the same as
368 the inputs of BrainGNN and the hyper-parameter settings for the graphconv,
369 pooling and MLP layers were the same as BrainGNN. The comparison results
370 are shown in Fig. 5

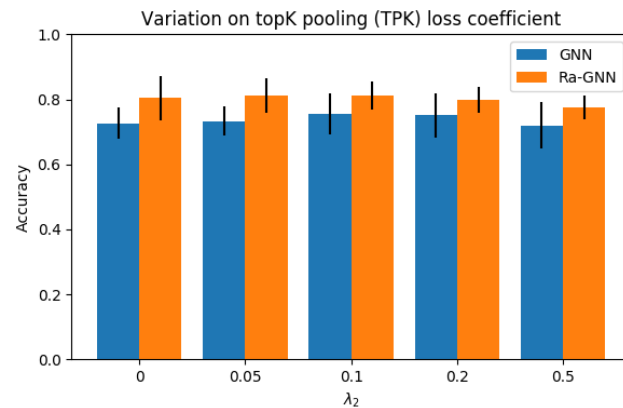
371 Our BrainGNN outperformed alternative models. The improvement may re-
372 sult from two causes. First, due to the intrinsic complexity of fMRI, complex
373 models with more parameters are desired, which also explains why CNN and
374 GNN-based methods were better than SVM and random forest. Second, our
375 model utilized the properties of fMRI and community structure in the brain net-
376 work and thus potentially modeled the local integration more effectively. Com-
377 pared to alternative machine learning models, BrainGNN achieved significantly
378 better classification results on two independent task-fMRI datasets. What is
379 more, BrainGNN does not have the burden of feature selection, which is needed
380 in traditional machine learning methods. Also, BrainGNN needs only 10 – 30%



(a) Variation on group-level consistency (GLC) loss coefficient λ_3 , when setting $\lambda_2^{(1)} = \lambda_2^{(2)} = 0$

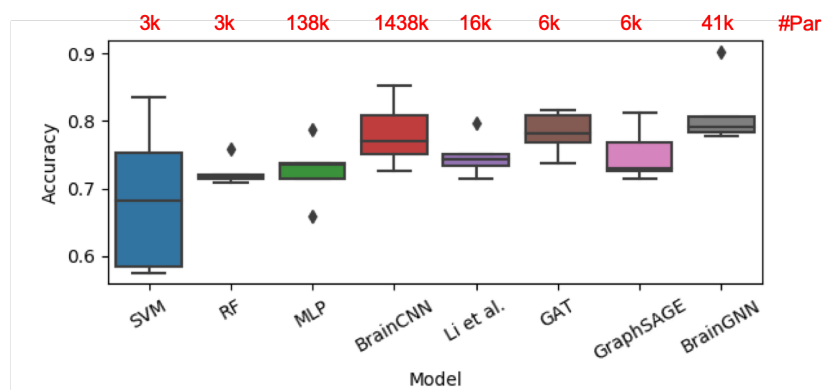


(b) Variation on group-level consistency (GLC) loss coefficient λ_3 , when setting $\lambda_2^{(1)} = \lambda_2^{(2)} = 0.1$

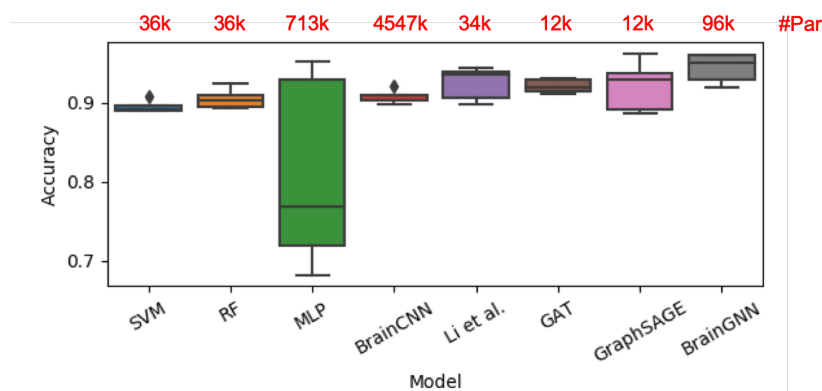


(c) Variation on topK pooling (TPK) loss coefficient $\lambda_2^{(1)}$ and $\lambda_2^{(2)}$ ($\lambda_2^{(1)} = \lambda_2^{(2)}$), when setting $\lambda_3 = 0.1$

Fig. 4: Ablation study comparison of ROI-aware GNN (Ra-GNN) and GNN without ROI embedding (GNN).



(a) Biopoint



(b) HCP

Fig. 5: Comparison of the classification accuracy of different baseline models. Classification accuracies of a 5-fold cross-validation study are depicted. The number of trainable parameters (#Par) of the deep learning models are denoted on the top of each model in red.

381 of the number of parameters compared to MLP and less than 3% of the number
382 of parameters compared to BrainNetCNN. Hence, BrainGNN is more suitable
383 as a deep learning tool for fMRI analysis, as sample size is often limited.

384 **3.3 Different Levels of Interpretation On Salient ROIs**

385 Our proposed R-pool can prune the uninformative nodes and their connections
386 from the brain graph based on the learning tasks. In other words, only the
387 salient nodes would be kept/selected. We investigated how to control the level
388 of interpretation by tuning the coefficient λ_3 that was associated with GLC loss.
389 As we discussed in Section 2.4, large λ_3 led to group-level interpretation and
390 small λ_3 led to individual-level interpretation. As we discussed in Section 3.1,
391 when λ_3 is too large, the regularization might hurt the model accuracy. We
392 put forth the hypothesis that reliable interpretation can only be guaranteed in
393 terms of a model with high classification accuracy. Hence, the interpretation was
394 restricted to models with fixed $\lambda_2^{(1)}$, $\lambda_2^{(2)}$ and varying λ_3 from 0 to 0.5 based on
395 our experiments in Section 3.1. Without losing the generalizability, we showed
396 the salient ROI detection results of 3 randomly selected ASD instances from
397 the Biopoint dataset in Fig. 6. We showed the remaining 21 ROIs after the
398 2nd R-pool layer (with pooling ratio = 0.5, 25% nodes left) and corresponding
399 pooling scores. As shown in Fig. 6(a), when $\lambda_3 = 0$, overlapped areas among
400 the three instances were rarely to be found. In Fig. 6(b-c), we circled the big
401 overlapped areas across the three instances. By visually examining the salient
402 ROIs, we found three overlapped areas in Fig. 6(b) and five overlapped areas
403 in Fig. 6(c). As proposed in Section 2.4, by tuning λ_3 , BrainGNN could achieve
404 different levels of interpretation.

405 **3.4 Validating Salient ROIs**

406 To summarize the salient ROIs over the five cross-validation folds, we averaged
407 the node pooling scores after the 1st R-pool layer for all subjects across all
408 folds per class. The top 20 salient ROIs were kept. We did not interpret the
409 model from the 2nd R-pool layer as we did in Section 3.3, because the nodes
410 left after the 1st R-pool layer may not be the same for different graphs. There-
411 fore, it was infeasible to average the pooling scores without padding 0 scores to
412 the unselected nodes. To validate the neurological significance of the result, we
413 used Neurosynth [52], a platform for fMRI data analysis. Neurosynth collects
414 thousands of neuroscience publications and provides meta-analysis, finding the
415 keywords and their associated statistical images. The decoding function on the
416 platform calculates the correlation between the input image and each functional
417 keyword's meta-analysis images.

418 In Fig. 7(a-b), we displayed the salient ROIs associated with HC and ASD
419 separately. Putamen, thalamus, temporal gyrus and insular, occipital lobe were
420 selected for HC; frontal gyrus, temporal lobe, cingulate gyrus, occipital pole, and
421 angular gyrus were selected for ASD. Hippocampus and temporal pole were im-
422 portant for both groups. The bar-chart in Fig. 7(c) illustrated the meta-analysis

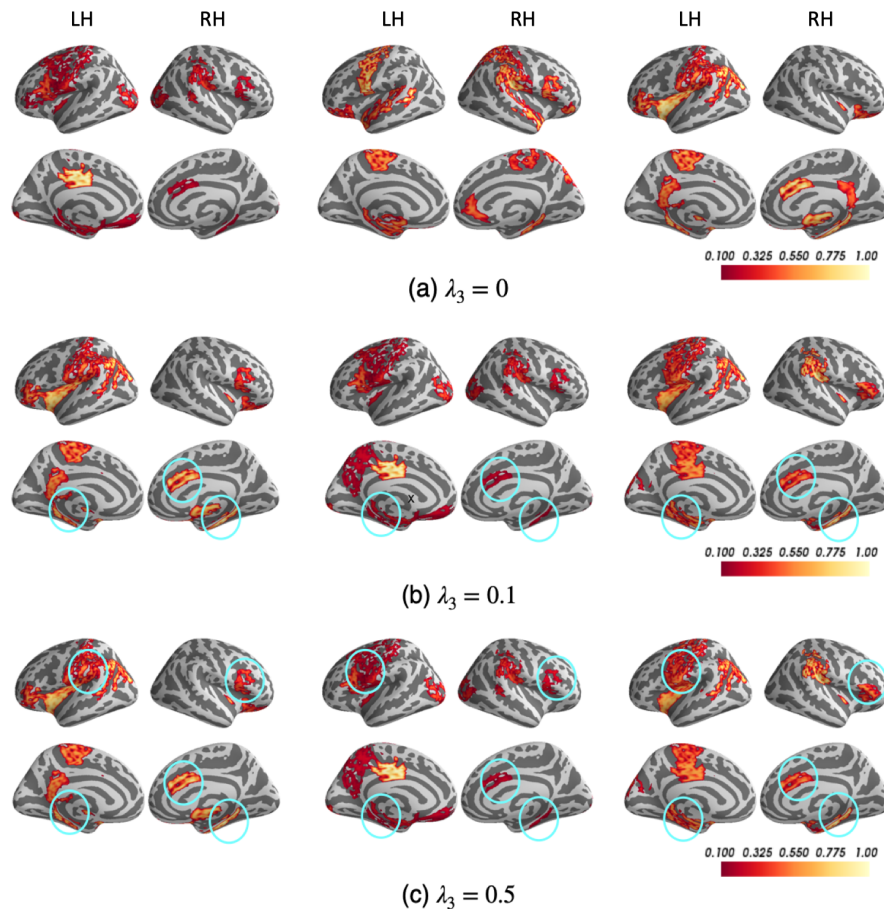
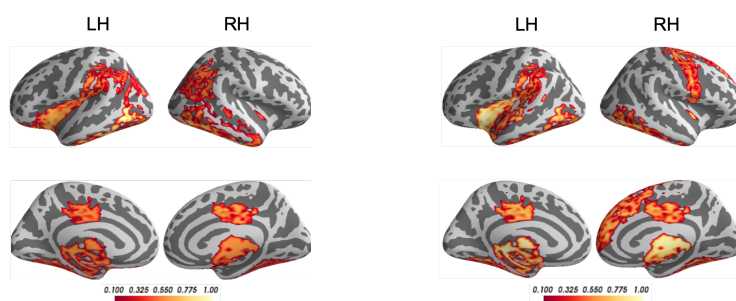


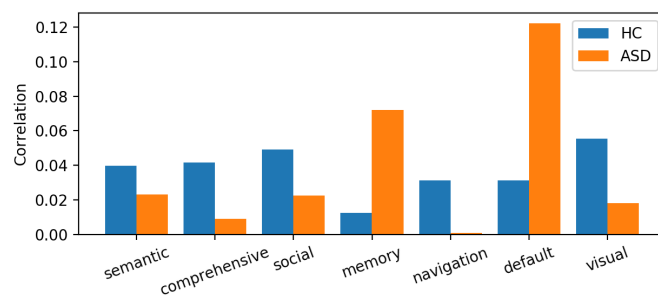
Fig. 6: The 21 selected salient ROIs of three different ASD individuals with different weights λ_3 associated with group-level consistency term L_{GLC} . The color bar ranges from 0.1 to 1. The bright-yellow color indicates a high score, while dark-red color indicates a low score. The common detected salient ROIs across different individuals are circled in blue.

423 on the functional keywords implied by the top 21 salient regions in HC and
424 ASD groups using Neurosynth. We selected ‘semantic’, ‘comprehension’, ‘so-
425 cial’, ‘memory’, ‘navigating’, ‘default’ and ‘visual’ as the functional keywords,
426 which were related to the Biopoint task [44]. We named the selected ROIs as the
427 biomarkers for identifying each group. Recall that these topics reflected unbiased
428 and aggregated findings across the fMRI literature. The functional dimensions
429 in Fig. 7(c) exposed a clear functional distinction between the two groups in task
430 fMRI decoding results. A higher value indicated a larger correlation to the func-
431 tional keywords. Specifically, the biomarkers for HC corresponded to the areas of

18 Li, X. et al



(a) Salient ROIs associated with HC. (b) Salient ROIs associated with ASD



(c) Functional keywords decoding.

Fig. 7: Interpreting salient ROIs for classifying HC vs. ASD using BrainGNN (a-b) and decoding correlation scores of ROIs associated with the functional keywords using Neurosynth (c).

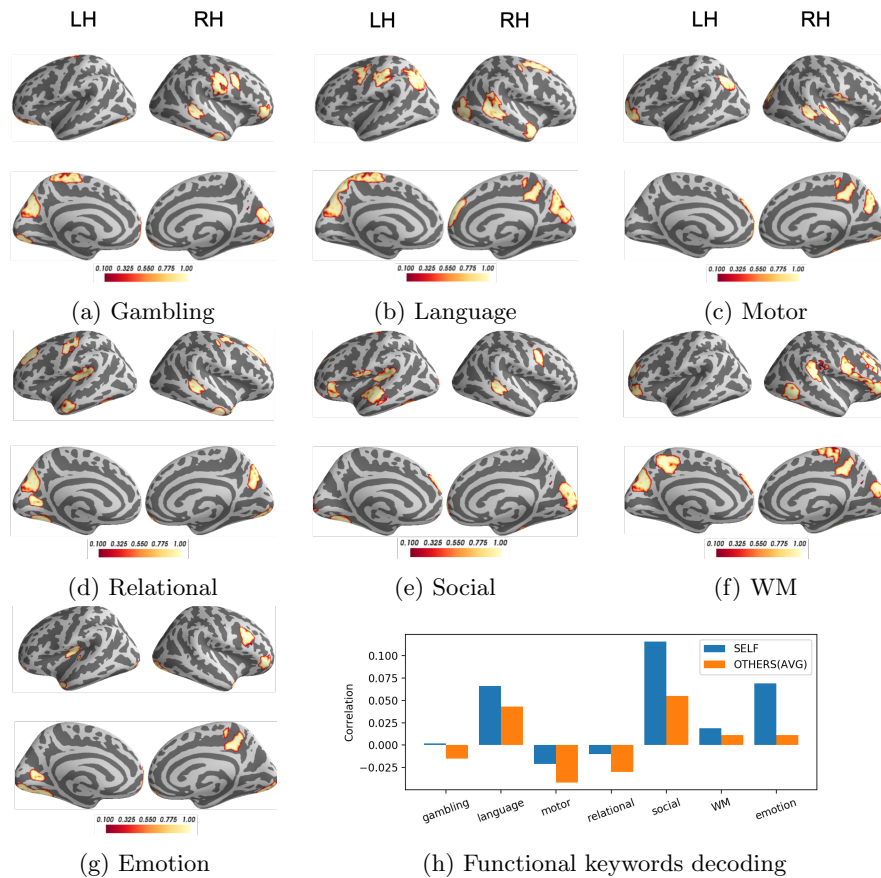


Fig. 8: Interpreting salient ROIs associated with classifying seven tasks (a-g) and decoding their correlation scores associated with the functional keywords using NeuroSynth (h). ‘SELF’ indicates the correlation score to the fMRI’s real task category, and ‘AVG(Others)’ indicates the average of the scores to the other task categories.

432 clear deficit in ASD, such as social communication, perception, and execution. In
 433 contrast, the biomarkers of ASD mapped to implicated activation-exhibited areas
 434 in ASD: default mode network [5] and memory [4]. This conclusion is consistent
 435 both with behavioral observations when administering the fMRI paradigm and
 436 with a prevailing theory that ASD includes areas of cognitive strengths amidst
 437 the social deficits [37,41,21].

438 In Fig. 8(a-g), we listed the salient ROIs associated with the seven tasks for
 439 the HCP dataset. We selected ‘gambling’, ‘language’, ‘motor’, ‘relational’, ‘so-
 440 cial’, ‘working memory’ (WM) and ‘emotion’ as the functional keywords, which

441 were exactly the functional design of the 7 tasks. The bar-chart in Fig. 8 (h) il-
442 lustrated the meta-analysis on functional keywords implied by the top 21 salient
443 regions corresponding to the seven tasks using Neurosynth. In all the seven tasks,
444 salient ROIs corresponding to each task had higher Neurosynth score than the
445 average of other tasks. The finding suggests that our algorithm identified ROIs
446 that are key to distinguish between the 7 tasks. For example, the anterior tempo-
447 ral lobe and temporal parietal regions are selected for the social task, which are
448 typically associated with social cognition [31,38]. Our findings also have overlaps
449 with the task decoding results in recent works [47].

450 3.5 Node Clustering Patterns in Ra-GNN layer

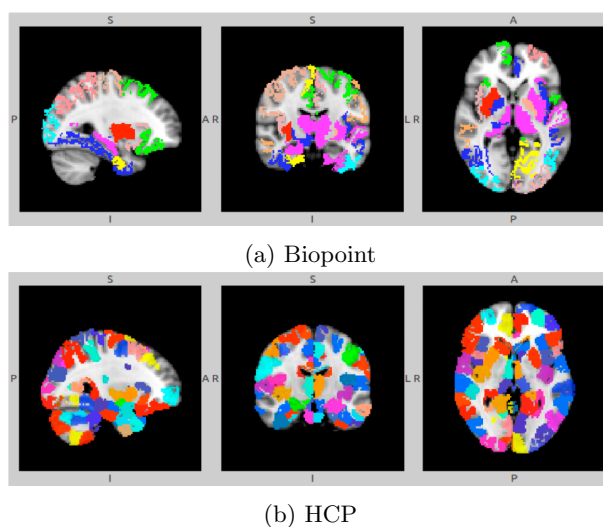


Fig. 9: ROI clustering learned from θ_1 parameter from Ra-GNN layer. Different colors denote different communities.

451 We clustered all the ROIs based on the kernel parameter α_{ij}^+ (learned in Eq.
452 (3)) of the 1st Ra-GNN layer and showed the node clustering results for Biopoint
453 and HCP data in Fig. 9a and Fig. 9b respectively. We used t-SNE [30] to visual-
454 ize the raw node features of ASD in each community in Fig. 10a and their latent
455 space embedded by the first Ra-GNN layer in Fig. 10b. The node representa-
456 tions in different communities were distinguishable, and the difference of node
457 representations within the same community were magnified, which corroborated
458 our assumption in Section 2.2 that different kernels accentuated diversified node
459 representations in different communities. Similar patterns were observed in the
460 HCP dataset.



(a) Node feature embedding for raw node representation $\mathbf{h}^{(0)}$. (b) Node feature embedding for node representation $\mathbf{h}^{(1)}$ after the 1st Ra-GNN layer.

Fig. 10: t-SNE embedding for node features of ASD group in different communities, each color indicates a community. Before Ra-GNN layer, the node features are not distinguishable, while the nodes in different communities are more distinguishable after convolution by the Ra-GNN layer. Also, the spurting pattern in (b) demonstrates that the nodes in the same community have different representation that are more distinguishable from the other communities.

461 4 Discussion

462 In this paper, we propose a graph learning model, BrainGNN, for brain network
463 analysis, that not only can perform prediction but also can be interpretable.
464 We tested the algorithms on two datasets, Biopoint and HCP, to classify brain
465 networks. Our main **contributions** are summarized as follows: 1) We formulate
466 an end-to-end framework for fMRI prediction and biomarker interpretation. The
467 methods can be generalized to general neuroimaging analysis; 2) We propose an
468 ROI-aware GNN for brain graph node (ROI) embedding, which is parameter-
469 efficient (Section 2.2) and interpretable for node clustering. Unlike other fMRI
470 analysis methods that employ clustering as a preprocessing step to reorder nodes,
471 BrainGNN learns the node grouping and extracts graph features jointly; 3) We
472 modify topK pooling [15] for informative node selection and introduce a novel
473 regularization term, topK pooling (TPK) loss (Section 2.3), to encourage more
474 reasonable node selection; 4) By regulating intermediate outputs with a novel
475 loss term, group-level consistency (GLC) loss, BrainGNN provides the flexibility
476 to choose between individual-level and group-level explanations. To the best of
477 our knowledge, we have not seen any previous research that provides flexible
478 individual-level to group-level interpretation in GNN (Section 3.3).

479 4.1 Deep Learning Methods for fMRI Prediction

480 Deep learning is a promising data-driven tool to automatically learn complex
481 feature representations in large data. Several deep learning approaches exist to

482 understand the human brain network [13,36,49]. A variety of deep methods have
483 been applied to fMRI connectome data, such as the feedforward neural net-
484 works (FNN) [34], long short-term memory (LSTM) recurrent neural networks
485 [9], and 2D convolutional neural networks (CNN) [24]. However, these exist-
486 ing deep learning methods for fMRI analysis usually require around millions of
487 parameters to learn due to the high dimensionality of fMRI connectome, thus
488 larger datasets are required to train the models. Compared with the above men-
489 tioned deep learning methods, GNNs require many fewer parameters and are
490 designed for graph-structured data analysis. Hammond et al. [20] proposed a
491 spectral graph convolution which defines convolution for graphs in the spectral
492 domain. Later, Defferrard et al. [10] simplified spectral graph convolution to a
493 local form and Kipf et al. introduced the Graph Convolutional Neural Network
494 (GCN) [27] which provided an approximated fast computation. Hamilton et al.
495 [19] proposed another variant of graph convolution in the spatial domain that
496 improves GCN’s scalability by using sampling-based neighborhood aggregation
497 and applies GCN to inductive node embedding. Different from the GNN methods
498 mentioned above, our proposed BrainGNN includes novel ROI-aware Ra-GNN
499 layers that efficiently assign each ROI an unique kernel, revealing ROI commu-
500 nity patterns and novel regulation terms (unit loss, GLC loss and TPK loss)
501 for pooling operation that regulate the model to select salient ROIs. BrainGNN
502 shows superior prediction accuracy for ASD classification and brain states de-
503 coding compared to the alternative machine learning, FCN, CNN and GNN
504 methods. As it is shown in Fig 5), BrainGNN improves average accuracy be-
505 tween 3% and 20% for ASD classification on Biopoint dataset and achieves an
506 average accuracy of 93.4% on a seven-class task states classification on HCP
507 dataset.

508 **4.2 Group-level and Individual-level Biomarker Analysis**

509 Despite the high accuracy achieved by deep learning models, a natural ques-
510 tion that arises is if the decision making process in deep learning models can
511 be interpretable. One common property of linear regression and random forest
512 is their interpretability of feature importance. For example, the coefficient of
513 linear regression model and the Gini impurity gain associated with each feature
514 in random forest can be seen as the importance scores of features. Data fea-
515 ture importance estimation is an important approach to understand both the
516 model and the underlying properties of data. From the brain biomarker detec-
517 tion perspective, understanding salient ROIs associated with the prediction is
518 an important approach to find the biomarkers, where the indicative ROIs could
519 be candidate biomarkers.

520 Although deep learning model visualization techniques have been developed
521 for convolution neural networks (CNNs), those methods are not directly appli-
522 cable to explain weighted graphs with node features for the graph classification
523 task. A few works [26,50,51] have discussed interpretable GNN models, where
524 the internal model information such as weights or structural information can be
525 accessed and inferred as group-level patterns for training instances only. Other

526 works have been used for explaining GNNs using post-hoc interpretation meth-
527 ods [35,3,53]. These post-hoc methods usually work by analyzing individual fea-
528 ture input and output pairs, which limits their explainability to individual-level
529 only. Few GNN studies have explored both individual-level and group-level ex-
530 planations, which are critical in neuroimaging research.

531 Here, we use model interpretability to address the issue of group-level and
532 individual-level biomarker analysis. In contrast, without additional post-processing
533 steps, the existing methods of fMRI analysis can only either perform individual-
534 level or group-level functional biomarker detection. For example, general linear
535 model (GLM), principal component analysis (PCA) and independent component
536 analysis (ICA) are group-based analysis methods. Some deterministic models
537 like connectome-based predictive modeling (CPM) [40,16] and other machine
538 learning based methods provide individual level-analysis. However, the model
539 flexibility for different-levels of biomarkers analysis might be required by dif-
540 ferent users. For precision medicine, individual-level biomarkers are desired for
541 planning targeted treatment, whereas group-level biomarkers are essential for
542 understanding the common characteristic patterns associated with the disease.
543 To fill the gap between group-level and individual-level biomarker analysis, we
544 introduce a tunable regularization term for our graph pooling function. By ex-
545 amining the pairs of inputs and intermediate outputs from the pooling layer,
546 our method can switch freely between individual-level and group-level explan-
547 ation under the control of the regularization term by end-to-end training. A large
548 regularization parameter encourages interpreting common biomarkers for all the
549 instances, while a small regularization parameter allows different interpretation
550 for different instances.

551 We believe that BrainGNN is the first work that uses a single framework
552 to transition between individual- and group-level analysis, filling the critical
553 interpretation gap in fMRI analysis. The biomarker interpretation results can
554 further help research in ASD and possibly generalize to rare diseases where there
555 are few patients available, as it provides individual- to group-level biomarker
556 associations.

557 4.3 BrainGNN as A Tool for Neuroimaging Analysis

558 Our proposed BrainGNN can be a research tool to identify autism biomarkers
559 using whole-brain fMRI. Our proposed method will help support efforts to better
560 understand the neural underpinnings of ASD, which is much needed in the field.
561 A more precise understanding of the neural underpinnings will guide treatment
562 approaches and help with the development of novel treatments, particularly in-
563 novative pharmacological interventions. It will also support the classification of
564 subjects in research towards more homogeneous samples, which will increase
565 power. The proposed method also provides researchers with the opportunity to
566 study neural network decisions. The challenge in applying deep models to neu-
567 roimaging research is the black box feature of this approach: no one knows what
568 the deep network is doing. Our proposed method is not only helpful for under-
569 standing the model mechanism, but also crucial for deciphering the human brain

570 network. The highly accurate results can furthermore help with classification and
571 diagnosis of neuropsychiatric diseases [33].

572 **4.4 Limitation and Future Work**

573 The pre-processing procedure performed in Section 2.5 was one possible way of
574 obtaining graphs from fMRI data, as demonstrated in this work. One meaning-
575 ful next step is to use more powerful local feature extractors to summarize ROI
576 information, such as embedding raw fMRI time series. A joint end-to-end train-
577 ing procedure that dynamically extracts graph node features from fMRI data
578 is challenging, but an interesting direction. Also, in the current work, we only
579 tried a single atlas for each dataset. In brain analysis, the reproducibility and
580 consistency of the methods are important [48,1]. For ROI-based analysis, usu-
581 ally different atlases lead to different results [8]. It is worth further investigating
582 whether the classification and interpretation results are robust to the choice of
583 the atlas. Although we discussed a few variations of hyperparameters in Sec-
584 tion 3.1, more variations should be studied, such as pooling ratio, the number
585 of communities, the number of convolutional layers, and different readout op-
586 erations. In future work, we will explore the connections between the Ra-GNN
587 layer and the tensor decomposition-based clustering methods and the patterns
588 of ROI selection and ROI clustering. For better understanding of the algorithm,
589 we aim to work on quantitative evaluations and theoretical studies to explain
590 the experimental results.

591 **5 Conclusions**

592 In this paper, we propose BrainGNN, an interpretable graph neural network for
593 fMRI analysis. BrainGNN takes graphs built from neuroimages as inputs, and
594 then outputs prediction results together with interpretation results. We applied
595 BrainGNN on Biopoint and HCP fMRI datasets. With the built-in interpretabil-
596 ity, BrainGNN not only performs better on prediction than alternative methods,
597 but also detects salient brain regions associated with predictions and discovers
598 brain community patterns. Overall, our model shows superiority over alternative
599 graph learning and machine learning classification models. By investigating the
600 selected ROIs after R-pool layers, our study reveals the salient ROIs to identify
601 autistic disorders from healthy controls and decodes the salient ROIs associated
602 with certain task stimuli. Certainly, our framework is generalizable to analysis of
603 other neuroimaging modalities. The advantages are essential for developing pre-
604 cision medicine, understanding neurological disorders, and ultimately benefiting
605 neuroimaging research.

References

1. Abraham, A., Milham, M.P., Di Martino, A., Craddock, R.C., Samaras, D., Thirion, B., Varoquaux, G.: Deriving reproducible biomarkers from multi-site resting-state data: an autism-based example. *NeuroImage* **147**, 736–745 (2017)
2. Baker, J.T., Holmes, A.J., Masters, G.A., Yeo, B.T., Krienen, F., Buckner, R.L., Öngür, D.: Disruption of cortical association networks in schizophrenia and psychotic bipolar disorder. *JAMA psychiatry* **71**(2), 109–118 (2014)
3. Baldassarre, F., Azizpour, H.: Explainability techniques for graph convolutional networks. arXiv preprint arXiv:1905.13686 (2019)
4. Boucher, J., Bowler, D.M.: Memory in autism. *Citeseer* (2008)
5. Buckner, R.L., Andrews-Hanna, J.R., Schacter, D.L.: The brain’s default network: anatomy, function, and relevance to disease. (2008)
6. Cangea, C., et al.: Towards sparse hierarchical graph classifiers. arXiv preprint arXiv:1811.01287 (2018)
7. Carroll, J.D., Chang, J.J.: Analysis of individual differences in multidimensional scaling via an n-way generalization of “eckart-young” decomposition. *Psychometrika* **35**(3), 283–319 (1970)
8. Dadi, K., Rahim, M., Abraham, A., Chyzhyk, D., Milham, M., Thirion, B., Varoquaux, G., Initiative, A.D.N., et al.: Benchmarking functional connectome-based predictive models for resting-state fmri. *Neuroimage* **192**, 115–134 (2019)
9. Dakka, J., Bashivan, P., Gheiratmand, M., Rish, I., Jha, S., Greiner, R.: Learning neural markers of schizophrenia disorder using recurrent neural networks. arXiv preprint arXiv:1712.00512 (2017)
10. Defferrard, M., Bresson, X., Vandergheynst, P.: Convolutional neural networks on graphs with fast localized spectral filtering. In: *Advances in neural information processing systems*. pp. 3844–3852 (2016)
11. Desikan, R.S., Ségonne, F., Fischl, B., Quinn, B.T., Dickerson, B.C., Blacker, D., Buckner, R.L., Dale, A.M., Maguire, R.P., Hyman, B.T., et al.: An automated labeling system for subdividing the human cerebral cortex on mri scans into gyral based regions of interest. *Neuroimage* **31**(3), 968–980 (2006)
12. Du, Y., Fu, Z., Calhoun, V.D.: Classification and prediction of brain disorders using functional connectivity: promising but challenging. *Frontiers in neuroscience* **12**, 525 (2018)
13. Eickenberg, M., Varoquaux, G., Thirion, B., Gramfort, A.: Convolutional network layers map the function of the human visual cortex. *ERCIM NEWS* (108), 12–13 (2017)
14. Finn, E.S., Shen, X., Scheinost, D., Rosenberg, M.D., Huang, J., Chun, M.M., Papademetris, X., Constable, R.T.: Functional connectome fingerprinting: identifying individuals using patterns of brain connectivity. *Nature neuroscience* **18**(11), 1664 (2015)
15. Gao, H., Ji, S.: Graph u-nets. arXiv preprint arXiv:1905.05178 (2019)
16. Gao, S., Greene, A.S., Constable, R.T., Scheinost, D.: Combining multiple connectomes improves predictive modeling of phenotypic measures. *Neuroimage* **201**, 116038 (2019)
17. Gong, L., Cheng, Q.: Exploiting edge features for graph neural networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 9211–9219 (2019)
18. Greene, A.S., Gao, S., Scheinost, D., Constable, R.T.: Task-induced brain state manipulation improves prediction of individual traits. *Nature communications* **9**(1), 1–13 (2018)

19. Hamilton, W., Ying, Z., Leskovec, J.: Inductive representation learning on large graphs. In: *Advances in neural information processing systems*. pp. 1024–1034 (2017)
20. Hammond, D.K., Vandergheynst, P., Gribonval, R.: Wavelets on graphs via spectral graph theory. *Applied and Computational Harmonic Analysis* **30**(2), 129–150 (2011)
21. Iuculano, T., Rosenberg-Lee, M., Supekar, K., Lynch, C.J., Khouzam, A., Phillips, J., Uddin, L.Q., Menon, V.: Brain organization underlying superior mathematical abilities in children with autism. *Biological Psychiatry* **75**(3), 223–230 (2014)
22. Kaiser, M.D., Hudac, C.M., Shultz, S., Lee, S.M., Cheung, C., Berken, A.M., Deen, B., Pitskel, N.B., Sugrue, D.R., Voos, A.C., et al.: Neural signatures of autism. *Proceedings of the National Academy of Sciences* **107**(49), 21223–21228 (2010)
23. Karwowski, W., Vasheghani Farahani, F., Lighthall, N.: Application of graph theory for identifying connectivity patterns in human brain networks: a systematic review. *frontiers in Neuroscience* **13**, 585 (2019)
24. Kawahara, J., Brown, C.J., Miller, S.P., Booth, B.G., Chau, V., Grunau, R.E., Zwicker, J.G., Hamarneh, G.: Brainnetcn: Convolutional neural networks for brain networks; towards predicting neurodevelopment. *NeuroImage* **146**, 1038–1049 (2017)
25. Kazi, A., Shekarforoush, S., Krishna, S.A., Burwinkel, H., Vivar, G., Kortüm, K., Ahmadi, S.A., Albarqouni, S., Navab, N.: Inceptiongc: receptive field aware graph convolutional network for disease prediction. In: *International Conference on Information Processing in Medical Imaging*. pp. 73–85. Springer (2019)
26. Kim, B.H., Ye, J.C.: Understanding graph isomorphism network for brain mr functional connectivity analysis. *arXiv preprint arXiv:2001.03690* (2020)
27. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016)
28. Li, X., Dvornek, N.C., Zhou, Y., Zhuang, J., Ventola, P., Duncan, J.S.: Graph neural network for interpreting task-fmri biomarkers. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 485–493. Springer (2019)
29. Loe, C.W., Jensen, H.J.: Comparison of communities detection algorithms for multiplex. *Physica A: Statistical Mechanics and its Applications* **431**, 29–45 (2015)
30. Maaten, L.v.d., Hinton, G.: Visualizing data using t-sne. *Journal of machine learning research* **9**(Nov), 2579–2605 (2008)
31. Mar, R.A.: The neural bases of social cognition and story comprehension. *Annual review of psychology* **62**, 103–134 (2011)
32. Moğultay, H., Alkan, S., Yarman-Vural, F.T.: Classification of fmri data by using clustering. In: *2015 23rd Signal Processing and Communications Applications Conference (SIU)*. pp. 2381–2383. IEEE (2015)
33. Nickerson, L.D.: Replication of resting state-task network correspondence and novel findings on brain network activation during task fmri in the human connectome project study. *Scientific reports* **8**(1), 1–12 (2018)
34. Patel, P., Aggarwal, P., Gupta, A.: Classification of schizophrenia versus normal subjects using deep learning. In: *Proceedings of the Tenth Indian Conference on Computer Vision, Graphics and Image Processing*. p. 28. ACM (2016)
35. Pope, P.E., Kolouri, S., Rostami, M., Martin, C.E., Hoffmann, H.: Explainability methods for graph convolutional neural networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 10772–10781 (2019)

36. Rajalingham, R., Issa, E.B., Bashivan, P., Kar, K., Schmidt, K., DiCarlo, J.J.: Large-scale, high-resolution comparison of the core visual object recognition behavior of humans, monkeys, and state-of-the-art deep artificial neural networks. *Journal of Neuroscience* **38**(33), 7255–7269 (2018)
37. Robertson, C.E., Kravitz, D.J., Freyberg, J., Baron-Cohen, S., Baker, C.I.: Tunnel vision: sharper gradient of spatial attention in autism. *Journal of Neuroscience* **33**(16), 6776–6781 (2013)
38. Ross, L.A., Olson, I.R.: Social cognition and the anterior temporal lobes. *Neuroimage* **49**(4), 3452–3462 (2010)
39. Schlichtkrull, M., Kipf, T.N., Bloem, P., Van Den Berg, R., Titov, I., Welling, M.: Modeling relational data with graph convolutional networks. In: *European Semantic Web Conference*. pp. 593–607. Springer (2018)
40. Shen, X., Finn, E.S., Scheinost, D., Rosenberg, M.D., Chun, M.M., Papademetris, X., Constable, R.T.: Using connectome-based predictive modeling to predict individual behavior from brain connectivity. *nature protocols* **12**(3), 506 (2017)
41. Turkeltaub, P.E., Flowers, D.L., Verbalis, A., Miranda, M., Gareau, L., Eden, G.F.: The neural basis of hyperlexic reading: An fmri case study. *Neuron* **41**(1), 11–25 (2004)
42. Van Essen, D.C., Smith, S.M., Barch, D.M., Behrens, T.E., Yacoub, E., Ugurbil, K., Consortium, W.M.H., et al.: The wu-minn human connectome project: an overview. *Neuroimage* **80**, 62–79 (2013)
43. Veličković, P., et al.: Graph attention networks. In: *ICLR* (2018)
44. Venkataraman, A., Yang, D.Y.J., Pelphrey, K.A., Duncan, J.S.: Bayesian community detection in the space of group-level functional differences. *IEEE transactions on medical imaging* **35**(8), 1866–1882 (2016)
45. Von Luxburg, U.: A tutorial on spectral clustering. *Statistics and computing* **17**(4), 395–416 (2007)
46. Wang, J., Zuo, X., He, Y.: Graph-based network analysis of resting-state functional mri. *Frontiers in systems neuroscience* **4**, 16 (2010)
47. Wang, X., Liang, X., Jiang, Z., Nguchu, B.A., Zhou, Y., Wang, Y., Wang, H., Li, Y., Zhu, Y., Wu, F., et al.: Decoding and mapping task states of the human brain via deep learning. *Human Brain Mapping* (2019)
48. Wei, X., Warfield, S.K., Zou, K.H., Wu, Y., Li, X., Guimond, A., Mugler III, J.P., Benson, R.R., Wolfson, L., Weiner, H.L., et al.: Quantitative analysis of mri signal abnormalities of brain white matter with high reproducibility and accuracy. *Journal of Magnetic Resonance Imaging: An Official Journal of the International Society for Magnetic Resonance in Medicine* **15**(2), 203–209 (2002)
49. Yamins, D.L., DiCarlo, J.J.: Using goal-driven deep learning models to understand sensory cortex. *Nature neuroscience* **19**(3), 356 (2016)
50. Yan, Y., Zhu, J., Duda, M., Solarz, E., Sripada, C., Koutra, D.: Groupinn: Grouping-based interpretable neural network for classification of limited, noisy brain data. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. pp. 772–782 (2019)
51. Yang, H., Li, X., Wu, Y., Li, S., Lu, S., Duncan, J.S., Gee, J.C., Gu, S.: Interpretable multimodality embedding of cerebral cortex using attention graph network for identifying bipolar disorder. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 799–807. Springer (2019)
52. Yarkoni, T., Poldrack, R.A., Nichols, T.E., Van Essen, D.C., Wager, T.D.: Large-scale automated synthesis of human functional neuroimaging data. *Nature methods* **8**(8), 665 (2011)

28 Li, X. et al

53. Ying, R., Bourgeois, D., You, J., Zitnik, M., Leskovec, J.: Gnn explainer: A tool for post-hoc explanation of graph neural networks. arXiv preprint arXiv:1903.03894 (2019)