

Discovering robust biomarkers of neurological disorders from functional MRI using graph neural networks: A Review

Yi Hao Chan, Deepank Girish, Sukrit Gupta, Jing Xia, Chockalingam Kasi, Yinan He, Conghao Wang, Jagath C. Rajapakse, *Fellow, IEEE*

Abstract—Graph neural networks (GNN) have emerged as a popular tool for modelling functional magnetic resonance imaging (fMRI) datasets. Many recent studies have reported significant improvements in disorder classification performance via more sophisticated GNN designs and highlighted salient features that could be potential biomarkers of the disorder. In this review, we provide an overview of how GNN and model explainability techniques have been applied on fMRI datasets for disorder prediction tasks, with a particular emphasis on the robustness of biomarkers produced for neurodegenerative diseases (dementia, Parkinson’s disease) and neuropsychiatric disorders (attention deficit hyperactivity disorder, autism spectrum disorder, major depressive disorder, and schizophrenia). We found that while most studies have performant models, salient features highlighted in these studies vary greatly across studies on the same disorder and little has been done to evaluate their robustness. To address these issues, we suggest establishing new standards that are based on objective evaluation metrics to determine the robustness of these potential biomarkers. We further highlight gaps in the existing literature and put together a prediction-attribution-evaluation framework that could set the foundations for future research on improving the robustness of potential biomarkers discovered via GNNs.

Index Terms—Attention deficit hyperactivity disorder, biomarker discovery, dementia, graph neural networks, major depressive disorder, model explainability, robustness, Parkinson’s disease, schizophrenia

I. INTRODUCTION

NEUROLOGICAL disorders often manifest as changes in the functional characteristics of the brain. Functional magnetic resonance imaging (fMRI) has been widely used to objectively quantify functional aberrations and identify the underlying neural substrates in the human brain. This has led to decades of research that documented the associations

This work has been submitted to the IEEE for possible publication. This work was supported in part by AcRF Tier-2 grant MOE T2EP20121-0003 of Ministry of Education, Singapore.

Yi Hao Chan, Deepank Girish, Jing Xia, Chockalingam Kasi, Yinan He, Conghao Wang and Jagath C. Rajapakse are with the School of Computer Science and Engineering, Nanyang Technological University, Singapore.

Sukrit Gupta is with Department of Computer Science and Engineering, Indian Institute of Technology, Ropar, India.

Corresponding author: Jagath C. Rajapakse (e-mail: AS.Jagath@ntu.edu.sg)

between neurological diseases and disruptions in whole-brain functional connectivity (FC) [1]. However, these characterisations have yet to reveal strong and reproducible biomarkers for most disorders [2], [3]. This lack of success is commonly attributed to limitations such as disease heterogeneity [4], inter-individual variability [5], small effect size [6], [7], noise [8], limited dataset sizes [9], site effects (also known as batch effects) when combining data from multiple sites [10] and variability introduced by the choice of pre-processing pipeline [11], [12]. Over the past years, larger datasets have emerged through inter-institution collaborations [13], more mature pre-processing pipelines are available (for example, [14]) and better harmonisation tools have been developed to reduce batch effects [15], alleviating some of the above-mentioned issues in fMRI studies. Coupled with the development of more sophisticated modelling tools, model performances have improved and potential biomarkers have been discovered in recent years, warranting the need for another review.

Many modern modelling tools involve the use of machine learning (ML) algorithms, in part due to the high dimensionality of fMRI datasets. Mass univariate approaches (where each voxel is modelled independently) have traditionally been used, but they do not capture inter-region functional relationships [16]. To address this limitation, multivariate techniques (e.g. multi-voxel pattern analysis [17]) have been proposed, involving ML algorithms such as support vector machines (SVM). More recently, deep learning models have been shown to outperform SVM in disease classification tasks [18]. While convolutional neural networks (CNN) customised for connectome datasets [19], [20] were proposed as an improvement over vanilla deep neural networks (DNN) [21], graph neural networks (GNN) have since emerged as the state-of-the-art deep learning model used in network neuroscience studies [22]. Besides being an intuitive fit to FC matrices (which are best represented as graphs), the flexibility afforded by GNNs makes it possible to design techniques that capture intra-modular relationships [23] or even inter-patient relationships [24]. While vanilla GNNs do not seem to do better than DNNs and CNNs [25], more carefully designed GNN architectures have demonstrated significant improvements in disease classification performance [26], [27].

However, disorder prediction is rarely the end goal as clinical adoption of such models is rare [28]. Instead, these

models can be used to provide neurological insights. This is most commonly demonstrated via model explainability techniques (henceforth termed as ‘explainers’), ranging from attribution methods such as Integrated Gradients (IG) [29] or perturbation-based approaches such as GNNExplainer [30]. Alternatively, GNN-specific mechanisms such as graph pooling [31] can simultaneously train the model and explain the model’s decision. These explainers typically assign attribution scores to each feature, or identify important subgraphs that contribute most to the model’s predictions. While several explainers have been used in existing studies, many GNN-specific explainers remain unexplored. Furthermore, little has been studied about the robustness of these explanations. Existing studies often provide very limited evaluation of their biomarkers, only reporting the top few features and cross-referencing other studies that had similar findings. To ensure that salient features are truly representative of disorder traits (and not mere artefacts), it would be prudent to establish new standards of reporting salient features.

In this review, we summarize recent progress on neurological disorder prediction via GNNs, with a focus on reviewing explainers used and potential biomarkers discovered by these fMRI studies. Neurological disorders included in this review include attention-deficit hyperactivity disorder (ADHD), autism spectrum disorder (ASD), major depressive disorder (MDD), schizophrenia (SZ), dementia (including Alzheimer’s disease (AD), mild cognitive impairment (MCI) and other forms of non-AD dementia) and Parkinson’s disease (PD). We investigate the reproducibility of these biomarker discovery studies and propose numerous promising research directions to improve their robustness.

A. Related studies

Several review papers have been written on topics such as the use of GNNs for fMRI data analysis, model explainability in GNNs and biomarker discovery for various neurodegenerative diseases and neuropsychiatric disorders. However, none has attempted to study these topics in a cohesive manner. Recent reviews on GNN applications in network neuroscience [22], [32] summarised various graph-based and population-based models that have been proposed. Several benchmarking studies [25], [33], [34] have also been conducted to analyse the efficacy of various GNN designs on fMRI data in a fair and controlled manner that minimises the effect of covariates. Our study builds on top of their findings by synthesizing their key takeaways and introduces additional considerations on model interpretability, so as to access the suitability of these methods for biomarker discovery.

Model explainability methods have been thoroughly reviewed in previous works [35]–[37], including methods that are specialised for GNNs [38], [39]. However, to the best of our knowledge, no review has been done to assess the relevance of these methods in the context of biomarker discovery from fMRI datasets. Existing reviews on biomarkers [1], [2] tend to summarise various types of biomarkers (often going beyond fMRI) and do not focus on evaluating the reliability of the computational techniques [40], [49] used to derive the biomarkers. In this review, we aim to address this gap.

B. Review methodology

To identify papers that applied GNNs on fMRI datasets, we performed a search on 15 May 2023 via PubMed and Scopus with the following search query: (“graph neural networks” OR “graph convolutional networks” OR “GNN” OR “GCN”) AND (“fmri” OR “functional MRI” OR “functional connectivity”). This produced 76 papers and 162 papers respectively. Manual filtering was then done to remove overlaps and exclude search results that were irrelevant (no disorder prediction, fMRI not used in experiments, no experiments conducted). For papers on disorder prediction, only papers within our defined scope (ADHD, ASD, MDD, SZ, dementia, PD) were included. After this filtering process, a total of 89 papers were reviewed.

II. MODELLING FUNCTIONAL MRI DATASETS FOR DISORDER PREDICTION

Blood-oxygen-level-dependent signals captured in fMRI datasets are fundamentally represented as time series data from individual voxels. Even at relatively low resolutions (e.g. 5mm), the number of voxels (> 10,000) far outnumbers typical dataset sizes. Coupled with the issue of low signal-to-noise ratio (SNR) in fMRI data [41], these problems have motivated researchers to group clusters of related voxels together to improve SNR. Examples of such techniques are atlas-based approaches, independent component analysis (ICA) [42] and functional gradients/manifolds [43]. The former applies atlases developed by delineating boundaries following anatomical landmarks (such as sulci [44] and gyri [45]) or task-fMRI experiments that identify regions of interests (ROI) that are activated when performing various tasks [46]. This provided an informed way of feature selection to reduce data dimensionality, with the disadvantage of neglecting the voxels that are not part of the atlas’ ROIs. On the other hand, ICA and functional gradients take a different approach by learning lower-dimensional representations of the original data.

A. Functional connectivity

Besides modelling the mean time series of each ROI or component directly, another common way to analyse fMRI data is to study the relationship between pairs of time series data. Pearson correlation has been the most common approach of computing such FC matrices. However, it is limited to capturing linear correlations and could have weaker connections suppressed by noise or imaging artifacts. Alternative metrics introduced to address these limitations include various forms of partial correlation and sparse representation [47].

There are two key paradigms of modelling FC: static FC (sFC) and dynamic FC (dFC). sFC assumes that FC is stable throughout the duration of the scan (i.e. Pearson’s correlation is computed from the entire time series, without splitting it into parts). On the other hand, dFC does not make such an assumption. The most common approach involves using a sliding window across the time series, generating multiple FC matrices in the process. More sophisticated approaches perform clustering and decomposition on these matrices to produce a single dFC matrix that better captures FC dynamics than the sFC matrix. Du *et al.* [48] provides a more detailed review of dFC (as well as ICA-based methods).

B. Graph representations of fMRI datasets

fMRI datasets, especially FC matrices, are most naturally represented as graphs. A graph data structure consists of nodes that are linked by edges, which represent the relationship between connected nodes. There are two major paradigms of modelling fMRI datasets: brain graph (BG) and population graphs (PG). Fig. 1(a) illustrates the differences between them and lists down several possible options when choosing the node features and adjacency matrix of the BG and PG.

In a BG, each node in the graph represents an ROI and the choice of node features varies across studies - most commonly, the connection profile for that ROI (i.e. the row in the FC matrix that corresponds to that ROI) is used. Edges store quantitative measures of the relationship between ROI pairs (e.g. Pearson’s correlation of the mean time series). Thus, each subject (or scan) is represented as a graph and graph classification is typically performed. On the other hand, in a PG, each node represents a subject (or scan) and the graph represents the population of interest. Node classification is typically performed. Node features typically store some representation of the imaging data (often after dimensionality reduction via recursive feature elimination (RFE) or principal component analysis (PCA)). Edges store measures of similarity between scans (usually demographic information such as age and gender, or any vector from which similarity can be computed). Another alternative is to construct a k-nearest neighbours (k-NN) graph [51], [52].

PGs allow a much wider variety of data (including metadata) to be incorporated into the analysis. However, such graphs are typically pre-defined rather arbitrarily and tend to be static [53]. Thus, learnable [54] and adaptive [26], [55] methods of PG construction have been proposed. Recent studies [27], [56]–[58] have also explored ways to use BG and PG simultaneously.

C. Encoding fMRI data with graph neural networks

Let $G = (V, A, X)$ represents a graph used by the GNN, where V represents the set of nodes in the graph, $A \in \mathbb{R}^{|V| \times |V|}$ represents the adjacency matrix used and $X \in \mathbb{R}^{|V| \times K}$ represents the node features, each of length K .

Traditionally, network-based analyses have been performed on fMRI datasets for disease studies, revealing insights such as lower clustering coefficient, global efficiency and node degree for patients with mild cognitive impairment [1]. With the advent of ML, researchers started training models that distinguish between healthy subjects and patients. Since graphs cannot be used as input to many of these models, features were handcrafted (e.g. using network metrics [59]) or in the case of FC matrices, vectorised by flattening the lower triangular [21]. Doing so loses the graph structure and leads to a high dimensional input. Thus, models tend to overfit and feature selection methods (such as two-sample t-test and recursive feature elimination) are often used to address these issues [60].

Recently, GNNs - neural networks that are designed to be applied directly to graphs - have been used for encoding fMRI datasets. GNNs provide a parameter efficient [31] manner of modelling FC matrices, reducing the problem of overfitting

and allowing for more sophisticated analyses involving modular brain networks [61] and multimodal data [57]. GNNs can be broadly categorized into spectral GNNs and spatial GNNs.

Spectral GNNs perform convolution by transforming the graph signal and filtering to the spectral domain before convolving. The convolution operation is defined as:

$$g \star x = U(U^T g \otimes U^T x), \quad (1)$$

where $U^T x$ converts a signal x to the spectral domain using graph Fourier transform and Ux transforms the signal x back to the spatial domain using inverse graph Fourier transform. The operation can be simplified as:

$$g_\theta \star x = U g_\theta U^T x, \quad (2)$$

where g_θ denotes a learnable diagonal matrix. Examples of spectral GNN include ChebNet [62] and a recent improvement of it called ChebNetII [63] which noted that ChebNet can learn inappropriate Chebyshev coefficient which results in overfitting and sub-optimal performance. ChebNetII uses Chebyshev interpolation to overcome this issue and demonstrated better performances.

Spatial GNNs, on the other hand, apply convolutions to the graph based on its topology. Adopting a message-passing paradigm, node features in node i are iteratively updated by aggregating node features from its neighbours $\mathcal{N}(i)$.

$$X'_i = \alpha_{i,i} X_i U + \sum_{j \in \mathcal{N}(i)} \alpha_{i,j} X_j W, \quad (3)$$

where $U, W \in \mathbb{R}^{K \times K'}$ represent learnable weights and α represents coefficients learnt to balance between retaining information from the original node vector and using information gathered from its neighbours.

In general, variants of spatial GNN differ in how they weigh and aggregate information. For example, graph attention network (GAT) [64] learns attention scores to choose which neighbours to focus on, while graph isomorphism network (GIN) [65] learns to balance between information from neighbours and the node’s own node features. Note that graph convolutional network (GCN) [66] can be seen as both spatial GNN and spectral GNN, as shown in Fig. 1(a).

Beyond baseline GNNs, several graph convolution layers and pooling layers have been proposed to be used for connectome datasets. Many baseline GNNs focus on node-level aggregation, but FC is often used as the graph of BGs. This motivates the use of edge features as well as edge-based convolutions such as EdgeConv [67] or GraphConv [68]. Fig. 1(a) also provides a summary of common baseline pooling layers. More details about pooling layers will be discussed in Section III-C.3 as they are often used as a means to improve model explainability.

Details about the above GNN architectures have been reviewed in [32] and several benchmarking studies have evaluated the efficacy of GNNs on fMRI datasets. Thus, we will summarise key points from these papers and focus on aspects pertaining to model explainability. BrainGB [33] splits the GNN methodology into 4 parts: node feature construction,

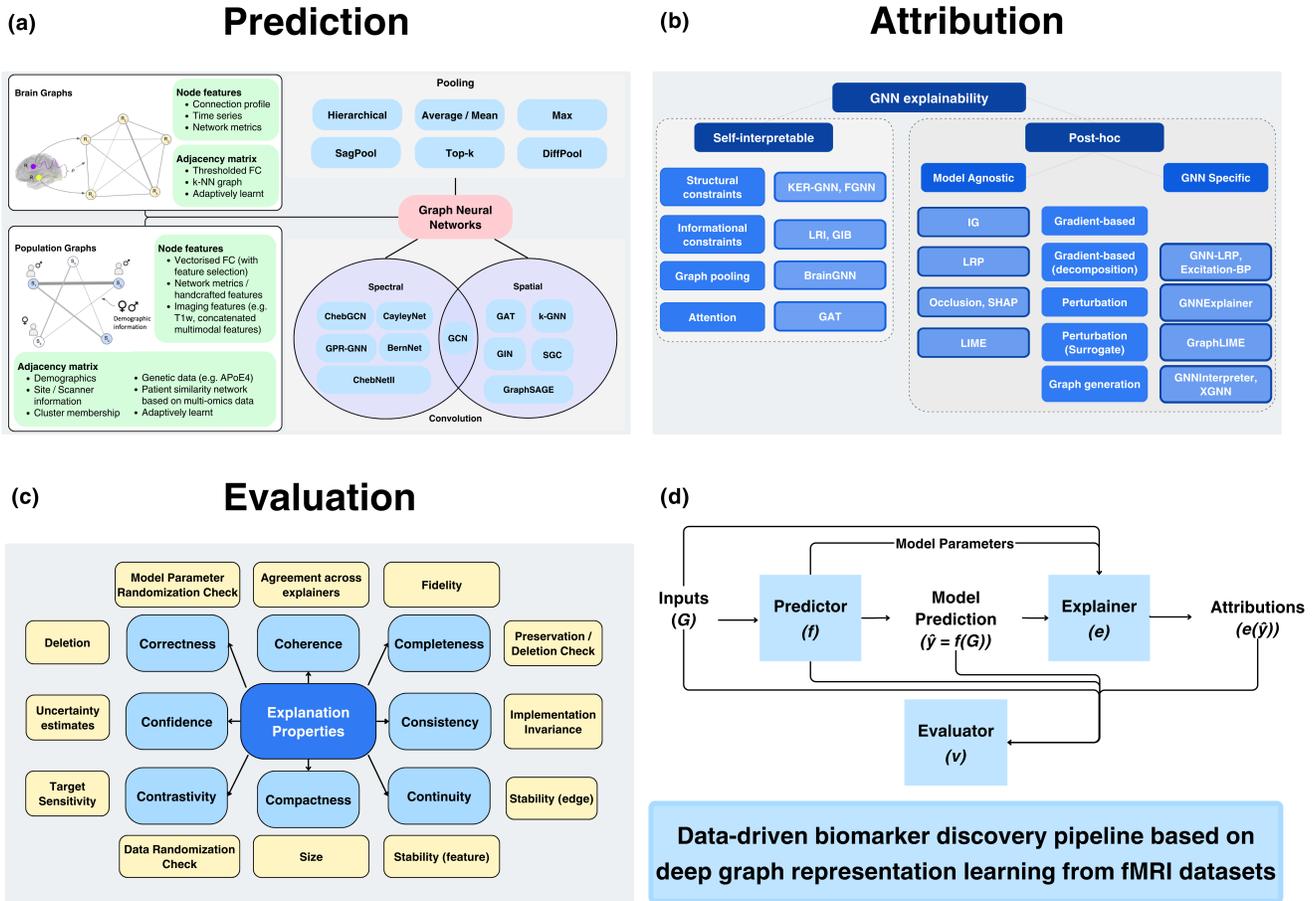


Fig. 1. Overview of the three key stages in a typical data-driven biomarker discovery pipeline that uses graph neural networks on fMRI datasets. (a) Two main types of GNNs are spectral GNN and spatial GNN. For both types of GNNs, pooling layers are often used to aggregate representation learnt from each node, especially in brain graphs (BG). Combinations of node features and adjacency matrices used in BG and population graphs (PG) varies widely across studies. (b) Taxonomy of explainers applicable to GNNs. Generally, post-hoc methods (typically ‘black box’) have a larger range of algorithms than self-interpretable (typically ‘grey box’) ones. (c) A subset of the Co-12 properties [49] used to evaluate explanations produced by explainers, deemed to be relevant for biomarker discovery from fMRI datasets. (d) An illustration showing where the three key stages occur in a typical pipeline that uses a post-hoc explainer. For intrinsically interpretable models, the explainer would become part of the predictor (e.g. pooling layers can be used for producing explanations, but they form part of the GNN architecture). We further note that it is possible for explanations to be made based on the ground truth labels (i.e. phenomenon) instead of the predictions (i.e. phenomenon) [50]. Overall, to ensure that robust biomarkers are produced from the GNN models, the choice of predictors, explainers and evaluators needs to be made carefully.

messaging passing mechanism, attention-enhanced message passing, and pooling strategies. On multiple datasets (both healthy subjects and patients), they showed that connection profile (i.e. use the corresponding row of the FC matrix), node concat message passing with attention (i.e. multiply learnt attention weights to the neighbour’s node feature before concatenating it with the node’s node feature vector, followed by a multi-layer perceptron), and concat pooling (i.e. concatenate node features when performing graph pooling) works best.

Neurograph [34] applied various spectral and spatial GNNs (GCN, GAT, GIN, etc.) on healthy subjects from the Human Connectome Project dataset and showed that their GNN* architecture (which has 3 graph convolution layers with skip connections) works best across various baseline problems such as task classification, gender classification and age prediction. They also experimented with various settings such as number of ROIs used, sparsity of graphs and how the node features were created. They found the following factors to drive

model performance: including more ROIs (400 and 1000), using sparser graphs (5%) and using Pearson correlation for node features. It is notable that another benchmarking study [25] suggested that GNNs do not even outperform 1D CNN for MDD and ASD classification. However, they opted for a relatively dense FC matrix (50-90%) and binarised the graph.

Overall, all three benchmarking papers noted that sparsity of the graph used by the GNN impacts model performance and lower sparsities (below 50%) were shown to be beneficial. Hyperparameter and GNN construction choices that lead to more model parameters (e.g. more ROIs, concatenation of node features) seem to enhance performance, but this should be done with care even though GNNs are parameter efficient (relative to other DNNs). Finally, vanilla GNNs do not seem to clearly outperform non-graph baselines such as SVM and CNN, but more carefully designed GNNs do improve model performance. We note that these benchmarking studies (and this introduction on GNNs) have not considered state-of-the-

art GNN architectures that are customised for fMRI datasets (e.g. GNNs that incorporate information about brain modularity [69]). These are covered in more detail in Section IV.

III. BIOMARKER DISCOVERY FROM FMRI DATA

Biomarkers are biological signatures that can be objectively measured and used to indicate physiological processes, pharmacological responses or disease status [70]. In this paper, we focus our discussion on potential biomarkers of brain disorders derived from fMRI via GNNs. Since fMRI based classification models are rarely adopted clinically [28], studies involving such models often provide additional insights via various model explainability approaches. These salient features detected from the trained models are potentially disease biomarkers, but it is uncertain whether they capture true signals of the disorder, or are mere artefacts of noise. To ascertain the validity of these potential biomarkers, they need to be evaluated against the following traits: easily accessible, reproducible, specific and sensitive. Most importantly, it should only change depending on the state of the disease and be unchanged when unrelated factors are varied [71]. Most existing studies on biomarkers for brain disorders have not attained such high standards yet and bridging this gap has been a priority for neuroimaging research in recent years [72].

A. Overview of existing biomarker discovery studies

In this section, we first give a broad overview of existing biomarker discovery studies using fMRI data. Early studies analyzed neuroimaging data to study disorders using univariate [73] or multivariate statistical methods (e.g. multivariate pattern analysis) [74], [75]. These methods provide coefficients of variables showing the significance of biomarkers and enable the generation of statistical models with high diagnostic or predictive potential by focusing on patterns of brain changes that are distributed across multiple regions in disordered states.

ML techniques have also been commonly used to select important features for brain disorder classification [76], [77]. A model that generalizes well to the test set and other unseen data points would be expected to have learnt optimal coefficients for each feature and feature importance could be inferred from these values. Features with the highest importance scores could represent potential biomarkers. For example, SVM is a supervised learning model with associated learning algorithms (and kernels for learning linear or non-linear relationships) that allows for classification and regression analysis. SVMs have yielded high accuracies ($\sim 90\%$) when applied to small fMRI datasets ($n < 100$) of SZ patients and highlighted ROIs with significantly lower brain activities [78]. Random forest (RF) exhibits significant advantages over other ML methods in terms of their capacity to handle highly non-linearly correlated data, simplicity of hyperparameter tuning, and robustness to noise [79], [80]. Their variable importance can be evaluated using a variety of methods, including calculation of the prediction power of selected features in classification using the impurity reduction principle.

With the advent of deep learning techniques, multiple layers of non-linearities are introduced to learn complex relationships

between the input and outputs. Furthermore, GNNs make it possible to model FC directly, obviating the need for feature engineering. Although they are often viewed as black boxes, the DNN's decisions could be analysed via model explainability algorithms such as gradients [81], IG or class activation mapping (CAM) [82].

To consolidate the findings from these studies, one could conduct meta-analyses by calculating the contribution of identified biomarkers to specific disorders. For example, BrainMap [83] is a database of MNI coordinates for activation foci consolidated from thousands of experiments. Researchers can extract existing relevant studies from BrainMap using specific keywords, revealing experiments where both activated and non-activated ROIs were pinpointed and indicating their relevance to various disorders. Subsequently, methods such as the Naive Bayes classifier can be used to determine the probability of disorders associated with these ROIs.

Despite such progress in modelling techniques and efforts to curate larger datasets, biomarkers for brain disorders remain elusive. For instance, efforts to identify diagnostic biomarkers for depression could not arrive at any consistent depression biomarker despite extensive efforts [84]. Thus, in the next sections, we introduce the literature of model explainability techniques (both model-agnostic and GNN-specific) with the goal of identifying research gaps (in existing fMRI studies on biomarker discovery) that could be addressed to improve the quality of these potential biomarkers.

B. Properties of model explainability techniques

Many model explainability algorithms, also known as 'explainers', have been proposed and they have been covered by numerous review papers [35]–[38]. We provide a brief summary of these explainers and place greater focus on key insights that are relevant to biomarker discovery. Before delving into the details of each explainer, we note several properties of explainers that are useful to characterise them.

Existing research on model explainability can be separated into three paradigms: (i) 'glass box' (intrinsically interpretable), (ii) 'black box' (reliant on post-hoc explainability methods), (iii) 'grey box' (some interpretation possible, with careful design of the model) [85]. On one end of the spectrum, 'shallow' models such as linear regression are intrinsically interpretable. In the case where all input features have the same scale, biomarkers can be extracted by identifying features that have the largest coefficients assigned to them by the model fitting process. On the other end of the spectrum, deep learning models learn complex and non-linear relationships that cannot be easily interpreted. They rely on post-hoc model explainability algorithms that are applied after model training. These algorithms typically generate scores based on some form of gradient computation (with respect to the input) or perturbation. In between these two extremes, some complex models such as fuzzy rule-based systems and Bayesian networks can provide a limited extent of interpretability [85]. The use of attention scores as well as graph pooling could also be grouped under this category.

While 'glass box' methods are desirable, most of these methods only capture linear relationships, which is likely to

be of limited use for disease studies. Unlike how studies on using fMRI data to predict phenotypic information have shown that non-linearities do not give much improvement over linear models, deep learning models have shown better performance for disease classification and prediction of clinical test scores [18]. However, many of these models fall under the ‘black box’ category and more research is needed to create ‘grey box’, or even ‘glass box’ alternatives.

All three paradigms of explainers produce attribution scores, which can be classified into two categories: local (‘instance-level’) and global (‘model-level’). Local explanations are specific to the input data provided to the model (i.e. each sample has its own attribution scores), while global explanations apply broadly to the entire model (all samples share the same attribution scores). In the context of biomarker discovery, local explanations could be more desirable for clinical use if individual insights are found to be reliable. Additionally, global explanations are unlikely to be helpful for very heterogeneous diseases since heterogeneous diseases would not be fully described by a single global explanation.

C. Model explainability for graph neural networks

GNNs can be difficult to interpret due to the non-Euclidean contextual information in the graph (node features and edge weights) that needs to be taken into account when computing explanations. For BG, explanations can be produced at the level of nodes (i.e. which ROIs contribute the most to the disorder), edges (i.e. which pairs of ROIs contributed most) and node features. For PG, explanations can be produced at such granularity too, but they carry a different meaning: nodes would correspond to patients, edges correspond to pairs of patients and node features would typically correspond to some imaging data. Not all explainers can produce all 3 types of explanations. Thus, it is important to first understand which types of explanations the method can provide. Table S2 in the supplementary materials summarises these key characteristics of explainers highlighted in Fig. 1(b).

Existing model explainability methods for GNNs can be split into two major groups: self-interpretable and post-hoc. Self-interpretable methods provide explanations simultaneously with the model predictions, while post-hoc methods are only applied after model training is complete. Self-interpretable methods typically include constraints to extract an informative subgraph, or architectural designs where weights that prioritise a subgraph are learnt. Post-hoc techniques can be further split into model agnostic and GNN specific techniques. Model agnostic algorithms can be applied to any deep learning model (and for GNNs, regardless of their internal structure), but some of them have been further extended to capture the graph structure [86]. GNN-specific methods explicitly consider the graph structure when generating the explanations. Model agnostic methods can be categorised into two major groups: gradient-based and perturbation-based. GNN specific methods follow such a characterisation too, but also has unique ones such as techniques based on graph generation. Fig. 1(b) summarises the taxonomy of GNN explainability methods and in the following subsections, each subcategory will be explained in detail.

1) *Self-interpretable - Structural constraints*: Drawing inspiration from CNNs, kernel GNN (**KerGNN**) introduces learnable graph filters (akin to convolutional filters in CNNs) into the messaging passing process in GNNs. Contrary to the rooted subtree approach used in GNNs that are based on the message-passing paradigm, KerGNN updates each node’s embedding based on subgraphs that are centered on the node. This is done via the use of graph kernels (specifically, the random walk kernel) to measure the similarity between subgraphs and the graph filters. Not only does such an approach makes GNNs more expressive (going beyond the 1-dimensional Weisfeiler-Leman (1-WL) limit), it also provides additional interpretability via the learnt graph filters, just like how convolutional filters can be visualised.

Such an approach would be useful in the case of disorders that are poorly understood as novel insights could be drawn from these visualisation. However, if there is some existing biological knowledge available, Factor GNN (**FGNN**) [87] could be considered. FGNN directly incorporates biological knowledge as the inductive bias into the model via a factor graph that involves two types of variables: observable variables and latent variables. In the original paper, the factor graph was built using gene ontology (GO) and genes, which are the latent variables and observable variables respectively. Unlike usual deep learning models where the hidden nodes in deep learning models do not have a physical meaning, a node in FGNN represents a biological unit (i.e. each feature represent a gene, while each hidden node represents a GO term). Instead of being fully connected, genes are only connected to GO terms if and only if the gene is included in the GO term. Then, the hidden layer could be used as input to a fully connected layer to make predictions (e.g. clinical outcomes), or be passed to another stack of factor graphs, forming a deep network. In the context of biomarker discovery from fMRI, brain modules could possibly be used as latent variables (while ROIs represent the observable variables).

2) *Self-interpretable - Informational constraints*: Instead of introducing architectural designs that guide the learning process (and identify salient subgraphs), methods based on information constraints introduce information bottlenecks such that the mutual information (MI) between the labels and the discovered subgraph is maximised, while keeping the MI between the original graph and subgraph below a predefined threshold. While the former can be approximated via cross entropy loss, the latter is estimated via a plethora of techniques such as learnable randomness injection (**LRI**) [88] and graph information bottleneck (**GIB**) [89]. The latter has been further developed in the fMRI biomarker discovery literature by BrainIB [90]. It extends GIB by considering the effects of edges (not just nodes) during subgraph discovery.

3) *Self-interpretable - Graph pooling*: Graph pooling is often performed in GNN architectures, especially in graph classification tasks where node features have to be condensed to a lower dimensionality or a single vector. Pooling techniques can be grouped into two categories: (i) flat pooling, where a graph-level representation is generated in one step, and (ii) hierarchical pooling, which gradually coarsens the graph by clustering nodes together or dropping some of them [91].

Several pooling techniques customised for FC data have been proposed. Rather than using max/mean pooling operations, Li *et al.* [31] proposed a node / ROI pooling layer (R-pool) in their **BrainGNN** architecture. R-pool projects the node feature embeddings to a learnable weight vector and retains nodes with the highest scores. Hierarchical pooling approaches that consider functional modules have also been proposed [61]. In this work, three levels of hierarchy were used: (i) ROIs belonging to the same sub-network (e.g. Yeo 7-network parcellation [92]) and brain hemisphere, (ii) the pair of matching sub-networks from each hemisphere, (iii) combining all sub-networks into a whole brain network. Weights from the final pooling layer were used to identify the sub-networks that contributed most to the model’s decision.

4) Self-interpretable - Attention: The widespread use of attention for model interpretability has also been present in the GNN literature, most popularly via self-attention in **GAT**. Attention scores have been used to identify salient FC features as well [93], [94]. However, there have been much debate in natural language processing (NLP) research about whether attention scores provide meaningful explanations [95]–[98]. In NLP applications, attention scores were found to not correlate well with multiple gradient-based approaches on recurrent neural networks and different attention distributions can lead to equivalent predictions [95]. However, Wiegrefe *et al.* [96] argue that explanations should be further categorised into ‘plausible’ or ‘faithful’ explanations. Their results provide additional support that attention does not provide faithful explanations, but does not invalidate claims that attention provides plausible explanations. These results suggest that greater care should be taken when using attention scores to discover potential biomarkers. Further research is needed to examine the validity of attention scores for biomarker discovery applications.

5) Post-hoc - Gradient-based: A large variety of gradient-based approaches have been proposed. Integrated Gradients (**IG**) [99] will be discussed here due to its versatility (works on most deep learning models where gradients can be calculated) and widespread use. IG was developed to address the issue of saturation in gradient-based attribution methods. With saturation, the output is no longer sensitive to small changes in the input features, making it difficult to interpret which features are responsible for the prediction of the correct class [100]. IG avoids this issue by accumulating gradients from interpolated points (controlled by α) between a baseline (x') and the input data (x). In the context of disorder classification, the baseline can be the average data of all healthy subjects. For a given feature k , IG is defined as:

$$IG_k(x) = (x_k - x'_k) \times \int_{\alpha=0}^1 \frac{\partial f(x'_k + \alpha \times (x - x'_k))}{\partial x_k} d\alpha.$$

Attribution scores from IG provide a local measure of how much the feature contributed to the model’s prediction and could be potentially useful for producing personalised insights. However, IG could also produce noisy pixel attributions in features unrelated to the predicted class. Modified versions of IG, such as GuidedIG [101], have been proposed to address

these issues. Features with high gradient scores have a greater impact in predicting the class of the model than features with low gradient scores. Thus, GuidedIG reduces noise in the attribution results by using an adaptive path technique that only incorporates a subset of features with high gradient scores.

6) Post-hoc - Decomposition: Decomposition-based approaches have some overlaps with gradient-based approaches, but differ in the way the scores are computed. Instead of computing the gradients directly with respect to inputs, scores are decomposed starting from the output layer and propagated backwards in a layerwise manner based on pre-defined rules. Layerwise relevance propagation (**LRP**) is one such example. Many forms of LRP exist and the variant called ϵ -LRP will be discussed [102]. Starting from the output layer, a score s is assigned to a neuron based on the logit, i.e.

$$s_i^l = \frac{h_{ji}}{\sum_i h_{ji} + \epsilon(\sum_i h_{ji})} s_i^{l+1},$$

where h_{ji} refers to the output from neuron i in layer l to neuron j in layer $l+1$. This is computed layerwise, reallocating the prediction score until the input layer is reached. The total relevance score of s is always preserved for each layer.

LRP does not consider the adjacency matrix in its computations. To address this, **GNN-LRP** [103] distributes scores to different graph walks (and thus have higher computational complexity). **Excitation-BP** is very similar to LRP, but it views the decomposition process from a probability standpoint. Overall, recent decomposition-based approaches like GNN-LRP have not been studied much in fMRI, with usage mainly limited to LRP. [104]

7) Post-hoc - Perturbation: Perturbation-based approaches introduce changes to the input with the motivation that if important features are still retained, the outputs should remain similar. In its simplest implementation (**Occlusion**), this involves masking features one by one and the feature that results in the largest change in output would be deemed to be the most important. **SHAP** [105] takes this idea to completion by considering all feature subsets (i.e. 2^k combinations), so as to compute Shapely scores which have been proven to be the unique solution that fulfills the criteria of local accuracy (model training on best feature subset should have similar predictions with the original model), missingness (features not in the best subset should have no impact on the model output) and completeness (attribution score should not decrease when a different model, where the feature contribution does not decrease, is used). Computing this would be too computationally expensive, thus it is achieved via approximation techniques.

In the context of graphs, such perturbations can be performed by discovering subgraphs. **GNNExplainer** [106] produces local explanations for GNN predictions by selecting a small subgraph from a given input graph and identifying important node features. Subgraphs are generated by randomly masking nodes in the graph and observing the resulting changes in the model’s prediction. A soft mask (i.e. continuous values, not binarised) containing learnable weights is used. Important node features (that are in the nodes within the subgraph) are identified by a binary feature selector. The mask and feature selector are optimised by maximising the mutual

information between the original model predictions and the model's predictions given the masked graph. One limitation of GNNExplainer is that the subgraph must be connected, which might not always be applicable to disease biomarkers.

8) *Post-hoc - Surrogate*: Surrogate-based approaches have some overlaps with perturbation-based approaches as they tend to rely on perturbations too. However, a distinctive feature is the use of simpler and interpretable models (often linear) to approximate the original complex model. This is possible as it limits the approximation to a local neighbourhood and analyse the model predictions of perturbed inputs within this neighbourhood. For instance, local interpretable model-agnostic explanations (LIME) [107] trains a surrogate model on the dataset of perturbed points, weighing them based on their proximity to the chosen data point.

GraphLIME [108] is a non-linear version of LIME, a key difference being that it uses Hilbert-Schmidt Independence Criterion (HSIC) lasso, a kernel-based non-linear interpretable feature selection algorithm. The algorithm first computes the importance of each feature in each node by considering the features of the target node and features in the N-hop neighbouring nodes. The given target node will aggregate the information from N-hop network neighbours to identify the most significant features. The HSIC lasso method is used to train a linear interpretable model to represent the relationship between the features and target node prediction. Subsequently, the coefficients from the linear interpretable model will be used to identify the top few features that are important for model prediction based on the coefficients. Thus far, no research on FC has used such models for biomarker discovery.

9) *Post-hoc - Graph generation*: Explainers based on graph generation bear some similarities with GNNExplainer as both identify salient subgraphs. However, generation-based approaches arrives at the subgraph via generative approaches, instead of perturbation. **XGNN** [109] interprets GNNs using a graph generator to identify important graph subgraphs. The generator is trained using reinforcement learning (RL) and validity rules are defined by pre-existing knowledge, making them unsuitable for biomarker discovery. On the other hand, **GNNinterpreter** [110] do not require pre-existing knowledge as it optimises the choice of subgraph by maximising the similarity between embeddings from the important subgraph with that of the average graph embeddings in the target class. Both explainers are different from all other explainers discussed above as they produce global explanation (i.e. one set of explanations for the whole model, across all data points). Thus far, no studies on FC have used explainers based on graph generation for biomarker discovery.

10) *Summary*: Overall, explainers based on informational constraints, pooling, attention, gradients and perturbation have been used on FC datasets. However, there remains much room to explore alternative explainers based on graph generation, decomposition, surrogates and structural constraints. We note that our review do not include temporal GNNs (which might be useful for dFC applications), counterfactual-based approaches and causality-based approaches. We refer interested readers to Kakkad *et al.* [39] for a review on these topics.

D. Evaluation of explanations

Better biomarker discovery tools are needed as few potential biomarkers turn out to be effective in clinical settings [72]. To bring us closer to this goal, one solution could involve developing objective means of assessing the robustness of saliency scores produced by explainers. Several methods of evaluating generic explainers [49] as well as GNN-specific explainers [40], [111] have been proposed. Nauta *et al.* [49] proposed a framework named 'Co-12' which encompasses 12 desired properties of explainers (such as Correctness, Completeness, Consistency, etc.). In this section, we discuss properties that are relevant to biomarker discovery and extend the analysis to GNN-specific explainers. Fig. 1(c) illustrates the 8 chosen properties and metrics that can be used to measure a model's performance with respect to the corresponding property. Notably, we excluded Controllability (since it is applicable only to interactive explanations), Context (since it pertains to user-specific needs which do not seem relevant for biomarker discovery), Composition (concerns about format and organisation of explanations are less apparent in graphs than in text and images) and Covariate complexity (which relates to having human-understandable explanations, but FC studies typically uses ROIs as features which is already understandable, unlike individual pixels in images).

Each property below encompasses one or more evaluation metric. Many metrics involve comparing two distributions and measuring the distance between them. This can be computed via Hellinger distance, which has an easily interpretable range (0 = perfectly similar, 1 = completely different distributions).

1) *Correctness*: Sanity checks on the explainers can be performed via **model parameter randomisation check**, which verifies whether the explanation of a trained model is different from a randomly initialised untrained model. Other forms of evaluating correctness involve **deletion**: changes in model outputs are computed for each feature subset (in the simplest case, independently removing each feature in the node vector) and then correlated with the importance scores. However, this could be too tedious for deep learning models applied on connectome datasets especially if many ROIs are used.

2) *Completeness*: Salient subgraphs or feature subsets identified by the explainers should not lose too much information as compared to the original input. This can be assessed via **preservation check** (whether using the selected features as input results in the same model prediction) and **deletion check** (not using the important features results in a different prediction).

Fidelity is based on a similar idea, but quantifies the difference by measuring the change in probabilities (as opposed to a simple change in prediction of classes). Note that this can be done both at the level of node features and (in the case of graph datasets like FC) at the level of subgraphs.

3) *Consistency*: Explanations should be robust across most variations in model implementations. While deep learning models can involve many hyperparameters that have profound and poorly understood changes to the model, a basic form of assessment for **implementation invariance** could entail checking whether two models with the same architecture,

but having different initialisation (i.e. randomly initialised weights), will produce similar outputs and explanations.

4) *Continuity*: Given a small change in the input that still produces a similar model output, it would be reasonable to expect that the explanation will remain similar too. **Stability** ensures this by introducing random noise to the input features. For graphs, additional changes can be introduced by randomly switching edges while keeping the number of edges constant. Both explanations (before and after perturbation) can then be compared to determine if the scores have good continuity.

5) *Contrastivity*: Explanations produced for data samples belonging to different classes would be expected to differ. This is verified via a **target sensitivity** check, where the mean explanations for each class can be computed and then compared across classes. Alternatively, a **data randomisation check** can be performed by randomising the labels and verifying that the explanations are changed. Given that biomarker discovery techniques based on ML techniques are often based on disorder classification tasks, it would be important to check if the scores produced have high contrastivity.

6) *Compactness*: Explanations should not overwhelm the user and this is especially relevant for FC datasets, where hundreds of ROIs and thousands of connections are analysed simultaneously in multivariate studies. For instance, while severe diseases might affect widespread areas of the brain, having an explanation that gives similar scores to an overwhelming number of ROIs or edges might not be as useful as explanations that correctly emphasises on a small group of features. Computing the **size** of explanations is straightforward for explainers that produce subgraphs (e.g. number of connections divided by number of edges in a complete graph), but would require thresholding in the case of most gradient-based approaches (since a score is given to each features, but some scores could be very close to 0).

7) *Confidence*: While almost all explainers discussed above do not produce uncertainty estimates and many of their scores have no natural variations (e.g. Saliency is based on gradients with respect to inputs, which would be fixed given the same model and input), measures of confidence could be computed based on the logits of the predictor. One example of an implementation is proposed by Atanasova *et al.* [112], where class probabilities from the predictor are compared against a predicted confidence value estimated by fitting a logistic regression model to saliency distance values (which computes the distance between saliency across in each class). A similar approach could be used for biomarker discovery applications.

8) *Coherence*: Since the ground truth is often unknown in biomarker discovery, coherence would be focused on the extent of agreement between explainers. Biomarkers that are present across explainers based on different approaches could be deemed to be more robust. However, further studies need to be conducted to determine the explainers to use (e.g. explainers with low stability should not be considered).

In summary, several evaluation metrics have been proposed to assess the robustness of explainers. However, most existing studies do not thoroughly evaluate the biomarkers they propose, often limiting their evaluation to cross-referencing with existing literature. Addressing this significant research gap by

incorporating evaluation metrics in their analysis would help to ensure the robustness of highlighted potential biomarkers.

IV. APPLICATION OF GNNs ON DISORDER PREDICTION AND BIOMARKER DISCOVERY FROM FMRI DATA

In the following subsections, we summarise key insights obtained from the studies considered in this review. Besides highlighting predictors that perform well and their feature attributions, a key goal in this review is to identify salient features that are reproduced across multiple studies. Such features would then be more promising biomarkers of the disorder as they exhibit greater robustness than non-reproducible ones.

This is a challenging task as salient features identified in these studies vary in type: ROI, connections, module, modular connections, and temporal features. Furthermore, the type and quality of these potential biomarkers are influenced by numerous factors: (i) architecture design of predictor, (ii) choice of node features and adjacency matrix, (iii) type of GCN used (Defferrard [62], Kipf [66], etc.) (iv) graph construction paradigm (BG or PG). (v) dataset used (size, composition, class distribution), (vi) model performance, (vii) choice of atlas and (viii) choice of explainer. Thus, in each study summarised below, we took note of these details as concisely as possible.

To aid readability, the following conventions are adopted.

- We focus on studies that performed biomarker discovery, but also highlight (at the end of each subsection) notable studies that did not do so. Details of other remaining studies without biomarker analysis can be found in the supplementary materials.
- Disorders with few studies (ADHD, MDD, SZ, PD) are discussed in full while insights from ASD and dementia studies are condensed to keep the discussion concise.
- Since the subsections are grouped by neurological disorders, we provide the details of studies involving multiple disorders at the first instance of mentioning them.
- X represents node features and A represent the adjacency matrix used for the GNN. Connection profile refers to the ROI's corresponding row in FC matrices that are built using Pearson's correlation, unless otherwise specified.
- Names of brain atlases used are abbreviated to make it easy to note how many ROIs it involves. Table I provides a mapping of the abbreviations to the full names.
- Healthy populations are referred as normal controls (NC) or typical developing children (TDC).
- While accuracy is not the best measure of model performance, it was chosen due to its availability relative to other performance metrics. To provide a better gauge of performance, we note the size of the dataset used and the performance of baseline models reported by these studies.
- Explainers are highlighted in bold for easy reference.
- Full names for abbreviations of functional modules and ROIs can be found in Table II and Table III.

For an introduction to each disorder, we refer readers to the review by Du *et al.* [48] (as PD is not included in that review, an fMRI-specific review on PD can be considered instead [113]) as well as more recent reviews on psychiatric disorders [2] and neurodegenerative disorders [1].

TABLE I

MAPPING OF ABBREVIATIONS USED FOR ATLASES TO THEIR FULL NAMES. MORE INFORMATION ABOUT THE ATLASES CAN BE FOUND IN THE CORRESPONDING PAPERS THAT USED THE ATLAS IN THEIR STUDY.

Atlas	Full name	ROI count
AAL116	Automated anatomical labelling	116
AAL166	Automated anatomical labelling v3	166
AAL90	Automated anatomical labelling	90
BASC325	Bootstrap Analysis of Stable Clusters	325
BM82	Broadmann	82
BN273	Brainnetome	273
BNA246	BrainNet Atlas	246
CC200	Craddock	200
CC400	Craddock	400
DK308	Desikan-Killiany, backtracking from 66 ROI	308
DK86	Desikan-Killiany, with subcortical ROIs	86
DOS160	Dosenbach	160
DX148	Destrieux	148
EZ115	Eickoff-Zilles	115
HO110	Harvard-Oxford	110
HO112	Harvard-Oxford	112
JHU81	JHU ICBM-DTI-81	81
MODL128	Dictionaries of functional modes	128
Power264	Power <i>et al.</i>	264
SF200	Schaefer	200
SHEN268	Shen, group-wise spectral clustering algorithm	268
TT93	Talairach and Tournoux	93
TT97	Talariach	97
YEO114	Yeo 17-network	114

TABLE II

MAPPING OF ABBREVIATIONS USED FOR BRAIN NETWORKS TO THEIR FULL NAMES.

Abbreviation	Full name
DAN	Dorsal Attention Network
CEN/FPN	Central Executive / Frontoparietal Network
DMN	Default Mode Network
SMN	Somatomotor Network
SN	Salience Network
VAN	Ventral Attention Network

A. Attention Deficit Hyperactivity Disorder

8 studies on ADHD were found (6 sFC, 2 multimodal), but only 2 studies highlighted salient features (2 sFC). The first sFC study by Yu *et al.* [94] proposed ATT-GCN, a GCN (Kipf) with self-attention (similar to Bahdanau attention, i.e. tanh) introduced before the graph convolution layer such that the adjacency matrix is formed as a weighted combination of the attention-weighted matrix and the original FC matrix (BG). Node features underwent feature selection via PageRank scoring, keeping only the top 180 nodes with highest scores. Using the CC200 atlas on a private dataset with 240 subjects (120 ADHD, 120 TDC), ATT-GCN achieved an accuracy of 68.5% (GAT: 58.6%). Salient features based on **attention scores** highlighted connections between frontal and temporal lobes as well as frontal and parietal lobes. Specifically, this involves weakened connections between the posterior central gyrus and the frontal gyri (especially SFG and MFG). The connection between MFG and frontal pole (FP) also stands out as another significantly weakened connection. On the other hand, the connection between cerebellar regions (CEREB) with fusiform gyrus (FUS) is enhanced.

The second study on sFC by Zhao *et al.* [114] proposed dynamic GCN (dGCN), which involves a GNN with customised

TABLE III

MAPPING OF ABBREVIATIONS USED FOR BRAIN REGIONS TO THEIR FULL NAMES.

Abbreviation	Full name
ACC	Anterior Cingulate Cortex
FOG	Frontal Orbital Gyrus
FP	Frontal Pole
IFG	Inferior Frontal Gyrus
IPL	Inferior Parietal Lobule
MFG	Middle Frontal Gyrus
MOG	Middle Occipital Gyrus
MTG	Middle Temporal Gyrus
PARAH	Parahippocampal Gyrus
PFC	Prefrontal Cortex
SFG	Superior Frontal Gyrus
SOG	Superior Occipital Gyrus
SPL	Superior Parietal Lobule
STG	Superior Temporal Gyrus
TP	Temporal Pole

convolution and readout layers. A NormConv layer performs multi-hop feature aggregation (2 hops) while an EdgeConv layer performs convolution on a dynamic graph structure where the closest k-neighbors in terms of feature similarity are considered, motivated by the sparsity of connections in the brain. The readout layer involves concatenation of max and sum operators over all nodes. The graph used in dGCN is the FC matrix (BG), while the node features were set as the partial correlation coefficients. Using the AAL116 atlas on a subset of the ADHD-200 dataset with 603 subjects (260 ADHD, 343 TDC), they obtained an accuracy of 72.0% (GAT: 68.0%). They investigated the **model weights** from their proposed convolutional layers to identify the salient brain ROI. ROIs found to have the largest weights include the frontal, occipital, subcortical and temporal lobes as well as the posterior fossa / cerebellum. Thereafter, FC was studied for this subset of ROIs and connections between frontal and temporal lobe were generally found to be weakened in ADHD subjects. Numerous intra-connections within the posterior fossa and subcortical regions were also found to be weakened.

Overall, most studies on ADHD relied on subsets of ADHD-200 (5/8) and model performance is moderately high (mean accuracy of 72% across the studies). Considering the moderate size of the datasets (mean: 469) used, these results suggest that fMRI has a moderate ability of discerning between typical controls and ADHD patients. Since most studies adopted the BG approach (6/8), future studies could go further into PG or fusion of BG and PG approaches to examine whether classification performance can be improved.

However, analysis of salient features was less encouraging. While both studies discussed above (trained on different datasets) agree that ADHD subjects have weaker connections between the frontal lobe and temporal lobe, there is no agreement at the level of connections between ROIs. For instance, Yu *et al.* [94] highlighted MTG-IFG as a weakened connection, while Zhao *et al.* [114] highlighted left temporal pole of STG - right SFG (medial). More ADHD studies that look into model explainability needs to be conducted before coming to a conclusion. Furthermore, since ADHD-200 comprises 3 subtypes (combined, inattentive and hyperactive/impulsive), it could be more insightful for future

TABLE IV

SUMMARY OF FINDINGS FROM STUDIES THAT IDENTIFIED POTENTIAL BIOMARKERS. ‘SIZE’ REFERS TO SIZE OF DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL). ‘TYPE’ REFERS TO THE TYPE OF EXPLANATION.

Reference	Dataset (Size)	Atlas	Explainer (Type)	Salient features identified
ADHD				
[94]	Private (240)	CC200	Attention (ROI, connection)	Regions (-): FP, IFG, MFG, MTG, posterior central gyrus, SOG; Regions (+): Brainstem, FOG, precuneus, putamen, TP; Connections (-): IFG-MTG, MFG-FP, MFG-posterior central gyrus, posterior central gyrus-SFG; Connections (+): CEREB-FUS, FOG-Precuneus, PARAH-STG
[114]	ADHD-200 (603)	AAL116	Weights (ROI, connection)	Regions: Frontal lobe, occipital lobe, subcortical lobe, temporal lobe, posterior-fossa (cerebellum); Connections (-): right rolandic operculum - right Heschl gyrus; Connections (+): left precuneus - right CEREB, left STG/TP - right medial SFG
ASD				
[115]	ABIDE (866)	MODL128	Finetuning (ROI)	Regions: STG ; Cerebellum IV and V, central parieto-occipital sulcus, left superior temporal sulcus, left anterior intraparietal sulcus, cerebrospinal fluid (between superior part of SFG and skull)
[58]	ABIDE (871)	HO112	Attention (ROI)	Regions: IFG, precentral gyrus, frontal orbital cortex, PARAH.
[116]	ABIDE (949)	Multiple	Unclear (ROI)	Regions: angular gyrus, precentral gyrus, precuneus and thalamus
[117]	ABIDE (1112)	AAL116	Pooling (ROI)	Regions: lateral PFC, lateral dorsal PFC, superior parietal lobule
[118]	ABIDE (871)	HO111	Feature selection (connection)	Connections: evenly distributed across the brain, lower in ASD for 25/30 FC e.g. right cuneal cortex and right parietal operculum cortex
[51]	ABIDE (1057)	CC200	Occlusion (module)	Modules: DMN, FPN, VAN
[119]	ABIDE (871)	Multiple	Clustering (module, modular connection)	Modules: CEN, DMN, SN ; Modular connections: CEN-SMN, CEN-SN, DMN-SMN, DMN-CEN, DMN-Visual
[120]	ABIDE (613) [d]	HO110	Feature selection (ROI)	Regions: lingual gyrus, MFG, SFG
[121]	ABIDE (1035) [d]	CC200	Gradients (ROI)	Regions: FP, precuneus, brain stem, PCG, lingual, OP, lateral occipital cortex, frontal orbital cortex
[122]	ABIDE (867) [d]	HO110	Feature selection (connection)	Connections: right pallidum-right IFG, left frontal orbital cortex-left central opercular cortex, left supramarginal gyrus - right ITG
[123]	ABIDE (1007) [M]	AAL116	Gradients (ROI, connection)	Regions: (higher T1w): DMN, reward, memory and motor ; (higher ALFF) reward and motor ; Connections: inter >intra, low homotopic interhemispheric connection in limbic regions
[81]	ABIDE (1007) [M]	AAL116	Gradients (connection)	Connections (-): right MTG and multiple ROIs in the frontal, parietal and occipital lobes; mix of higher and lower FCs between ROIs in the limbic regions to multiple other regions

studies to analyse salient features separately for each subtype. The presence of heterogeneity in the disorder suggests that biomarker reproducibility might not always be possible (e.g. different datasets are dominated by different subtypes of ADHD) and a subtype-level analysis would have more potential of highlighting replicable biomarkers.

B. Autism Spectrum Disorder

39 papers on ASD were found (30 sFC, 7 dFC, 2 multimodal). 12 sFC, 5 dFC and 2 multimodal studies performed model explainability. For conciseness, we only summarise studies with larger datasets (more than 500 subjects) and elaborate on the others in the supplementary materials. Thus, 7 sFC, 3 dFC and 2 multimodal studies will be discussed here.

Li *et al.* [115] proposed a functional graph discriminative network (FGDN) which consists of 5 layers: a functional graph construction layer, two graph convolutional layers (Deferrard, BG), fully connected layer, and one output layer. The functional graph construction layer involves building a template for each class, which are subsequently used as the

graphs for all scans. The template is created by computing the mean FC matrix for the class and then further processing it via k-NN with Gaussian kernel as the distance function. Each test sample will then have an ASD graph and a TDC graph. Both graphs have the same node features (based on the FC matrix of the test sample) but different adjacency matrices (based on the templates). Both graphs are separately used as input to the remaining layers of the architecture. If the output from using the ASD graph is larger than that of the TDC graph, it is categorised as ASD. Using 3 atlases (AAL116, HO118 and MODL128) on the ABIDE dataset with 866 subjects (402 ASD, 464 TDC), the best performance was obtained by FGDN when the MODL128 atlas was used with tangent space embedding to calculate FC features (71.8% accuracy, while correlation-based feature obtained 68.0%). Model explainability was done by a less routine manner: given a target ROI, the values in all other ROIs were set to 0 and the model is **finetuned**. The finetuned models with higher accuracies were determined to be more important. They identified 5 most discriminative ROIs and highlighted that the

STG is the most salient ROI.

Zhang *et al.* [58] proposed a local-to-global GNN (LG-GNN) that focused on both analyzing disease-related local brain regions and biomarkers (BG, Kipf) and modeling both imaging and non-imaging subject-level relationship between the subjects (PG, Defferrard). They model a brain (local) graph where they learn the feature embeddings of local brain regions and identify biomarkers by introducing the self-attention based pooling mechanism. The adjacency matrix is modeled as Pearson's Correlation, and the node features are given as constants. The population (global) graph takes into consideration the embeddings from the previous layers (X for the population graph) and non-imaging information (used to compute A along with embeddings) to generate multi-scale feature embeddings corresponding to each subject. Using the HO112 atlas on the ABIDE dataset with 871 subjects (C-PAC), LG-GNN obtained an accuracy of 81.8% (GCN: 67.8%, population GCN: 71.4%) The **self-attention** mechanism from the BG helps in identifying biomarkers for classification. Salient regions identified include IFG, precentral gyrus, frontal orbital cortex and PARAH.

To address the issue of large variations in model performance when different atlases are used, Wang *et al.* [116] proposed a multi-atlas graph convolutional network method (MAGCN) that involves a stacking ensemble learning method. The GCN (Kipf, PG) involves node features based on the vectorised FC matrix after RFE (retaining 2000 features), while edge features follow Parisot *et al.* [24] (gender and site). The embeddings from GCN are passed through a softmax layer to make a preliminary decision and these probabilities are concatenated together (across atlases) and passed to a ridge classifier to produce the final prediction. 6 atlases were used: AAL116, EZ115, HO110, TT93, CC200, DOS160). On the ABIDE dataset with 949 subjects (419 ASD, 530 TDC), MAGCN achieves an accuracy of 75.9% (MLP: 75.2%). Model explainability was done but the details of the technique was not clear. Nevertheless, brain regions highlighted include angular gyrus, precentral gyrus, precuneus and thalamus.

Zhu *et al.* [117] proposed a triple-pooling GNN (TPGNN) which comprises three branches to capture multi-scale patterns (global, community, ROI). Each branch has its own GraphSAGE layers (BG), where the original node vector (connection profile, but negative values were set to 0) is concatenated with the updated vector learnt via weighted aggregation of its neighbours. The global branch involves multiple graph convolution layers for local aggregation before a global max-pooling operation. The community branch uses two hierarchical pooling layers based on soft clustering graph pooling. The ROI branch uses two top-K pooling layers before passing them through a graph convolution layer. Finally, the three views are integrated via concatenation, followed by two fully connected layers. Using the AAL116 atlas on the ABIDE dataset with 1112 subjects (splits and pipeline unclear), TPGNN obtained an accuracy of 72.5% (population GCN: 69.7%). Ablation studies showed that the ROI branch contributed the most to model performance while the community branch does not have much impact. For model interpretability, each branch has own biomarkers from its **pooling** layer. Regions highlighted

included lateral PFC, lateral dorsal PFC and SPL.

Shao *et al.* [118] utilised GCNs (Defferrard, PG) for ASD classification. For node features, deep feature selection (DFS) was performed. Instead of using vanilla DNN [24], they proposed to use a sparse 1-1 layer between the input and first hidden layer. Sparsity is introduced via regularisation similar to elastic net. Using the HO111 atlas on the ABIDE dataset with 871 subjects (403 ASD, 468 TDC), they obtained an accuracy of 79.5% (MLP: 78.1%). Comparing various feature selection strategies (DFS, RFE, Fisher), they found that DFS outperforms RFE and Fisher. They found that DL models are more compatible than ML Models for datasets with high dimensionalities. For model explainability, they assumed that features with higher weights are more important and they computed the absolute sum of the weights' magnitude across cross-validation folds. Based on the top 30 FC features from **DFS**, they found that they are evenly distributed across the brain (left, right and inter hemisphere). ASD has lower connectivity in 25 of these features (e.g. connection between right cuneal cortex and right parietal operculum cortex) and higher connectivity in the other 5 (e.g. connection between right thalamus and the anterior division of the right ITG).

Wang *et al.* [51] presented a connectivity-based GCN (cGCN), which uses a k-NN graph from the groupwise FC matrix for the graph convolution layer. EdgeConv was used for the GCN (BG) to compute the maximal activation from a node's top-k neighbours. The input to EdgeConv is the concatenation of the node's connection profile and the difference between the node's connection profile with its neighbour. Using the CC200 atlas on the ABIDE dataset with 1057 subjects (525 ASD, 532 TDC), cGCN obtained an accuracy of 70.7% with $k = 3$ (DNN: 70%). Model explainability was performed via **occlusion**. Salient networks identified include DMN, FPN and VAN. These findings were consistent across different values of k .

Li *et al.* [119] proposed an ensemble framework called TE-HI-GCN that used HI-GCN (see supplementary materials under ASD for more details) in a transfer learning framework. Transfer learning was introduced by having the target domain to go through another HI-GCN (with parameter sharing), i.e. T-HI-GCN. Furthermore, they use connectivity matrices with different levels of weight thresholds and sparsity as the graph and each matrix has its own branch of T-HI-GCN. These models were trained separately and combined via majority voting, in an ensemble framework. Various atlases were used, including AAL116, TT97, HO110, EZ115, CC200 and DOS160. On the ABIDE dataset with 871 subjects (403 ASD, 468 TDC), they obtained an accuracy of 76.5% (GCN: 59.8%, HI-GCN: 67.2%). For model explainability, they aggregated **cluster membership scores** obtained during the clustering process (f-GCN) to reveal important clusters and inter-cluster connections. CEN, DMN and SN were to be crucial for ASD. For inter-modular connections, they found CEN-SMM, CEN-SN, DMN-SMM, DMN-CEN, and DMN-Visual to be crucial.

A dFC study by Zhu *et al.* [120] raised the issue that many existing approaches are limited to low-level correlations. They proposed a contrastive multi-view composite GCN (CMV-CGCN) to address this. CMV-CGCN integrates the attention

mechanism, multi-view contribution learning and contrastive multi-view learning. It comprises two channels: one for low-order FC (LOFC) and another for high-order FC (HOFC). HOFC is computed by calculating the Pearson correlation between every pair of rows in the original FC matrix, producing a HOFC matrix with the same dimensionality. In each channel, node features are obtained via RFE to select 2000 FCs from the initial sFC matrix vector. Edges are constructed by incorporating gender and site information, similar to Parisot *et al.* [24] (PG). Each channel contains a node-level graph encoder that generates node embeddings using composite GCNs (Kipf, PG) made up of three layers: graph convolution, graph attention and graph convolution.

Graph contribution learning is introduced to the outputs of the graph attention layer. This module encourages nodes from each channel to have different contributions. It does so by computing weights for each node via a softmax-like expression with learnable weight vectors (separate for LOFC and HOFC) multiplied to a concatenation of the node LOFC and HOFC embeddings. These weights are used to compute a global node representation which is then optimised via a cross-entropy loss function. Contrastive learning is introduced to the outputs of the last graph convolution layer to encourage consistent representation of two views for the same subject, while encouraging representations of different subjects to be mapped to distant points in the embedding space. Both contribution learning and contrastive learning are implemented by introducing additional terms to the loss function (on top of the usual cross entropy loss). The learnt embeddings are then combined and fed into the downstream network for classification. Using the HO110 atlas on a subset of the ABIDE dataset with 613 subjects (286 ASD, 327 TDC), CMV-CGCN yielded an accuracy of 75.2% (MV-GCN: 72.0%) For model explainability, the frequency of each brain region involved in the 2000 FCs selected via **RFE** was computed, and 20 most salient regions were selected. This included the lingual gyrus, MFG and SFG.

Cui *et al.* [121] proposed a dual graph based dynamic multi-site GCN (DG-DMSGCN). It comprises a sliding window dual GCN (SW-DGCN) which is formed by two GCNs: one for a common graph (constructed via one-sample t-test) and another for a diversity graph (constructed via two-sample t-test between ASD and TDC). Each GCN has 3 EdgeConv layers (BG). The outputs of the GCN are concatenated and passed through a 1x1 convolution before passing it to a feature fusion block which computes the mean and variance of the representations across sliding windows. This produces a subject feature vector as the output of SW-DGCN. The SW-DGCN is followed by DMSGCN, which encourages representations to be similar for subjects with the same label or same/similar sites, while representations for different classes should be different. This is realised via three evaluation metrics for feature, site and label. Feature similarity follows Parisot *et al.*'s formulation [24], while the other 2 metrics introduce a learnable parameter that is assigned depending on site and label similarity. The scores from these metrics are multiplied together to form the population adjacency matrix (PG) that is used by another GCN. The learnt representations are then passed into a classification block implemented via MLP. Using

the CC200 atlas on the ABIDE dataset with 1035 subjects (505 ASD, 530 TDC), DG-DMSGCN obtained an accuracy of 73.1% (SVM: 64.0%). Model explainability was performed via **saliency** (gradients) and revealed salient regions such as FP, precuneus cortex, brain-stem and paracingulate gyrus.

Chen *et al.* [122] proposed an invertible dynamic GCN (ID-GCN), which involves an invertible block made up of 2 GCNs (Defferrard): one to capture spatial information and another for temporal information. While the functional graph is the typical BG formed from the FC matrix, the spatial graph is constructed based on the correlation of a distance vector constructed for each ROI, which measures the spatial distance between ROI pairs. However, these original graph are further processed to form their k-NN graph. For both graphs, features selected from the connection profile is used as node features. Feature selection is done via random forests and this incorporates features from both sFC and dFC (sliding window). The invertible block involves a series of summation operations between two inputs, such that it is possible to recover the original inputs despite the non-linearities (GCN) used in the summation. Using the HO110 atlas on a subset of the ABIDE dataset with 867 subjects (416 ASD, 451 TDC, 13 sites, C-PAC), ID-GCN achieved an accuracy of 76.3% (GCN: 73.2%). Model explainability was done based on the features selected by **random forest**. This includes the connections between the right pallidum and right IFG, left frontal orbital cortex and left central opercular cortex, and connections involving the left supramarginal gyrus and right ITG.

Two multimodal studies were performed by Chen *et al.* on ASD. In [123], they proposed an attention-based node-edge GCN (ANEGCN) to model multimodal datasets. The node vector contains 3 features: 2 from ALFF (amplitude of low frequency fluctuations, from fMRI) and 1 feature based on average T1w intensity at that ROI. Edges are initialised as the FC matrix. However, unlike most GCNs that focus on updating node features, ANEGCN also updates edge features simultaneously. ANEGCN comprises node edge graph convolution (NE-GC) blocks that utilises attention to aggregate information only from neighbors with similar node characteristics and edge characteristics, so as to avoid over-smoothing. After the last block, the node feature maps and edge feature maps are concatenated together before another MLP. The loss function involves cross-entropy with label smoothing, L2 regularisation and also an adversarial loss term (fast gradient sign method). Using the AAL116 atlas on the ABIDE dataset with 1007 subjects (481 ASD, 526 TDC), they obtained an accuracy of 72.7% (GCN: 70.4%). For model interpretability, they directly computed the **gradients** with respect to the input. These saliency scores are normalised and averaged across the class to obtain group-wise saliency maps. The analysis revealed increased T1w intensities in DMN, reward, memory and motor function areas, as well as higher ALFF values in regions associated with reward and motor function. Edge features revealed the dominance of inter-network connections relative to intra-network ones as well as lowered homotopic interhemispheric connections in the limbic regions. In an extension of this work, Chen *et al.* [81] further proposed an adversarial learning-based

node–edge GAT (AL-NEGAT) architecture. AL is introduced to improve generalisability via the fast gradient sign method and projected gradient descent. The rest of the architecture is similar to ANEGCN. Using the same dataset as above, they obtained a slightly higher accuracy of 74.7%. For model explainability, they also computed **gradients** with respect to nodes and edges. This revealed that ASD subjects have lower FC between the right middle temporal gyrus and multiple ROIs in the frontal, parietal and occipital lobes, as well as a mix of higher and lower FCs between ROIs in the limbic regions to multiple other regions in the brain.

Overall, ASD studies have relatively larger dataset sizes (mean: 698) due to the availability of the ABIDE dataset (used in 36 out of 39 studies). Model performance is moderately high (73.3%). Both BG (24 studies) and PG (11 studies) are rather well investigated and 4 studies used a mix of them. The best performing models utilised techniques such as adaptive PG construction [124] and combining BG/PG [58], [119]. Put together, future studies should consider evaluating the generalisability of these models to datasets other than ABIDE.

Out of 12 sFC studies that incorporated explainers, 8 studies analysed features at the level of ROIs. 4 studies highlighted the prefrontal regions (prefrontal cortex, frontal gyrus), 3 studies found thalamus to be salient and 2 studies mentioned parahippocampal gyrus and precentral gyrus. 3 of the sFC studies involved modular-level analysis, of which 2 pointed out common salient networks (DMN, CEN/PFN). From the 5 dFC studies, 3 ROI-level studies have reproducible potential biomarkers mentioned in 2/3 studies: lingual gyrus, superior frontal gyrus, frontal pole, insular cortex. However, FC-level analysis did not reveal any common salient features. The 2 multimodal studies (from the same first author) are noted to have very similar experiment setup, but very few salient features were replicated (e.g. T1w intensity at the left rolandic operculum was found to be a salient feature in both studies).

Summarising across the 19 studies that analysed salient features, one salient region found to be reproducible across multiple sFC (4) and dFC (2) studies was the prefrontal cortex. However, this is still a rather broad area and more thorough studies are needed in the future to narrow this down, as well as to study connection-level features related to this region. In such studies, it would be prudent to adopt GNN architectures that have been demonstrated to perform well (instead of baseline GNNs), such as those proposed by Zhang *et al.* [58] (BG/PG mix) and Mao *et al.* [124] (adaptive PG).

C. Major Depressive Disorder

14 studies on MDD were found (9 sFC, 4 dFC and 1 multimodal). 4 sFC studies and 2 dFC studies highlighted salient features. Notably, some studies analyse first-episode drug-naive (FEDN) patients and recurrent MDD patients separately. Qin *et al.* [82] applied a 3-layer GCN (Defferrard, BG) for MDD classification. The graph produced via k-NN algorithm on FC matrix was used for the GCN and connection profile was used as node features. Using the DOS160 atlas on the REST-meta-MDD dataset with 1586 subjects (821 MDD, 765 CN; across 16 sites, harmonised via ComBat,

they achieved an accuracy of 81.5%. Notably, they found that other classification tasks had lower performance (74.1% for FEDN vs NC, 78.1% for recurrent vs NC and 70.9% for FEDN vs recurrent). For MDD vs NC classification, the most salient regions revealed via **CAM** were found within the DMN, FP and Cingulo-Opercular networks. Specific brain regions include the right dorsal anterior cingulate cortex, right ventrolateral prefrontal cortex, left inferior parietal lobule (found to be associated with depressive severity) and left posterior insula.

Gallo *et al.* [30] utilises a vanilla 2-layer GCN (Kipf, BG) followed by an average pooling layer for their architecture. The adjacency matrix was thresholded to retain 50% of edges and then binarised, while the node features contain the connection profile. Using the HO112 atlas on two of the largest consortia (REST-meta-MDD and psymri), amounting to 2498 subjects (1249 MDD, 1249 NC after downsampling to balance the classes), an accuracy of 61.3% was attained (SVM-RBF: 61.2%). Notably, they perform separate experiments for non-medicated and medicated patients and also noted that using CC200 atlas gave comparable performance. Multiple techniques of model explainability were conducted (**GNNExplainer**, **ablation study** and **univariate t-test**). GNNExplainer revealed that stronger inter-hemispheric thalamic connection is a discriminative feature across both datasets, while region-level ablation study did not show such consistency (but the left thalamus was the second most salient region in the psymri dataset and the right thalamus was the most salient region in REST-meta-MDD). However, t-test results only showed thalamic hyperconnectivity in REST-meta-MDD but not for psymri.

Kong *et al.* [125] proposed a multi-stage graph fusion network (MSGFN) which introduces a graph construction module. Using a multi-layer autoencoder to learn latent representations of the inputs, they construct graphs from the output of each layer (including the input layer). Each graph is then used as input to a separate branch of GCN (Kipf, BG) and the subsequent outputs are then combined by taking their average. Adjacency matrices were derived from the binarization of the affinity matrices produced from the autoencoder and node features are the connection profile from FC. Using two atlases (gray matter features from the BM82 atlas and white matter features from the JHU81 atlas) on a private dataset of 218 subjects (129 MDD, 89 NC), an accuracy of 70.9% was achieved. Salient features were identified via averaging of the relevant **weights** (details not specified). They highlighted 10 WM-GM connections, with the connection between the right anterior corona radiata and the left dorsolateral PFC having the highest score.

Jun *et al.* [126] uses a spectral GCN (Defferrard, PG) with phenotypic information like age and gender used to construct the graph. One difference from other studies is the use of both effective connectivity (EC) and FC features as node vector, with feature selection via lasso. EC uses group sparse representation derived from structured equation modelling. Using the Yeo114 atlas on a private dataset of 75 subjects (29 MDD, 44 NC), they achieved an accuracy of 74.1% (SVM: 69.8%) when EC was used and 56.4%

TABLE V

SUMMARY OF FINDINGS FROM STUDIES THAT IDENTIFIED POTENTIAL BIOMARKERS. ‘SIZE’ REFERS TO SIZE OF DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (dFC) OR [M] (MULTIMODAL). ‘TYPE’ REFERS TO THE TYPE OF EXPLANATION.

Reference	Dataset (Size)	Atlas	Explainer (Type)	Salient features
MDD				
[82]	REST-meta-MDD (1586)	DOS160	CAM (ROI, module)	Regions: right dorsal ACC, right ventrolateral PFC, left IPL, left posterior insula. Modules: DMN, FPN, Cingulo-Opercular network.
[30]	Multiple (2498)	Multiple	Multiple (connection)	Connections: stronger inter-hemispheric thalamic connection across most datasets and explainers
[125]	Private (218)	Multiple	Weights (connection)	Connections: right anterior corona radiata - left dorsolateral PFC
[126]	Private (75)	YEO114	Gradients (connection)	Connections (reduced in EC): left dorsal PFC - left precentral ventral region, striate cortex, parietal medial region, IPL, PARAH
[127]	Private (277) [d]	Unclear	Unclear (ROI)	Regions: bilateral pallidum, right putamen, bilateral MFG, right postcentral gyrus, right Heschl gyrus, right caudate, right olfactory cortex, right IFG, triangular part.
[128]	REST-meta-MDD (681) [d]	AAL116	Attention (connection, temporal)	Connections: cross-hemisphere connections within the insula and lingual gyrus ; lingual gyrus - calcarine sulcus ; Temporal: middle & end of time series
SZ				
[129]	Multiple (1412)	Multiple	CAM (ROI)	Regions: decreased nodal efficiency in bilateral putamen and pallidum in SZ, across both atlases
[131]	Private (345) [M]	Multiple	Pooling (ROI)	Regions: bilateral rectus gyrus, bilateral lingual gyrus, bilateral cuneus; right medial orbitofrontal cortex, medial SFG, calcarine cortex, anterior cingulate gyrus
[133]	COBRE (154) [M]	DK293	Pooling (ROI)	Regions: ROI relevance scores had a stronger correlation with SC than FC

(SVM: 60.3%) when FC was used. For model explainability, they applied **sensitivity analysis** for EC features, showing reduced connectivity between the left dorsal PFC with the left precentral ventral region, striate cortex, parietal medial region, inferior parietal lobule and parahippocampal cortex.

A dFC study by Kong *et al.* [127] introduced a novel spatio-temporal GCN (STGCN) which comprises spatial graph attention convolution (SGAC) and temporal fusion. In SGAC, a linear filter is first used to balance between an identity matrix with an adjacency matrix built from a time window (BG). fMRI time series was used as the node features. These are passed through a GAT layer, followed by hierarchical anatomy-guided pooling using 3 brain parcellations (90, 54, 14 ROIs). This is implemented by learning an appropriately-shaped matrix (based on the desired number of nodes) that linearly transforms the representation to a lower dimensionality. Finally, a temporal fusion module based on long short-term memory (LSTM) was proposed to capture the dFC features. Outputs of the temporal fusion module are concatenated before making the final prediction. Using the AAL90 atlas on a private dataset of 277 subjects (180 MDD, 97 NC), they obtained an accuracy of ~ 84%. While the most discriminating regions were identified, the method used is unclear. ROIs highlighted include pallidum and MFG.

Fang *et al.* [128] proposed an unsupervised cross-domain fMRI adaptation framework (UFA-Net) to address the issue of inter-site heterogeneity. Using fMRI time series as features, data from both the source site and target site are put through an attention-guided spatio-temporal graph convolution module and 3 convolution layers (Kipf, BG). The outputs of the convolution layers for both source and target data are passed to a maximum mean discrepancy constrained module to perform

cross-site feature alignment. Using the AAL116 atlas on 681 subjects (356 MDD, 325 CN) from two imaging sites (Site 20 being the source, Site 1 being the target) of the REST-meta-MDD consortium, their model achieves an accuracy of 59.7%. **Attention** scores revealed that cross-hemisphere connections within the insula and lingual gyrus, as well as connections between the lingual gyrus and calcarine sulcus were the most discriminative connections. A spatio-temporal feature map was also presented, where they showed that a large inter-group difference is present in the middle and end of the scan period.

Overall, classification accuracies for MDD (mean: 71.0%, after excluding 3 studies with unusually high accuracies) are generally lower than other disorders. Notably, some MDD studies have much larger datasets than other disorders (over 1000 subjects) due to the availability of the REST-meta-MDD dataset. As most studies used BG (12/14), future studies could consider exploring PG-based approaches.

Interestingly, both Gallo *et al.* [30] and Fang *et al.* [128] reported a common salient connection (right lingual gyrus - right calcarine cortex/sulcus) despite using different atlases and performing a different analysis (sFC and dFC respectively). However, both studies have rather low classification accuracies (around 60%) on the REST-meta-MDD dataset. Notably, Qin *et al.* [82] achieved a much higher accuracy on the same dataset, with a key difference being the choice of adjacency matrix (Qin *et al.* used a k-NN graph, while Gallo *et al.* used a binarised thresholded matrix). However, Qin *et al.* only highlighted salient ROIs and not connections. Future studies on MDD could adopt Qin *et al.*'s modelling approach and further verify the robustness of this salient connection.

D. Schizophrenia

7 papers on SZ were found (3 sFC, 1 dFC, 3 multimodal), out of which 3 studies identified salient features (1 sFC and 2 multimodal). The sFC study used vanilla GNNs on a large ($n > 1000$) multi-site dataset, while the multimodal studies used both T1w and fMRI data on smaller datasets.

For the study on sFC, Lei *et al.* [129] used a vanilla GNN architecture with 3 ChebConv layers for SZ classification. Graph used for the GCN was formed via a k-NN approach (BG) based on the Euclidean distance between the ROIs and retaining the 10 nearest neighbors. Connection profile was used as node features. Using the AAL90 atlas (also replicated using the DOS160 atlas) on a multi-site dataset with 1412 subjects (505 SZ, 907 NC, 6 sites, harmonised via ComBat [130]), they achieved a balanced accuracy ranging from 65.7% to 79.2% on individual sites (SVM: 67.2% to 75.6%) and 85.8% on a pooled and harmonised dataset (SVM: 80.9%). Model interpretability was performed using CAM. Analysis of the top 10 most salient ROIs revealed significantly decreased nodal efficiency in bilateral putamen and pallidum in SZ patients, which was further verified to be significantly associated with negative symptom scores. This finding was consistent across the two brain atlases used in this study.

Chen *et al.* [131] extended their earlier fMRI study [132] (see supplementary materials for details) to consider multimodal datasets. They used GCNs (Kipf, BG) with 18 different configurations of selecting the node features and graph (3 groups of node features, 3 types of edge features, and 2 brain atlases). The GCNs were similar to their earlier study (i.e. top-k pooling, multi-scale readout). 3 groups of node features used for each ROI include: (i) structural measures (average GM volume, average WM volume and structural degree centrality), (ii) functional measures (ReHo (Regional Homogeneity), ALFF and degree centrality), (iii) multimodal measures (combination of structural and functional measures). 3 types of edge features include: (i) structural brain network (KL divergence of GM volume map), (ii) functional brain network (absolute Pearson correlation), (iii) multimodal brain network formed by adding structural and functional brain networks. These networks were further sparsified via thresholding and binarisation. Using the AAL90 atlas and BNA246 atlas on a private dataset with 345 Han Chinese subjects (140 SZ, 205 NC), they obtained an accuracy of 95.8% when a combination of multimodal measures and functional brain network was used with the AAL90 atlas (SVM: 81.2%). Top 10 salient ROIs were identified via the last **pooling** layer and ranked based on the frequency of occurrence (within top 10) across the dataset. These are found to be dominated by the prefrontal and occipital cortices. Specifically, ROIs involved are: bilateral rectus gyrus, bilateral lingual gyrus, bilateral cuneus as well as the right medial orbitofrontal cortex, medial SFG, calcarine cortex and anterior cingulate gyrus.

A multimodal study by Sebenius *et al.* [133] proposed multimodal GNN (MM-GNN) which extends top-k pooling and SAGpool to allow for multimodal pooling across morphometric similarity networks (from T1w) and FC matrices. The graph used for the GNN (Kipf) was constructed via binarised

BG with 5%-10% sparsity, while node features were based on 10 handcrafted features for each modality (structural measures such as cortical thickness and network features such as node degree). Using an atlas with 308 ROIs (293 ROIs after quality checks) on the COBRE dataset with 154 subjects (67 SZ, 87 NC), MM-GNN achieved an accuracy of 75% (SVM-RBF: 71%) when a pooling ratio of 0.6 was used (beyond which, too much information is lost). Interpretability is enabled from ROI **pooling** scores, but salient ROIs were not discussed in detail. Instead, it was reported that the ROI relevance scores had a stronger correlation with structural connectivity than FC.

Overall, model performances for SZ studies are high (mean: 88.2%) but this comes with the caveat of small dataset sizes (mean: 230, excluding one study with more than 1000 subjects). However, it is notable that the study with large dataset reported a high balanced accuracy of 85.8%. Majority (6/7) of SZ studies have adopted a BG approach and there remains much room to explore for PG-based approaches.

Salient regions highlighted by the two studies discussed were disparate despite the use of COBRE in both studies. This could possibly be explained by how the study by Lei *et al.* [129] included multiple sites (COBRE only represents 10% of the dataset they used), as well as their use of ComBat to remove site effect, which might have removed site-specific information that were highlighted in the study by Chen *et al.* [131]. Given the high classification performance, future studies should go further into model explainability and identify salient features of SZ.

E. Dementia

29 studies on dementia were found (10 sFC, 14 dFC, 5 multimodal). A subset of these (5 sFC, 7 dFC and 3 multimodal) reported salient features. However, these studies span various stages and subtypes of dementia, with some categories only having 1 or 2 studies performed at incomparable degrees of granularity (including studies on subjective cognitive decline, MCI and non-AD dementia). These studies are discussed in the supplementary materials, while categories with sufficient studies such as significant memory concern (SMC), early MCI (EMCI), late MCI (LMCI) and AD are discussed here. Studies involving multiple stages of AD are discussed thereafter. This results in 11 studies (2 sFC, 6 dFC and 3 multimodal) discussed in this section.

A study on SMC by Zuo *et al.* [134] uses a vanilla GCN classifier to perform disease classification, but they also proposed an Adversarial Temporal-Spatial Aligned Transformer (ATAT) architecture that maps fMRI time series to FC matrices. While ATAT learns the ROIs boundaries, it is influenced by the original FC matrices via a discriminator (and its adversarial loss function). These synthesised matrices are used alongside the original FC matrices. Using the AAL90 atlas on the ADNI3 dataset with 168 subjects (82 SMC, 86 NC), they obtained around 87.5% when either SVM or GCN were used. Notably, using original FC produces for classification leads to around 80% for SVM and around 84% accuracy for GCN. For model explainability, **occlusion** was performed. This highlighted 8 ROIs found across all AD stages: the

TABLE VI

SUMMARY OF FINDINGS FROM STUDIES THAT IDENTIFIED POTENTIAL BIOMARKERS. ‘SIZE’ REFERS TO SIZE OF DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL). ‘TYPE’ REFERS TO THE TYPE OF EXPLANATION.

Reference	Dataset (Size)	Atlas	Explainer (Type)	Salient features
Dementia (SMC)				
[134]	ADNI (168) [d]	AAL90	Occlusion (ROI)	Regions: orbital part of the SFG, anterior cingulate, paracingulate gyri, calcarine fissure, lingual gyrus, precuneus, paracentral lobule, caudate nucleus, lenticular nucleus putamen
[141]	ADNI (138) [d]	AAL90	Unclear (ROI)	Regions: ITG, MFG, IFG, left hippocampus
[143]	Multiple (207) [M]	AAL90	Unclear (ROI)	Regions: right ITG, right insula, left olfactory cortex, left angular gyrus, right amygdala, right precuneus
[26]	ADNI (88) [M]	AAL90	RFE (connection)	Connections: function: precentral gyrus (left and right) - left superior TP ; structure: left precentral gyrus - left putamen
Dementia (EMCI)				
[61]	ADNI (910)	SF200	Pooling (module)	Modules: DAN, VAN, DMN
[135]	ADNI (101) [d]	YEO114	RL (ROI)	Regions: ROIs in DMN, including bilateral parahippocampal cortices; right insula, involving both SMN and SN/VAN
[140]	ADNI (88) [d]	AAL90	Weights (ROI)	Regions: medial part of left SFG, right putamen
[141]	ADNI (180) [d]	AAL90	Unclear (ROI)	Regions: ITG, MFG, left hippocampus and IFG
[142]	ADNI (154) [M]	AAL90	Occlusion (ROI)	Regions: IFG, olfactory cortex, PARAH, MOG, ITG
[143]	Multiple (249) [M]	AAL90	Unclear (ROI)	Regions: right ITG, right insula, left olfactory cortex, left angular gyrus, right amygdala, right precuneus
[26]	ADNI (88) [M]	AAL90	RFE (connection)	Connections: function: right precentral gyrus - right putamen ; structure: left superior frontal lobe - left thalamus
Dementia (LMCI)				
[140]	ADNI (82) [d]	AAL90	Weights (ROI)	Regions: medial part of left superior frontal gyrus and right putamen
[142]	ADNI (107) [M]	AAL90	Occlusion (ROI)	Regions: IFG, olfactory cortex, PARAH, MOG, ITG
[143]	Multiple (329) [M]	AAL90	Unclear (ROI)	Regions: right ITG, right insula, left olfactory cortex, left angular gyrus, right amygdala, right precuneus
[26]	ADNI (82) [M]	AAL90	RFE (connection)	Connections: function: orbital part of the MFG/SFG - right ITG ; structure: left superior frontal lobe - left thalamus
Dementia (AD)				
[136]	ADNI (83)	AAL90	Weights (ROI)	Regions: left MFG, left orbital SFG, right precentral gyrus
[138]	ADNI (292) [d]	AAL116	GradCAM (ROI)	Regions: bilateral hippocampus, right precuneus, right frontal mid-cortex, left precentral cortex
[139]	ADNI (107) [d]	AAL90	Attention (ROI)	Regions: frontal and temporal regions
PD				
[144]	Private (150) [MM]	DK86	Multiple (ROI)	Regions: cerebellum, precuneus, pallidum
[146]	PPMI (41) [MM]	BNA246	logreg (ROI)	Regions: precentral gyrus, postcentral gyrus, parietal lobe

orbital part of the SFG, anterior cingulate and paracingulate gyri, calcarine fissure, lingual gyrus, the precuneus, paracentral lobule, caudate nucleus and lenticular nucleus putamen.

A study on EMCI by Mei *et al.* [61] proposed a hierarchical functional brain network encoder (HFBN-GCN) that involves two branches to address the over-smoothing problem: one branch uses the whole topology while another hierarchical pooling module. Each branch has 4 GCN layers and the representations from the last GCN layer are added together and passed to a MLP to produce the final diagnosis. The hierarchical pooling module uses 4 layers of GCN (Kipf, BG) for a coarse to fine readout operation: node level, hemi-

spherical level, network level and whole brain level. Principal components of the FC matrix produced via PCA are used as node features. Using the SF200 atlas (as well as the Yeo-7 and Yeo-17 atlas for pooling) on the ADNI dataset with 910 scans (345 EMCI, 565 NC), they achieved an accuracy of 73.4% for EMCI classification (GCN: 66.9%). They detected biomarkers by extracting **weights** on each functional sub-network from the final pooling layers of the best model (HFBN-GCN with 7-sub-network). DAN, VAN and DMN were found to be the top 3 contributing sub-networks.

Lee *et al.* [135] hypothesised that each subject should have a different set of ROIs and proposed a set-input neural

network architecture that is able to automatically select ROIs for each subject from dFC data. This involves a temporal embedding module (3 layers of 1D CNN), an ROI selection module (an ROI-wise 1D CNN, trained via model-free RL) and a disease identification module (2 GCN layers (Kipf, BG)). Since different subjects have different subsets of ROIs, this necessitates a set-input neural network which in turn requires permutation-equivariance in the inputs and permutation-invariance in the outputs. This is achieved by using a CNN before the ROI selection module and a GCN after the ROI selection module. Optimisation of the weights in the proposed architecture involves pre-training the temporal embedding and disease identification modules, before joint training of the whole architecture. The REINFORCE algorithm is used to train the ROI selection network. Using the Yeo114 atlas on 101 subjects from ADNI2 and ADNI GO (53 EMCI, 48 NC), an accuracy of 74.4% was achieved (BrainNetCNN: 67.6%). Ablation studies clearly showed that the ROI selection module has the greatest contribution to model performance. Individualised biomarkers were produced via **RL** in the ROI selection module. For ease of discussion, group-level comparisons are summarised here. Regions selected are dominated by ROIs in DMN, including bilateral parahippocampal cortices. The right insula, involving both SMN and SN/VAN networks, was also highlighted.

3 studies were solely focused on AD. Alorf *et al.* [136] adopted the brain connectivity-based GCN architecture (BC-GCN) proposed in [137] and used it for AD stage prediction. BC-GCN entails an edge-based graph convolution approach with edge pooling and node pooling. Using the AAL90 atlas on the ADNI dataset with 83 subjects (33 AD, 50 CN), they obtained an accuracy of 94.2%. Explainability via model **weights** highlighted left MFG, left orbital SFG and right precentral gyrus to be the most salient regions.

Xing *et al.* [138] uses a sliding window approach and an LSTM with GCN layers before it to model dFC data. They also introduce the concept of ‘Assistant Task Training’, which involves the model having additional branches to predict demographics like age and gender (besides disease class). Combination could be done via hard parameter sharing (each branch share the same parameters for the front layers) or soft parameter sharing (each branch have their own parameters and the feature maps are fused via a linear combination). BG approach is adopted and node features used are the volume of the corresponding ROI. Although they used the GCN from Defferrard *et al.*, they opined that $K = 1$ is sufficient since brain networks are densely connected. Using the AAL116 atlas on the ADNI2 dataset with 292 subjects (118 AD, 174 NC), they obtained an accuracy of 90% for AD/NC classification when soft parameter sharing is used (static FC, GCN: 81.3%). Model explainability is performed via **Grad-CAM**, revealing that ROIs in the bilateral hippocampus, right precuneus, right frontal mid-cortex and left precentral cortex have highest activation when predicting AD.

Wang *et al.* [139] proposed a GNN architecture that incorporates self-attention to adaptively learn the adjacency matrix used by the GCN. This architecture comprises a feature selection module extract representations from dFC sliding windows,

as well as a self attention model. In the feature selection module, the connection profile from each sliding window is passed through an MLP before being used in a GCN layer (Kipf). Graph used by the GCN is originally the FC matrix. Representations from the GCN are then concatenated across windows before being passed through another fully connected layer and a softmax layer. Embeddings from the GCN are also passed to a multi-head self-attention layer. The attention matrix is then combined with the original graph via dot product. This new matrix is then used for the GCN. Using the AAL90 atlas on the ADNI dataset with 107 subjects (59 AD, 48 NC), they obtained an accuracy of 89.8% (GCN: 87.5%). Model explainability was achieved via **attention** scores. For ADNI, salient regions include the amygdala, precentral gyrus and parahippocampal gyrus.

Several studies considered multiple stages of early AD (i.e. SMC, EMCI, LMCI). Yu *et al.* [140] built a dynamic HOFc, but extracted the local weighted clustering coefficients to use as node features. The graph was constructed by combining both image and non-image information. They proposed using GCN (Defferrard, PG) with an Inception module (spectral graph convolution with different kernel sizes and concatenation as aggregator). Using the AAL90 atlas on the ADNI dataset with 126 subjects (38 LMCI, 44 EMCI and 44 NC), they obtained an accuracy of 87.5% for EMCI vs NC (GCN: 79.6%) and 89.0% for LMCI vs NC (GCN: 84.2%). For model explainability, the authors selected the most important ROIs based on the **weight** of the brain network. They identified the medial part of the left SFG and the right putamen to be generally salient for both disease classes.

Zhu *et al.* [141] proposes a structure and feature based graph U-Net (SFG U-Net) architecture for classification of MCI and SMC patients using dFC data. Using a sliding window approach to construct dFC matrices, they are used to generate a PG based on similarity of dFC feature vectors (averaged across windows) as well as age and gender. This is passed to a U-Net architecture where the inputs go through multiple blocks of graph convolution (Kipf) and pooling operations. They opined that the top-k pooling operation in graph U-Net ignores the topological structure and proposed a novel adaptive pooling layer. This involves 2 branches: in branch 1, node vector goes through dense layer ; in branch 2, node vector goes through GCN w/ the usual PG. The outputs of both branches are combined via a weighted sum and used to do Top-K pooling, which is then applied on both the node vector and the graph. Since it is a graph U-Net, an unpooling layer is also introduced and achieved by recording the position of nodes during the pooling operation. Using the AAL90 atlas on the ADNI dataset with 291 subjects (44 SMC, 86 EMCI, 67 LCMI, 94 NC), they obtained an accuracy of 83.2% (GCN: 81.1) for NC vs SMC and 80.0% (GCN: 75.6%) for NC vs EMCI. Most discriminative regions for early AD (means of deriving importance scores were not clearly explained) were identified to be the ITG, MFG, left hippocampus and IFG.

A multimodal study by Lei *et al.* [142] proposes a multi-scale enhanced GCN (MSE-GCN) to tackle the over-smoothing problem encountered when using multiple GCNs (while still increasing the receptive field). MSE-GCN does

this by providing graph embedding information from random walks from the corresponding subject graph (PG). Multimodal information from both dFC (rs-fMRI) and structural connectivity (SC) from diffusion tensor imaging are first extracted via local weighted clustering coefficient and these feature vectors are encoded by using multiple GCN layers (Defferrard) in parallel before being used to construct the PG. Then, a multi-scale GCN setup is used by performing random walks of varying numbers of steps on the PG. Representations from various scales are fused via addition and passed through a fully connected layer for classification. Using the AAL90 dataset on the ADNI dataset with 184 subjects (40 LMCI, 77 EMCI, 67 NC), they obtained an accuracy of 85.4% for EMCI vs NC (GCN: 71.5%, MS-GCN: 81.3%), 93.5% for LMCI vs NC (GCN: 73.8%, MS-GCN: 85.1%). Model explainability was performed via **occlusion**, revealing that the following regions are salient across both classification tasks: IFG, olfactory cortex, PARAH, MOG and ITG.

Song *et al.* [143] raised the issues of confounders, heterogeneous protocols and small datasets as reasons for poor model performance. They introduced three mechanisms in the current GCN (PG) to address them: (i) a dual-modality fused connectivity network that introduces a penalty term based on SC to construct FC, (ii) multi-center attention graph that considers phenotype, scanner, site and disease status, (iii) a multi-channel mechanism that assigns different filters based on feature statistics, so as to address the problem of how existing GCNs apply the same convolutional coefficients on all features. They also proposed a pooling mechanism that considers the disease class during model training, which essentially does top-k pooling based on the difference between class-specific similarity matrices. Using the AAL90 atlas on 3 datasets (ADNI2, ADNI3, and an in-house dataset) with 459 subjects (44 SMC, 86 EMCI, 166 LMCI, 163 NC), they achieved an accuracy of 93.2% for NC vs SMC classification (GCN: 86.5%), 91.2% for NC vs EMCI (GCN: 85.5) and 94.2% (GCN: 87.5%) for NC vs LMCI. They identified the top 10 most discriminative features as obtained from the **dual-modality fused connectivity network**, highlighting regions such as the right ITG, right insula and left olfactory cortex.

Song *et al.* [26] proposed 3 mechanisms to improve GCNs (Defferrard, PG): (i) similarity-aware receptive fields in adjacency matrix, (ii) adaptive mechanism combined with pre-trained GCNs to score all subjects and (iii) calibration mechanism for fusing dual-modal information into the adjacency matrix. PG are first constructed separately using SC and FC matrices for node vectors (after RFE for feature selection). Edges are initialised with imaging and non-imaging data in a manner similar to Parisot *et al.* However, they introduce 3 receptive fields that consider disease status: (i) in training set, CN are connected to each other via [24], (ii) in training set, patients with disease are also connected to each other via [24], i.e. CN and patients are not connected, (iii) for test data, they are connected to everyone else as usual. During training, edges are then iteratively modified via a 2-stage ‘adaptive calibration mechanism’ that fuses structure and function. In Stage 1, they train GCNs on the original graphs (one GCN for each modality) which produce a score vector for all data samples.

In Stage 2, they used the difference in score vectors rather than imaging features, but everything else remained the same. Score vectors are the logits produced by the trained GCNs from Stage 1. Score vectors are used for the adaptive mechanism, replacing the correlation distance used for constructing A which they call a similarity-aware adaptive matrix. Finally, calibration mechanism for structure-function fusion involves using hadamard product of the SC and FC matrices, followed by normalisation. Using the AAL90 atlas on the ADNI dataset with 170 subjects (44 SMC, 44 EMCI, 38 LMCI, 44 NC), they achieved an accuracy of 84.1% (GCN: 76.1%) for SMC classification, 85.2% (GCN: 75.0%) for EMCI classification, 89.0% for LMCI classification (GCN: 80.7%). For model explainability, the most discriminative features identified via **RFE** were reported separately for both modalities. For SMC, precentral gyrus (left and right) - left superior TP were the most salient functional connections, while left precentral gyrus - left putamen was the most salient structural connection. For EMCI, right precentral gyrus - right putamen was the most salient functional connection, while left superior frontal lobe - left thalamus was the most salient structural connection. For LMCI, orbital part of the MFG/SFG - right ITG were the most salient functional connections. The most salient structural connection for LMCI was the same as EMCI. Other notable findings for functional connections include: inter-hemispheric connection between the thalamus in both SMC and EMCI ; left paracentral lobule - right caudate and right STG - right superior TP were in both EMCI and LMCI. Other notable findings for structural connections include: left SFG - left thalamus was most salient for both EMCI and LMCI, right SFG (medial) - left anterior cingulum, cuneus - left MOG.

Overall, classification performance seems high (mean: 83.6%) even when only looking at studies with larger datasets (above 250 subjects, mean: 81.3%). AD classification generally leads to higher performance (mean: 85.8%) than non-AD (mean: 83.1%), supporting the intuition that neurodegeneration due to AD leads to profound changes in FC that can be picked up more easily in later stages of AD. Similar to ASD, both BG and PG are well researched by existing studies.

There are a few reproducible ROI-level salient features across the 11 studies that reported salient features. For SMC, the inferior temporal gyrus was highlighted in 2 dFC studies [141], [143]. EMCI studies have numerous consistent findings. A sFC study [61] and a dFC study [135] agreed that the DMN is impacted in EMCI. Inferior temporal gyrus was highlighted in 1 dFC [141] and 2 multimodal studies [142], [143]. The 2 multimodal studies also identified the olfactory cortex as a salient feature (replicated in LMCI). For AD, the 3 studies broadly agree on frontal regions but do not seem to have any specific common ROIs. In summary, the inferior temporal gyrus was an ROI that was quite consistently highlighted across disease stages, except for AD.

F. Parkinson's Disease

Out of 3 studies on PD (all adopting a multimodal approach), 2 studies performed biomarker discovery from fMRI. One study fused multimodal data via concatenating hand-crafted features from various modalities in the node vector,

while another used a GAE approach that uses different modalities in the node features and the adjacency matrix.

In the former study, Safai *et al.* [144] used an architecture involving 2 GAT layers (6 attention heads in the first layer, 1 head in the second layer) followed by a readout layer retaining top-k node embeddings and a fully connected layer. The graph for the GAT layers was the SC matrix (BG). For node features, they extracted 7 network measures (e.g. degree, clustering coefficient) and 4 statistical measures (e.g. mean, skewness, kurtosis) from SC, FC, as well as volumetric and cortical thickness information. SC and FC were thresholded beforehand to remove weaker connections. The concatenation of these 24 features forms the input. Using the DK86 atlas on a private dataset with 109 subjects (75 PD, 34 NC ; 150 after oversampling via SMOTE [145] to reduce imbalance), they achieved an accuracy of 73.0%. **Attention weights** from the first GAT layer were extracted and fidelity scores were computed for each head to determine the head that contributed most to classification accuracy. As an alternative, they also computed saliency scores directly via **gradients** (with respect to input features). Both analyses revealed that ROIs from the cerebellum, precuneus and pallidum are important and they further verified that late onset of PD is associated with reduced structural influence in these areas.

In the latter study, Shi *et al.* [146] proposed a weighted graph auto-encoder (WGAE) architecture to model multimodal datasets. WGAE comprises a GCN layer (Kipf) with residual connection being used in the encoder. SC is used as the graph structure (BG) while functional connection profile is used as node features. The decoder has two branches. The first branch is based on an inner product as proposed in GAE [147], while the second branch involves a linear layer for edge weight reconstruction. The latter is also based on inner product, but it also introduces a learnable weight matrix. Training of the model involves optimising two losses together: a connection loss (based on reconstruction from the first branch) and an edge weight loss (based on reconstruction from the second branch). Thereafter, the latent representation at the bottleneck layer of the GAE is used for downstream classification via logistic regression. Using the BNA246 atlas on the Parkinson's Progression Markers Initiative (PPMI) dataset with 41 subjects (22 PD, 19 NC), they obtained an accuracy of 80.5% (MLP only: 58.5%, replacing GCN with MLP: 56.1%). Coefficients learnt by the **logistic regression** model were used to find the most prominent regions. This revealed that the precentral gyrus, postcentral gyrus and parietal lobe were the most salient regions.

Overall, the two studies highlighted ROIs in the parietal lobe to be important for PD but have do not have much in agreement for other regions. Such insights could be too coarse to be useful and finer details should be reported. However, all studies have small sample sizes and studies that worked on PPMI dataset noted the problem of data imbalance. More studies on PD are encouraged, especially with more recent and larger dataset released by PPMI in recent years.

V. DISCUSSION AND FUTURE WORK

Overall, we found that over two-thirds (Fig. 2(a)) of existing studies (on modelling fMRI datasets with GNNs for disorder prediction) are dominated by ASD and dementia, while other disorders are less thoroughly covered, making it challenging to assess the robustness of their potential functional biomarkers. While classification performance is generally high (mean accuracy of around 80%), only a few moderately reproducible salient features were discovered for each disorder. These are largely limited to ROI-level biomarkers, leaving much room for improvement. In this section, we will elaborate on 3 key challenges pertaining to each stage of the prediction-attribution-evaluation framework. A more detailed discussion can be found in the supplementary materials.

A. Challenges

1) **Prediction:** Although model performances were found to be generally high, a deeper analysis reveals that this might be boosted by studies that relied on small datasets (Fig. S1). A previous study [60] has demonstrated that classification accuracy drops when dataset size increases. This phenomenon is also observed in our analysis (Fig. S1). Furthermore, we observe that standard deviations are higher in studies that used smaller datasets too (Fig. S2). Standard deviation is also observed to decrease when larger datasets are used. Most studies do not demonstrate generalisability beyond the dataset used for training (i.e. to a new dataset, not just the test split).

To address these issues, we suggest that more benchmarking studies, utilising state-of-the-art GNN models and multiple datasets of the same disorder, should be conducted. Although several benchmarking studies have been done using GNNs on fMRI datasets, they were focused on baseline GNNs and not many of them performed disease classification. Another motivation for benchmarking studies is the possibility that several existing studies that proposed new architectures do not ensure fair comparisons with previous work (e.g. not carefully tuning the models being compared against, using hyperparameters that were tuned for other purposes in the original paper, not ensuring that number of parameters in all models are similar). Thus, a benchmarking study is needed to determine which GNN models are more effective. This would then pinpoint the models that should be used for biomarker discovery, since a performant model is more likely to produce explanations that reflect traits that are specific to the disorder. At present, model performance might be limited by how majority of existing studies used baseline GNNs as part of their architectures (Fig. 2(b)) instead of GNNs customised to handle connectome datasets (e.g. mix of BG/PG, adaptive methods of constructing the graph).

2) **Attribution:** Existing studies have used a diverse range of techniques to compute attribution scores. From Fig. 2(c), it is evident that attribution based on pooling (16.7%) and attention scores (11.7%) are most popular in fMRI studies that used GNNs. Notably, model agnostic approaches such as feature selection scores (10.0%), occlusion (8.3%) and gradients (6.7%) are also often used. However, there remains much room to explore GNN-specific explainers (as shown in

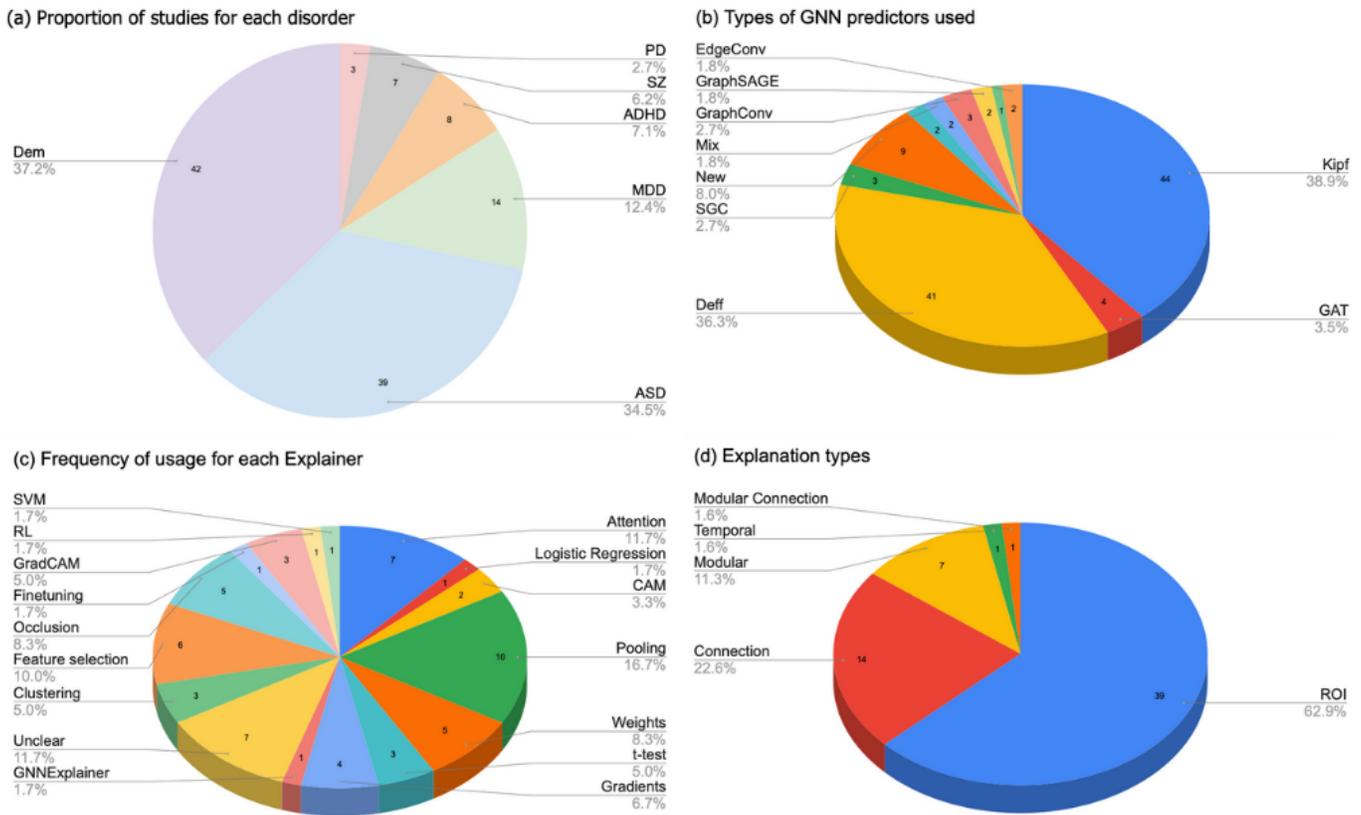


Fig. 2. Pie charts illustrating the number of studies spanning across disorders, predictors and explainers. (a) Number of studies for each disorder covered in this review paper. Studies on ASD and dementia dominate existing research. (b) Frequency of usage of GNN models. The two baseline GNNs (Kipf and Defferrard) are most commonly used in fMRI studies, but several studies have proposed GNNs customised for fMRI datasets. (c) Frequency of usage of various explainers. A myriad of explainers have been used in existing papers, but GNN-specific explainers are clearly under-explored. (d) Number of studies using each type of explanation. ROI-level analyses are the most common approach in existing biomarker discovery research.

Fig. 1(b)) as they are currently under-utilised. Another very under-explored area is temporal attribution (i.e. identifying salient time points).

Besides the choice of explainers, one major challenge faced while assessing the robustness of the potential biomarkers brought up by these studies is the lack of consistency in the way the most salient features are reported. This problem is caused by three factors: (i) different granularity of the biomarkers due to the design of the GNN architecture, (ii) different atlases used, (iii) different extent of thoroughness in reporting the salient features (e.g. listing top 10/20/30 features, or just mentioning a few salient features in passing). From Fig. 2(d), it is evident that majority of the studies (62.9%) report ROI-level features while another significant portion of studies report at the level of connections (22.6%) and functional brain modules (11.3%). Harmonising these findings manually is infeasible without improvements in reporting standards. To address this, we suggest the following guidelines as a preliminary step towards establishing a standard that future publications can refer to:

- Attribution scores for each feature should be recorded in a spreadsheet that is published alongside the publication. Minimally, they should contain the ranking of the fea-

tures, or be sorted according to importance.

- Key metadata about the most salient features (e.g. ROI names as defined by the atlas, along with their MNI coordinates) should be provided in the supplementary materials in the form of a spreadsheet, so as to facilitate future research and meta-analysis.
- Whenever explainers provide information about polarity (i.e. positive and negative scores, reflecting hyper/hypo-connectivity [149]), these raw values should be reported even if only the absolute value / magnitude is used in the analysis.

Having such guidelines would make it possible to develop computational tools to automate meta-analysis, helping research in this area to progress more quickly.

3) Evaluation: Out of the Co-12 properties, only very few existing studies have explored ‘Coherence’ by using multiple explainers in their analysis [30], [93], [144] and identifying common features that are picked up by all explainers. Much more remains to be done to assess the robustness of the attribution scores produced by explainers. At present, it is still unclear which combinations of predictors and explainers are suitable for biomarker discovery. Thus, future research should utilise the properties and metrics highlighted in Fig.

1(c) to objectively determine the robustness of explainers. Such scores would provide another set of information, on top of model prediction performances, to determine the robustness of the proposed architecture. The properties and metrics in Fig. 1(c) serve as a starting point for this research direction. Many novel evaluation metrics could also be proposed in future studies to address other desiderata not covered by these metrics. One example used in radiology applications is the prediction-saliency correlation metric [158] which computed the correlation between changes in model predictions and the corresponding saliency maps.

B. Outlook

In this section, we highlight three exciting developments in the field that concerns datasets, predictors and granularity of explanations.

1) *Growing problem of dataset heterogeneity*: The issue of small datasets has plagued the field in the past and still remains a problem, as discussed above. The availability of open-source datasets on brain disorders will increase with time as more research centres release their datasets to other researchers and this could alleviate the above problem. However, heterogeneity will also increase due to the use of newer scanners, imaging protocols and evolving standards in inclusion criteria used in studies. It is unlikely that fMRI datasets will be as large as datasets of natural images (where issues relating to heterogeneity might have been addressed due to the sheer scale of data used). Thus, this necessitates solutions to deal with these sources of variations to both data and labels.

We note that although many of the papers reviewed in this study (especially those using larger datasets) rely on datasets aggregated from multiple sites, only a few of them performed data harmonisation [129]. Existing works have demonstrated that salient features identified from such aggregated datasets tend to be biased towards the largest site [29]. The use of data harmonisation algorithms like ComBat could lead to changes to salient features [151] and these changes would require further study before they can be used for biomarker discovery. To address this, better data harmonisation techniques would need to be developed. Several approaches based on deep learning have been proposed in recent years [150], [151].

2) *Moving beyond class labels*: Binary classification serves as a useful way to study disorders since patients can be contrasted against ‘healthy’ / ‘normal’ controls. However, such class labels are a preliminary construct guided by our existing imperfect understanding of these disorders. Furthermore, datasets aggregated from multiple sites are often used without paying heed to variations in inclusion criteria, blurring the lines of what is defined as a control subject and a patient. Thus, it might not be the best approach to constrain our learning algorithms with these labels, especially when the goal is to go beyond our current understanding. For instance, psychiatric disorders have large number of possible combinations of symptoms that lead to the same diagnosis even though they have different underlying biology [156]. In these scenarios, class labels could be unreliable, especially for poorly understood disorders.

On the other hand, using a completely data-driven approach is very challenging. For instance, clustering-based approaches could be applied directly on the data but cluster interpretation is often fraught with subjectivity. Instead, we present two alternatives that are more feasible: (i) predicting test scores (such as MoCA, PANSS) or imaging-guided labels such as PET grading [153], (ii) adopting a transdiagnostic approach: psychiatric disorders often have overlaps and shared biomarkers have been identified across these disorders [154], [155].

3) *Realistic alternatives to generalisable biomarkers*: Most existing studies attempt to (i) demonstrate that their proposed architecture lead to higher model performance and (ii) produce one set of candidate biomarkers for the entire class and justify their relevance by identifying a subset of them that are also mentioned in other studies. However, our review has revealed few reproducible salient features at the ROI level and none at the level of functional connectivity features. Rather than fixating at the goal of maximising model generalisability and discovering generalisable biomarkers, identifying narrower forms of biomarkers could be a more promising direction.

This is especially so if the presence of heterogeneity is well-established for the disorder, but not well-characterised yet (else it would have been possible to use them as labels). This would suggest that there are sub-populations within the dataset that are likely to have different biomarkers. As such, it is (i) unrealistic to expect generalisation to these heterogeneous sub-populations and (ii) it would likely require going beyond the existing paradigm of producing class-wide biomarkers. Alternative approaches that produce more granular biomarkers include (i) subtype-specific biomarkers [152], (ii) studying salient features at the level of individual sites before carefully grouping sites together, considering the variability in terms of scanners, inclusion/exclusion criteria and demographics (age, gender, country, race), (iii) studying the robustness of individualised biomarkers, which would likely be the most clinically useful option if proven to be robust.

Another direction of research could involve looking beyond identifying salient nodes, connections and modules. On one hand, low-level features like edges could be too granular and vulnerable to noise and variability induced by pre-processing procedures. On the other hand, it is easy to find matches at the level of functional modules but these are often too high-level as they encompass a broad area of the brain. The scale between these extremes, i.e. motif-level analysis could prove to be more useful. This could be achieved by explainers that identify relevant sub-graphs, such as GNNExplainer [106].

VI. CONCLUSION

In summary, while there have been an abundance of novel GNN architectures designed for disorder prediction from fMRI data, there remains much room for further research to improve the robustness of the salient features highlighted by these predictors and explainers. Benchmarking studies that involve state-of-the-art GNN predictors customised for fMRI datasets are needed. Studies on optimal choices of predictors and explainers are also required as their robustness on fMRI datasets is still poorly understood. Existing evaluation metrics

were designed for generic (graph) datasets and more metrics appropriate for FC datasets are needed to determine whether explainers are sensitive to known properties of FC. The lack of standardised reporting of salient features has made it challenging to consolidate insights from existing studies. This is a complex issue that needs to be revisited once a better understanding of predictors and explainers is developed, and when better metrics have been created. Finally, the paucity of reproducible salient features (especially at the level of connections) motivates the search for alternative approaches. Possible directions for future research include looking beyond the use of class labels and pivoting away from the goal of solely chasing for generalisable biomarkers for the entire disease class, especially for heterogeneous disorders. Moving towards regression tasks, transdiagnostic studies and more fine-grained biomarkers could result in more robust biomarkers of neurological disorders being discovered from fMRI datasets in the near future.

REFERENCES

- [1] M. Filippi, E. G. Spinelli, C. Cividini, A. Ghirelli, S. Basaia, and F. Agosta, 'The human functional connectome in neurodegenerative diseases: relationship to pathology and clinical progression', *Expert Review of Neurotherapeutics*, vol. 23, no. 1, pp. 59–73, 2023.
- [2] A. Abi-Dargham *et al.*, 'Candidate biomarkers in psychiatric disorders: state of the field', *World Psychiatry*, vol. 22, no. 2, pp. 236–262, 2023.
- [3] F. Chollet and P. Payoux, 'Functional Imaging for Neurodegenerative Diseases', *La Presse Médicale*, vol. 51, no. 2, p. 104121, 2022.
- [4] S. Verdi, A. F. Marquand, J. M. Schott, and J. H. Cole, 'Beyond the average patient: how neuroimaging models can address heterogeneity in dementia', *Brain*, vol. 144, no. 10, pp. 2946–2953, 2021.
- [5] E. Canario, D. Chen, and B. Biswal, 'A review of resting-state fMRI and its use to examine psychiatric disorders', *Psychoradiology*, vol. 1, no. 1, pp. 42–53, 2021.
- [6] R. A. Poldrack *et al.*, 'Scanning the horizon: towards transparent and reproducible neuroimaging research', *Nature reviews neuroscience*, vol. 18, no. 2, pp. 115–126, 2017.
- [7] X.-Z. Jia *et al.*, 'Small effect size leads to reproducibility failure in resting-state fMRI studies', *BioRxiv*, p. 285171, 2018.
- [8] T. T. Liu, 'Noise contributions to the fMRI signal: An overview', *NeuroImage*, vol. 143, pp. 141–151, 2016.
- [9] S. Marek *et al.*, 'Reproducible brain-wide association studies require thousands of individuals', *Nature*, vol. 603, no. 7902, pp. 654–660, 2022.
- [10] J. M. M. Bayer *et al.*, 'Site effects how-to and when: An overview of retrospective techniques to accommodate site effects in multi-site neuroimaging analyses', *Frontiers in Neurology*, vol. 13, p. 923988, 2022.
- [11] K. Dadi *et al.*, 'Benchmarking functional connectome-based predictive models for resting-state fMRI', *NeuroImage*, vol. 192, pp. 115–134, 2019.
- [12] R. Botvinik-Nezer *et al.*, 'Variability in the analysis of a single neuroimaging dataset by many teams', *Nature*, vol. 582, no. 7810, pp. 84–88, 2020.
- [13] A. R. Laird, 'Large, open datasets for human connectomics research: Considerations for reproducible and responsible data use', *NeuroImage*, vol. 244, p. 118579, 2021.
- [14] Esteban, Oscar, Christopher J. Markiewicz, Ross W. Blair, Craig A. Moodie, A. Ilkay Isik, Asier Erramuzpe, James D. Kent *et al.* "fMRI-Prep: a robust preprocessing pipeline for functional MRI." *Nature methods* 16, no. 1 (2019): 111-116.
- [15] F. Hu *et al.*, 'Image harmonization: A review of statistical and deep learning methods for removing batch effects and evaluation metrics for effective harmonization', *NeuroImage*, p. 120125, 2023.
- [16] C. Habeck and Y. Stern, 'Multivariate data analysis for neuroimaging data: overview and application to Alzheimer's disease', *Cell biochemistry and biophysics*, vol. 58, no. 2, pp. 53–67, 2010.
- [17] M. E. Weaverdyck, M. D. Lieberman, and C. Parkinson, 'Tools of the Trade Multivoxel pattern analysis in fMRI: a practical introduction for social and affective neuroscientists', *Social Cognitive and Affective Neuroscience*, vol. 15, no. 4, pp. 487–509, 2020.
- [18] L. Zhang, M. Wang, M. Liu, and D. Zhang, 'A survey on deep learning for neuroimaging-based brain disorder analysis', *Frontiers in neuroscience*, vol. 14, p. 779, 2020.
- [19] J. Kawahara *et al.*, 'BrainNetCNN: Convolutional neural networks for brain networks; towards predicting neurodevelopment', *NeuroImage*, vol. 146, pp. 1038–1049, 2017.
- [20] R. J. Meszlényi, K. Buza, and Z. Vidnyánszky, 'Resting state fMRI functional connectivity-based classification using a convolutional neural network architecture', *Frontiers in neuroinformatics*, vol. 11, p. 61, 2017.
- [21] S. Gupta, Y. H. Chan, J. C. Rajapakse, A. D. N. Initiative, and Others, 'Obtaining leaner deep neural networks for decoding brain functional connectome in a single shot', *Neurocomputing*, vol. 453, pp. 326–336, 2021.
- [22] A. Bessadok, M. A. Mahjoub, and I. Rekek, 'Graph neural networks in network neuroscience', *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 5, pp. 5833–5848, 2022.
- [23] Q. Wang *et al.*, 'Leveraging Brain Modularity Prior for Interpretable Representation Learning of fMRI', *arXiv preprint arXiv:2306.14080*, 2023.
- [24] S. Parisot *et al.*, 'Disease Prediction using Graph Convolutional Networks: Application to Autism Spectrum Disorder and Alzheimer's Disease', *Medical image analysis*, 2018.
- [25] A. ElGazzar, R. Thomas, and G. Van Wingen, 'Benchmarking Graph Neural Networks for FMRI analysis', *arXiv preprint arXiv:2211.08927*, 2022.
- [26] X. Song *et al.*, 'Graph convolution network with similarity awareness and adaptive calibration for disease-induced deterioration prediction', *Medical Image Analysis*, vol. 69, p. 101947, 2021.
- [27] T. Xiao, L. Zeng, X. Shi, X. Zhu, and G. Wu, 'Dual-Graph Learning Convolutional Networks for Interpretable Alzheimer's Disease Diagnosis', in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2022, pp. 406–415.
- [28] M. Duda *et al.*, 'Reliability and clinical utility of spatially constrained estimates of intrinsic functional networks from very short fMRI scans', *Human Brain Mapping*, vol. 44, no. 6, pp. 2620–2635, 2023.
- [29] Y. H. Chan, W. C. Yew, and J. C. Rajapakse, 'Semi-supervised Learning with Data Harmonisation for Biomarker Discovery from Resting State fMRI', in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2022, pp. 441–451.
- [30] S. Gallo *et al.*, 'Functional connectivity signatures of major depressive disorder: machine learning analysis of two multicenter neuroimaging studies', *Molecular Psychiatry*, pp. 1–10, 2023.
- [31] X. Li *et al.*, 'Braingnn: Interpretable brain graph neural network for fmri analysis', *Medical Image Analysis*, vol. 74, p. 102233, 2021.
- [32] R. Li *et al.*, 'Graph signal processing, graph neural network and graph learning on biological data: a systematic review', *IEEE Reviews in Biomedical Engineering*, 2021.
- [33] H. Cui *et al.*, 'Braingb: A benchmark for brain network analysis with graph neural networks', *IEEE transactions on medical imaging*, vol. 42, no. 2, pp. 493–506, 2022.
- [34] A. Said *et al.*, 'NeuroGraph: Benchmarks for Graph Machine Learning in Brain Connectomics', *arXiv preprint arXiv:2306.06202*, 2023.
- [35] P. Linardatos, V. Papastefanopoulos, and S. Kotsiantis, 'Explainable ai: A review of machine learning interpretability methods', *Entropy*, vol. 23, no. 1, p. 18, 2020.
- [36] E. Tjoa and C. Guan, 'A survey on explainable artificial intelligence (xai): Toward medical xai', *IEEE transactions on neural networks and learning systems*, vol. 32, no. 11, pp. 4793–4813, 2020.
- [37] R. Marcinkevičs and J. E. Vogt, 'Interpretable and explainable machine learning: A methods-centric overview with concrete examples', *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, p. e1493, 2023.
- [38] H. Yuan, H. Yu, S. Gui, and S. Ji, 'Explainability in graph neural networks: A taxonomic survey', *IEEE transactions on pattern analysis and machine intelligence*, vol. 45, no. 5, pp. 5782–5799, 2022.
- [39] J. Kakkad, J. Jannu, K. Sharma, C. Aggarwal, and S. Medya, 'A Survey on Explainability of Graph Neural Networks', *arXiv preprint arXiv:2306.01958*, 2023.
- [40] C. Agarwal, O. Queen, H. Lakkaraju, and M. Zitnik, 'Evaluating explainability for graph neural networks', *Scientific Data*, vol. 10, no. 1, p. 144, 2023.
- [41] L. Vizioli *et al.*, 'Lowering the thermal noise barrier in functional brain mapping with magnetic resonance imaging', *Nature communications*, vol. 12, no. 1, p. 5181, 2021.
- [42] V. D. Calhoun, J. Liu, and T. Adalı, 'A review of group ICA for fMRI data and ICA for joint inference of imaging, genetic, and ERP data', *Neuroimage*, vol. 45, no. 1, pp. S163–S172, 2009.

- [43] S.-J. Hong *et al.*, 'Toward a connectivity gradient-based framework for reproducible biomarker discovery', *NeuroImage*, vol. 223, p. 117322, 2020.
- [44] E. T. Rolls, C.-C. Huang, C.-P. Lin, J. Feng, and M. Joliot, 'Automated anatomical labelling atlas 3', *Neuroimage*, vol. 206, p. 116189, 2020.
- [45] R. S. Desikan *et al.*, 'An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest', *Neuroimage*, vol. 31, no. 3, pp. 968–980, 2006.
- [46] B. A. Seitzman *et al.*, 'A set of functionally-defined brain regions with improved representation of the subcortex and cerebellum', *Neuroimage*, vol. 206, p. 116290, 2020.
- [47] C. Yang, P. Wang, J. Tan, Q. Liu, and X. Li, 'Autism spectrum disorder diagnosis using graph attention network based on spatial-constrained sparse functional brain networks', *Computers in Biology and Medicine*, vol. 139, p. 104963, 2021.
- [48] Y. Du, Z. Fu, and V. D. Calhoun, 'Classification and prediction of brain disorders using functional connectivity: promising but challenging', *Frontiers in neuroscience*, vol. 12, p. 525, 2018.
- [49] M. Nauta *et al.*, 'From anecdotal evidence to quantitative evaluation methods: A systematic review on evaluating explainable ai', *ACM Computing Surveys*, vol. 55, no. 13s, pp. 1–42, 2023.
- [50] K. Amara *et al.*, 'GraphframeX: Towards systematic evaluation of explainability methods for graph neural networks', *arXiv preprint arXiv:2206.09677*, 2022.
- [51] L. Wang, K. Li, and X. P. Hu, 'Graph convolutional network for fMRI analysis based on connectivity neighborhood', *Network Neuroscience*, vol. 5, no. 1, pp. 83–95, 2021.
- [52] L. Zhang, J.-R. Wang, and Y. Ma, 'Graph Convolutional Networks via Low-Rank Subspace for Multi-Site rs-fMRI ASD Diagnosis', in *2021 14th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMED)*, 2021, pp. 1–6.
- [53] K.-M. Binti, T. T. Mueller, S. Starck, V. Baltatzis, A. Hammers, and D. Rueckert, 'A Comparative Study of Population-Graph Construction Methods and Graph Neural Networks for Brain Age Regression', *arXiv preprint arXiv:2309.14816*, 2023.
- [54] L. Cosmo, A. Kazi, S.-A. Ahmadi, N. Navab, and M. Bronstein, 'Latent-graph learning for disease prediction', in *Medical Image Computing and Computer Assisted Intervention—MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part II* 23, 2020, pp. 643–653.
- [55] H. W. Park, S. Y. Kim, and W. H. Lee, 'Graph Convolutional Network with Morphometric Similarity Networks for Schizophrenia Classification', in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2023, pp. 626–636.
- [56] H. Jiang, P. Cao, M. Xu, J. Yang, and O. Zaiane, 'Hi-GCN: A hierarchical graph convolution network for graph embedding learning of brain network and brain disorders prediction', *Computers in Biology and Medicine*, vol. 127, p. 104096, 2020.
- [57] Y. He, Y. H. Chan, and J. C. Rajapakse, 'Predicting gender from structural and functional connectomes via brain and population graph neural networks', *bioRxiv*, pp. 2023–2011, 2023.
- [58] H. Zhang *et al.*, 'Classification of Brain Disorders in rs-fMRI via Local-to-Global Graph Neural Networks', *IEEE Transactions on Medical Imaging*, 2022.
- [59] W. Yin, L. Li, and F.-X. Wu, 'A graph attention neural network for diagnosing ASD with fMRI data', in *2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2021, pp. 1131–1136.
- [60] J. Teng, C. Mi, J. Shi, and N. Li, 'Brain disease research based on functional magnetic resonance imaging data and machine learning: a review', *Frontiers in Neuroscience*, vol. 17, 2023.
- [61] L. Mei *et al.*, 'Modular graph encoding and hierarchical readout for functional brain network based eMCI diagnosis', in *MICCAI Workshop on Imaging Systems for GI Endoscopy*, 2022, pp. 69–78.
- [62] M. Defferrard, X. Bresson, and P. Vandergheynst, 'Convolutional neural networks on graphs with fast localized spectral filtering', *Advances in neural information processing systems*, vol. 29, 2016.
- [63] M. He, Z. Wei, and J.-R. Wen, 'Convolutional neural networks on graphs with chebyshev approximation, revisited', *Advances in Neural Information Processing Systems*, vol. 35, pp. 7264–7276, 2022.
- [64] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, 'Graph attention networks', *arXiv preprint arXiv:1710.10903*, 2017.
- [65] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, 'How powerful are graph neural networks?', *arXiv preprint arXiv:1810.00826*, 2018.
- [66] T. N. Kipf and M. Welling, 'Semi-supervised classification with graph convolutional networks', *arXiv preprint arXiv:1609.02907*, 2016.
- [67] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, 'Dynamic graph cnn for learning on point clouds', *ACM Transactions on Graphics (tog)*, vol. 38, no. 5, pp. 1–12, 2019.
- [68] C. Morris *et al.*, 'Weisfeiler and leman go neural: Higher-order graph neural networks', in *Proceedings of the AAAI conference on artificial intelligence*, 2019, vol. 33, pp. 4602–4609.
- [69] Q. Wang, M. Wu, Y. Fang, W. Wang, L. Qiao, and M. Liu, 'Modularity-Constrained Dynamic Representation Learning for Interpretable Brain Disorder Analysis with Functional MRI', in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2023, pp. 46–56.
- [70] K. K. Jain and K. K. Jain, *The handbook of biomarkers*. Springer, 2010.
- [71] J. K. Aronson and R. E. Ferner, 'Biomarkers—a general review', *Current protocols in pharmacology*, vol. 76, no. 1, pp. 9–23, 2017.
- [72] L. Parkes, T. D. Satterthwaite, D. S. Bassett, 'Towards precise resting-state fMRI biomarkers in psychiatry: synthesizing developments in transdiagnostic research, dimensional models of psychopathology, and normative neurodevelopment', *Current Opinion in Neurobiology*, vol. 65, p. 120-8, 2020.
- [73] M.N. Coutanche, S. L. Thompson-Schill, R. T. Schultz, 'Multi-voxel pattern analysis of fMRI data predicts clinical symptom severity', *Neuroimage*, vol. 57, no. 1, p. 113-123, 2011.
- [74] S. Klöppel, A. Abdulkadir, C. R. Jack Jr, N. Koutsouleris, J. Mourão-Miranda, P. Vemuri, 'Diagnostic neuroimaging across diseases', *Neuroimage*, vol. 61, no. 2, p. 457-463, 2012.
- [75] T. Wolfers, J. K. Buitelaar, C. F. Beckmann, B. Franke, A. F. Marquand, 'From estimating activation locality to predicting disorder: a review of pattern recognition for neuroimaging-based psychiatric diagnostics', *Neuroscience & Biobehavioral Reviews*, vol. 57, p. 328-349, 2015.
- [76] G. Orru, W. Pettersson-Yeo, A. F. Marquand, G. Sartori, A. Mechelli, 'Using support vector machine to identify imaging biomarkers of neurological and psychiatric disease: a critical review', *Neuroscience & Biobehavioral Reviews*, vol. 36, no. 4, p. 1140-1152, 2012.
- [77] R. de Filippis, E. A. Carbone, R. Gaetano, A. Bruni, V. Pugliese, C. Segura-Garcia, P. De Fazio, 'Machine learning techniques in a structural and functional MRI diagnostic approach in schizophrenia: a systematic review', *Neuropsychiatric disease and treatment*, p. 1605-1627, 2019.
- [78] Y. Liu, Y. Zhang, L. Lv, R. Wu, J. Zhao, W. Guo, 'Abnormal neural activity as a potential biomarker for drug-naive first-episode adolescent-onset schizophrenia with coherence regional homogeneity and support vector machine analyses', *Schizophrenia research*, vol. 192, p. 408-415, 2018.
- [79] A. Sarica, A. Cerasa, A. Quattrone, 'Random forest algorithm for the classification of neuroimaging data in Alzheimer's disease: a systematic review', *Frontiers in aging neuroscience*, vol. 9, p. 329, 2017.
- [80] A. J. Fredo, A. Jahedi, M. Reiter, R. A. Müller, 'Diagnostic classification of autism using resting-state fMRI data and conditional random forest', *Age*, vol. 12, no. 2, p. 6-41, 2018.
- [81] Y. Chen *et al.*, 'Adversarial learning based node-edge graph attention networks for autism spectrum disorder identification', *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [82] K. Qin *et al.*, 'Using graph convolutional network to characterize individuals with major depressive disorder across multiple imaging sites', *EBioMedicine*, vol. 78, 2022.
- [83] P. T. Fox *et al.*, 'BrainMap taxonomy of experimental design: description and evaluation', *Human brain mapping*, vol. 25, no. 1, pp. 185–198, 2005.
- [84] N. R. Winter *et al.*, 'A Systematic Evaluation of Machine Learning-Based Biomarkers for Major Depressive Disorder', *JAMA psychiatry*, 2024.
- [85] S. Ali *et al.*, 'Explainable Artificial Intelligence (XAI): What we know and what is left to attain Trustworthy Artificial Intelligence', *Information Fusion*, vol. 99, p. 101805, 2023.
- [86] P. E. Pope, S. Kolouri, M. Rostami, C. E. Martin, and H. Hoffmann, 'Explainability methods for graph convolutional neural networks', in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 10772–10781.
- [87] T. Ma and A. Zhang, 'Incorporating biological knowledge with factor graph neural network for interpretable deep learning', *arXiv preprint arXiv:1906.00537*, 2019.
- [88] S. Miao, Y. Luo, M. Liu, and P. Li, 'Interpretable geometric deep learning via learnable randomness injection', *arXiv preprint arXiv:2210.16966*, 2022.
- [89] J. Yu, T. Xu, Y. Rong, Y. Bian, J. Huang, and R. He, 'Recognizing predictive substructures with subgraph information bottleneck', *IEEE transactions on pattern analysis and machine intelligence*, vol. 46, no. 3, pp. 1650–1663, 2021.

- [90] K. Zheng, S. Yu, B. Li, R. Jenssen, and B. Chen, 'Brainib: Interpretable brain network-based psychiatric diagnosis with graph information bottleneck', arXiv preprint arXiv:2205.03612, 2022.
- [91] C. Liu *et al.*, 'Graph pooling for graph neural networks: progress, challenges, and opportunities', in Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, 2023, pp. 6712–6722.
- [92] B. T. T. Yeo *et al.*, 'The organization of the human cerebral cortex estimated by intrinsic functional connectivity', *Journal of neurophysiology*, 2011.
- [93] Z. Zhang *et al.*, 'Identifying biomarkers of subjective cognitive decline using graph convolutional neural network for fMRI analysis', in 2022 IEEE International Conference on Mechatronics and Automation (ICMA), 2022, pp. 1306–1311.
- [94] A. Yu, L. Chen, and C. Qiao, 'Graph Convolutional Network with Attention Mechanism for Discovering the Brain's Abnormal Activity of Attention Deficit Hyperactivity Disorder', in 2022 15th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI), 2022, pp. 1–5.
- [95] S. Jain and B. C. Wallace, 'Attention is not Explanation', in Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), 2019, pp. 3543–3556.
- [96] S. Wiegrefe and Y. Pinter, 'Attention is not not Explanation', in Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), 2019, pp. 11–20.
- [97] S. Serrano and N. A. Smith, 'Is Attention Interpretable?', in Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, 2019, pp. 2931–2951.
- [98] B. Bai, J. Liang, G. Zhang, H. Li, K. Bai, and F. Wang, 'Why attentions may not be interpretable?', in Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining, 2021, pp. 25–34.
- [99] M. Sundararajan, A. Taly, and Q. Yan, 'Axiomatic attribution for deep networks', in International conference on machine learning, 2017, p. 3319–3328.
- [100] P. Sturmfels, S. Lundberg, and S.-I. Lee, 'Visualizing the impact of feature attribution baselines', *Distill*, vol. 5, no. 1, p. e22, 2020.
- [101] A. Kapishnikov, S. Venugopalan, B. Avci, B. Wedin, M. Terry, and T. Bolukbasi, 'Guided integrated gradients: An adaptive path method for removing noise', in Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2021, pp. 5050–5058.
- [102] G. Montavon, A. Binder, S. Lapuschkin, W. Samek, and K.-R. Müller, 'Layer-wise relevance propagation: an overview', *Explainable AI: interpreting, explaining and visualizing deep learning*, pp. 193–209, 2019.
- [103] T. Schnake *et al.*, 'Higher-order explanations of graph neural networks via relevant walks', *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 11, pp. 7581–7596, 2021.
- [104] W. Yan *et al.*, 'Discriminating schizophrenia from normal controls using resting state functional network connectivity: A deep neural network and layer-wise relevance propagation method', in 2017 IEEE 27th international workshop on machine learning for signal processing (MLSP), 2017, pp. 1–6.
- [105] S. M. Lundberg and S.-I. Lee, 'A unified approach to interpreting model predictions', *Advances in neural information processing systems*, vol. 30, 2017.
- [106] Z. Ying, D. Bourgeois, J. You, M. Zitnik, and J. Leskovec, 'Gnnexplainer: Generating explanations for graph neural networks', *Advances in neural information processing systems*, vol. 32, 2019.
- [107] M. T. Ribeiro, S. Singh, and C. Guestrin, "'Why should i trust you?'" Explaining the predictions of any classifier', in Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining, 2016, pp. 1135–1144.
- [108] Q. Huang, M. Yamada, Y. Tian, D. Singh, and Y. Chang, 'Graphlime: Local interpretable model explanations for graph neural networks', *IEEE Transactions on Knowledge and Data Engineering*, 2022.
- [109] H. Yuan, J. Tang, X. Hu, and S. Ji, 'Xggn: Towards model-level explanations of graph neural networks', in Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2020, pp. 430–438.
- [110] X. Wang and H. W. Shen, 'GNNInterpreter: A Probabilistic Generative Model-Level Explanation for Graph Neural Networks', in The Eleventh International Conference on Learning Representations, 2022.
- [111] Y. Li, J. Zhou, S. Verma, and F. Chen, 'A survey of explainable graph neural networks: Taxonomy and evaluation metrics', arXiv preprint arXiv:2207.12599, 2022.
- [112] P. Atanasova, J. G. Simonsen, C. Lioma, and I. Augenstein, 'A Diagnostic Study of Explainability Techniques for Text Classification', in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), 2020, pp. 3256–3274.
- [113] M. Tahmasian *et al.*, 'A systematic review on the applications of resting-state fMRI in Parkinson's disease: Does dopamine replacement therapy play a role?', *Cortex*, vol. 73, pp. 80–105, 2015.
- [114] K. Zhao, B. Duka, H. Xie, D. J. Oathes, V. Calhoun, and Y. Zhang, 'A dynamic graph convolutional neural network framework reveals new insights into connectome dysfunctions in ADHD', *Neuroimage*, vol. 246, p. 118774, 2022.
- [115] J. Li, F. Wang, J. Pan, and Z. Wen, 'Identification of autism spectrum disorder with functional graph discriminative network', *Frontiers in Neuroscience*, vol. 15, p. 729937, 2021.
- [116] Y. Wang, J. Liu, Y. Xiang, J. Wang, Q. Chen, and J. Chong, 'MAGE: automatic diagnosis of autism spectrum disorders using multi-atlas graph convolutional networks and ensemble learning', *Neurocomputing*, vol. 469, pp. 346–353, 2022.
- [117] Z. Zhu, B. Wang, and S. Li, 'A triple-pooling graph neural network for multi-scale topological learning of brain functional connectivity: Application to ASD diagnosis', in Artificial Intelligence: First CAAI International Conference, ICICAI 2021, Hangzhou, China, June 5–6, 2021, Proceedings, Part II 1, 2021, pp. 359–370.
- [118] L. Shao, C. Fu, Y. You, and D. Fu, 'Classification of ASD based on fMRI data with deep learning', *Cognitive Neurodynamics*, vol. 15, no. 6, pp. 961–974, 2021.
- [119] L. Li *et al.*, 'TE-HI-GCN: An ensemble of transfer hierarchical graph convolutional networks for disorder diagnosis', *Neuroinformatics*, pp. 1–23, 2022.
- [120] H. Zhu, J. Wang, Y.-P. Zhao, M. Lu, and J. Shi, 'Contrastive multi-view composite graph convolutional networks based on contribution learning for autism spectrum disorder classification', *IEEE Transactions on Biomedical Engineering*, 2022.
- [121] W. Cui *et al.*, 'Dynamic multi-site graph convolutional network for autism spectrum disorder identification', *Computers in Biology and Medicine*, vol. 157, p. 106749, 2023.
- [122] Y. Chen, A. Liu, X. Fu, J. Wen, and X. Chen, 'An invertible dynamic graph convolutional network for multi-Center ASD classification', *Frontiers in Neuroscience*, vol. 15, p. 828512, 2022.
- [123] Y. Chen *et al.*, 'Attention-based node-edge graph convolutional networks for identification of autism spectrum disorder using multi-modal mri data', in Pattern Recognition and Computer Vision: 4th Chinese Conference, PRCV 2021, Beijing, China, October 29–November 1, 2021, Proceedings, Part III 4, 2021, pp. 374–385.
- [124] J. Mao, Y. Sheng, W. Lan, X. Tian, J. Liu, and Y. Pan, 'Graph Convolutional Networks Based on Relational Attention Mechanism for Autism Spectrum Disorders Diagnosis', in International Conference on Intelligent Robotics and Applications, 2022, pp. 341–348.
- [125] Y. Kong *et al.*, 'Multi-stage graph fusion networks for major depressive disorder diagnosis', *IEEE Transactions on Affective Computing*, vol. 13, no. 4, pp. 1917–1928, 2022.
- [126] E. Jun, K.-S. Na, W. Kang, J. Lee, H.-I. Suk, and B.-J. Ham, 'Identifying resting-state effective connectivity abnormalities in drug-naïve major depressive disorder diagnosis via graph convolutional networks', *Human Brain Mapping*, vol. 41, no. 17, pp. 4997–5014, 2020.
- [127] Y. Kong *et al.*, 'Spatio-temporal graph convolutional network for diagnosis and treatment response prediction of major depressive disorder from functional connectivity', *Human brain mapping*, vol. 42, no. 12, pp. 3922–3933, 2021.
- [128] Y. Fang, M. Wang, G. G. Potter, and M. Liu, 'Unsupervised cross-domain functional MRI adaptation for automated major depressive disorder identification', *Medical image analysis*, vol. 84, p. 102707, 2023.
- [129] D. Lei *et al.*, 'Graph convolutional networks reveal network-level functional dysconnectivity in schizophrenia', *Schizophrenia Bulletin*, vol. 48, no. 4, pp. 881–892, 2022.
- [130] F. Orhac *et al.*, 'A guide to ComBat harmonization of imaging biomarkers in multicenter studies', *Journal of Nuclear Medicine*, vol. 63, no. 2, pp. 172–179, 2022.
- [131] X. Chen *et al.*, 'Discriminative analysis of schizophrenia patients using graph convolutional networks: A combined multimodal MRI and connectomics analysis', *Frontiers in Neuroscience*, vol. 17, p. 1140801, 2023.
- [132] X. Chen *et al.*, 'Classification of schizophrenia patients using a graph convolutional network: A combined functional MRI and connectomics analysis', *Biomedical Signal Processing and Control*, vol. 80, p. 104293, 2023.

- [133] I. Sebenius, A. Campbell, S. E. Morgan, E. T. Bullmore, and P. Liò, 'Multimodal graph coarsening for interpretable, MRI-based brain graph neural network', in 2021 IEEE 31st International Workshop on Machine Learning for Signal Processing (MLSP), 2021, pp. 1–6.
- [134] Q. Zuo, L. Lu, L. Wang, J. Zuo, and T. Ouyang, 'Constructing brain functional network by adversarial temporal-spatial aligned transformer for early AD analysis', *Frontiers in neuroscience*, vol. 16, p. 1087176, 2022.
- [135] J. Lee, W. Ko, E. Kang, H.-I. Suk, A. D. N. Initiative, and Others, 'A unified framework for personalized regions selection and functional relation modeling for early MCI identification', *NeuroImage*, vol. 236, p. 118048, 2021.
- [136] A. Alorf and M. U. G. Khan, 'Multi-label classification of Alzheimer's disease stages from resting-state fMRI-based correlation connectivity data and deep learning', *Computers in Biology and Medicine*, vol. 151, p. 106240, 2022.
- [137] Y. Li *et al.*, 'Brain connectivity based graph convolutional networks and its application to infant age prediction', *IEEE transactions on medical imaging*, vol. 41, no. 10, pp. 2764–2776, 2022.
- [138] X. Xing *et al.*, 'DS-GCNs: connectome classification using dynamic spectral graph convolution networks with assistant task training', *Cerebral Cortex*, vol. 31, no. 2, pp. 1259–1269, 2021.
- [139] L. Wang *et al.*, 'Dementia analysis from functional connectivity network with graph neural networks', *Information Processing & Management*, vol. 59, no. 3, p. 102901, 2022.
- [140] S. Yu, G. Yue, A. Elazab, X. Song, T. Wang, and B. Lei, 'Multi-scale graph convolutional network for mild cognitive impairment detection', in *Graph Learning in Medical Imaging: First International Workshop, GLMI 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 17, 2019, Proceedings 1*, 2019, pp. 79–87.
- [141] Y. Zhu, X. Song, Y. Qiu, C. Zhao, and B. Lei, 'Structure and feature based graph U-net for early Alzheimer's disease prediction', in *Multimodal Learning for Clinical Decision Support: 11th International Workshop, ML-CDS 2021, Held in Conjunction with MICCAI 2021, Strasbourg, France, October 1, 2021, Proceedings 11*, 2021, pp. 93–104.
- [142] B. Lei *et al.*, 'Multi-scale enhanced graph convolutional network for mild cognitive impairment detection', *Pattern Recognition*, vol. 134, p. 109106, 2023.
- [143] X. Song *et al.*, 'Multicenter and Multichannel Pooling GCN for Early AD Diagnosis Based on Dual-Modality Fused Brain Network', *IEEE Transactions on Medical Imaging*, vol. 42, no. 2, pp. 354–367, 2022.
- [144] A. Safai *et al.*, 'Multimodal brain connectomics-based prediction of Parkinson's disease using graph attention networks', *Frontiers in Neuroscience*, vol. 15, p. 741489, 2022.
- [145] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, 'SMOTE: synthetic minority over-sampling technique', *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
- [146] G. Shi, Y. Zhu, F. Zhang, W. Liu, Y. Yao, and X. Li, 'Fusion Learning of Multimodal Neuroimaging with Weighted Graph AutoEncoder', in *2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2022, pp. 2467–2473.
- [147] T. N. Kipf and M. Welling, 'Variational graph auto-encoders', *arXiv preprint arXiv:1611.07308*, 2016.
- [148] Y. H. Chan, C. Wang, W. K. Soh, and J. C. Rajapakse, 'Combining neuroimaging and omics datasets for disease classification using graph neural networks', *Frontiers in Neuroscience*, vol. 16, p. 866666, 2022.
- [149] S. Gupta, M. Lim, and J. C. Rajapakse, 'Decoding task specific and task general functional architectures of the brain', *Human Brain Mapping*, vol. 43, no. 9, pp. 2801–2816, 2022.
- [150] L. An *et al.*, 'Goal-specific brain MRI harmonization', *Neuroimage*, vol. 263, p. 119570, 2022.
- [151] Y. H. Chan, W. C. Yew, Q. H. Chew, K. Sim, and J. C. Rajapakse, 'Elucidating salient site-specific functional connectivity features and site-invariant biomarkers in schizophrenia via deep neural networks', *Scientific Reports*, vol. 13, no. 1, p. 21047, 2023.
- [152] Y. H. Chan, J. L. Ang, S. Gupta, Y. He, and J. C. Rajapakse, 'Subtype-Specific Biomarkers of Alzheimer's Disease from Anatomical and Functional Connectomes via Graph Neural Networks', in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2024, pp. 2195–2199.
- [153] C. Li *et al.*, 'Predicting Brain Amyloid- β PET Grades with Graph Convolutional Networks Based on Functional MRI and Multi-Level Functional Connectivity', *Journal of Alzheimer's Disease*, no. Preprint, pp. 1–15, 2022.
- [154] M. P. van den Heuvel and O. Sporns, 'A cross-disorder connectome landscape of brain dysconnectivity', *Nature reviews neuroscience*, vol. 20, no. 7, pp. 435–446, 2019.
- [155] S. C. de Lange *et al.*, 'Shared vulnerability for connectome alterations across psychiatric and neurological brain disorders', *Nature human behaviour*, vol. 3, no. 9, pp. 988–998, 2019.
- [156] J. Chen, K. R. Patil, B. T. T. Yeo, and S. B. Eickhoff, 'Leveraging machine learning for gaining neurobiological and nosological insights in psychiatric research', *Biological psychiatry*, vol. 93, no. 1, pp. 18–28, 2023.
- [157] S. Ghosal, 'Interpretable Machine Learning and Deep Learning Frameworks for Predictive Analytics and Biomarker Discovery from Multimodal Imaging Genetics Data', *Johns Hopkins University*, 2023.
- [158] J. Zhang, H. Chao, G. Dasegowda, G. Wang, M. K. Kalra, and P. Yan, 'Revisiting the Trustworthiness of Saliency Methods in Radiology AI', *Radiology: Artificial Intelligence*, vol. 6, no. 1, 2024.
- [159] S. S. Saboksayr, J. J. Foxe, and A. Wismüller, 'Attention-deficit/hyperactivity disorder prediction using graph convolutional networks', in *Medical imaging 2020: computer-aided diagnosis*, 2020, vol. 11314, pp. 430–437.
- [160] Z. Rakhimberdina and T. Murata, 'Linear graph convolutional model for diagnosing brain disorders', in *Complex Networks and Their Applications VIII: Volume 2 Proceedings of the Eighth International Conference on Complex Networks and Their Applications COMPLEX NETWORKS 2019 8*, 2020, pp. 815–826.
- [161] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger, 'Simplifying graph convolutional networks', in *International conference on machine learning*, 2019, pp. 6861–6871.
- [162] D. Yao *et al.*, 'Triplet graph convolutional network for multi-scale analysis of functional connectivity using functional MRI', in *Graph Learning in Medical Imaging: First International Workshop, GLMI 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 17, 2019, Proceedings 1*, 2019, pp. 70–78.
- [163] D. Yao *et al.*, 'A mutual multi-scale triplet graph convolutional network for classification of brain disorders using functional or structural connectivity', *IEEE transactions on medical imaging*, vol. 40, no. 4, pp. 1279–1289, 2021.
- [164] L. Liu, Y.-P. Wang, Y. Wang, P. Zhang, and S. Xiong, 'An enhanced multi-modal brain graph network for classifying neuropsychiatric disorders', *Medical image analysis*, vol. 81, p. 102550, 2022.
- [165] D. Yao, E. Yang, L. Sun, J. Sui, and M. Liu, 'Integrating multimodal MRIs for adult ADHD identification with heterogeneous graph attention convolutional network', in *Predictive Intelligence in Medicine: 4th International Workshop, PRIME 2021, Held in Conjunction with MICCAI 2021, Strasbourg, France, October 1, 2021, Proceedings 4*, 2021, pp. 157–167.
- [166] G. Wang, L. Zhang, and L. Qiao, 'The effect of node features on GCN-based brain network classification: an empirical study', *PeerJ*, vol. 11, p. e14835, 2023.
- [167] Y. Chu, G. Wang, L. Cao, L. Qiao, and M. Liu, 'Multi-scale graph representation learning for autism identification with functional MRI', *Frontiers in Neuroinformatics*, vol. 15, p. 802305, 2022.
- [168] Q. Mai, U. Nakarmi, and M. Huang, 'BrainVGAE: end-to-end graph neural networks for noisy fMRI dataset', in *2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2022, pp. 3852–3855.
- [169] T. Yang, M. A. Al-Duailij, S. Bozdog, and F. Saeed, 'Classification of autism spectrum disorder using rs-fMRI data and graph convolutional networks', in *2022 IEEE International Conference on Big Data (Big Data)*, 2022, pp. 3131–3138.
- [170] T. Eslami, V. Mirjalili, A. Fong, A. R. Laird, and F. Saeed, 'ASD-DiagNet: a hybrid learning approach for detection of autism spectrum disorder using fMRI data', *Frontiers in neuroinformatics*, vol. 13, p. 70, 2019.
- [171] M. Cao *et al.*, 'Using DeepGCN to identify the autism spectrum disorder from multi-site resting-state data', *Biomedical Signal Processing and Control*, vol. 70, p. 103015, 2021.
- [172] R. Anirudh and J. J. Thiagarajan, 'Bootstrapping graph convolutional neural networks for autism spectrum disorder classification', in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 3197–3201.
- [173] J. Pan, H. Lin, Y. Dong, Y. Wang, and Y. Ji, 'MAMF-GCN: Multi-scale adaptive multi-channel fusion deep graph convolutional network for predicting mental disorder', *Computers in Biology and Medicine*, vol. 148, p. 105823, 2022.
- [174] S. I. Ktena *et al.*, 'Metric learning with spectral graph convolutions on brain connectivity networks', *NeuroImage*, vol. 169, pp. 431–442, 2018.
- [175] G. Ma *et al.*, 'Deep graph similarity learning for brain data analysis', in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 2019, pp. 2743–2751.

- [176] K. Masood and R. Kashef, 'Integrating graph convolutional networks (gcnn) and long short-term memory (lstm) for efficient diagnosis of autism', in International Conference on Artificial Intelligence in Medicine, 2022, pp. 110–121.
- [177] H. Felouat and S. Oukid-Khouas, 'Graph convolutional networks and functional connectivity for identification of autism spectrum disorder', in 2020 Second International Conference on Embedded & Distributed Systems (EDiS), 2020, pp. 27–32.
- [178] L. Peng, N. Wang, J. Xu, X. Zhu, and X. Li, 'GATE: graph CCA for temporal SELF-supervised learning for label-efficient fMRI analysis', IEEE Transactions on Medical Imaging, vol. 42, no. 2, pp. 391–402, 2022.
- [179] L. Liu et al., 'BrainTGL: A dynamic graph representation learning model for brain network analysis', Computers in Biology and Medicine, vol. 153, p. 106521, 2023.
- [180] M. Zhu, Y. Quan, and X. He, 'The classification of brain network for major depressive disorder patients based on deep graph convolutional neural network', Frontiers in Human Neuroscience, vol. 17, p. 1094592, 2023.
- [181] S. Venkatapathy, M. Votinov, L. Wagels, S. Kim, U. Habel, and H.-G. Jo, 'Ensemble graph neural network model for classification of major depressive disorder using whole-brain functional connectivity', Frontiers in Psychiatry, vol. 14, p. 1125339, 2023.
- [182] E. Pitsik et al., 'A graph convolutional network for classification of resting-state fMRI data', in 2022 6th Scientific School Dynamics of Complex Networks and their Applications (DCNA), 2022, pp. 223–225.
- [183] E. N. Pitsik et al., 'The topology of fMRI-based networks defines the performance of a graph neural network for the classification of patients with major depressive disorder', Chaos, Solitons & Fractals, vol. 167, p. 113041, 2023.
- [184] D. Yao, J. Sui, E. Yang, P.-T. Yap, D. Shen, and M. Liu, 'Temporal-adaptive graph convolutional network for automated identification of major depressive disorder using resting-state fMRI', in Machine Learning in Medical Imaging: 11th International Workshop, MLMI 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 4, 2020, Proceedings 11, 2020, pp. 1–10.
- [185] Q. Wang, L. Li, L. Qiao, and M. Liu, 'Adaptive multimodal neuroimage integration for major depression disorder detection', Frontiers in Neuroinformatics, vol. 16, p. 856175, 2022.
- [186] J. Lee, I. Lee, and J. Kang, 'Self-attention graph pooling', in International conference on machine learning, 2019, pp. 3734–3743.
- [187] J. Huang, X. Li, M. Wang, and D. Zhang, 'Hierarchical Representation Learning of Dynamic Brain Networks for Schizophrenia Diagnosis', in Chinese Conference on Pattern Recognition and Computer Vision (PRCV), 2020, pp. 470–479.
- [188] X. Zhao, F. Zhou, L. Ou-Yang, T. Wang, and B. Lei, 'Graph convolutional network analysis for mild cognitive impairment prediction', in 2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019), 2019, pp. 1598–1601.
- [189] Z. Qin, Z. Liu, and P. Zhu, 'Aiding Alzheimer's Disease Diagnosis Using Graph Convolutional Networks Based on rs-fMRI Data', in 2022 15th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI), 2022, pp. 1–7.
- [190] X. An, Y. Zhou, Y. Di, and D. Ming, 'Dynamic functional connectivity and graph convolution network for Alzheimer's disease classification', in Proceedings of the 2020 7th International Conference on Biomedical and Bioinformatics Engineering, 2020, pp. 1–4.
- [191] X. Song, A. Elazab, and Y. Zhang, 'Classification of mild cognitive impairment based on a combined high-order network and graph convolutional network', Ieee Access, vol. 8, pp. 42816–42827, 2020.
- [192] M. Liu, H. Zhang, F. Shi, and D. Shen, 'Building dynamic hierarchical brain networks and capturing transient meta-states for early mild cognitive impairment diagnosis', in Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part VII 24, 2021, pp. 574–583.
- [193] R. Hu, L. Peng, J. Gan, X. Shi, and X. Zhu, 'Complementary graph representation learning for functional neuroimaging identification', in Proceedings of the 30th ACM International Conference on Multimedia, 2022, pp. 3385–3393.
- [194] J. Liu, G. Tan, W. Lan, and J. Wang, 'Identification of early mild cognitive impairment using multi-modal data and graph convolutional networks', BMC bioinformatics, vol. 21, pp. 1–12, 2020.
- [195] M. Yuan and Y. Lin, 'Model selection and estimation in regression with grouped variables', Journal of the Royal Statistical Society Series B: Statistical Methodology, vol. 68, no. 1, pp. 49–67, 2006.
- [196] L. Zhang, A. Zaman, L. Wang, J. Yan, and D. Zhu, 'A cascaded multi-modality analysis in mild cognitive impairment', in Machine Learning in Medical Imaging: 10th International Workshop, MLMI 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 13, 2019, Proceedings 10, 2019, pp. 557–565.
- [197] X. Li, N. C. Dvornek, J. Zhuang, P. Ventola, and J. Duncan, 'Graph embedding using infomax for ASD classification and brain functional difference detection', in Medical Imaging 2020: Biomedical Applications in Molecular, Structural, and Functional Imaging, 2020, vol. 11317, p. 1131702.
- [198] X. Li et al., 'Pooling regularized graph neural network for fmri biomarker analysis', in Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part VII 23, 2020, pp. 625–635.
- [199] S. Gupta and J. C. Rajapakse, 'Iterative consensus spectral clustering improves detection of subject and group level brain functional modules', Scientific reports, vol. 10, no. 1, pp. 1–15, 2020.
- [200] S. Yang, D. Jin, J. Liu, and Y. He, 'Identification of young high-functioning autism individuals based on functional connectome using graph isomorphism network: A pilot study', Brain Sciences, vol. 12, no. 7, p. 883, 2022.
- [201] Y. Chu, H. Ren, L. Qiao, and M. Liu, 'Resting-State Functional MRI Adaptation with Attention Graph Convolution Network for Brain Disorder Identification', Brain Sciences, vol. 12, no. 10, p. 1413, 2022.
- [202] F. Zhao et al., 'Multi-view feature enhancement based on self-attention mechanism graph convolutional network for autism spectrum disorder diagnosis', Frontiers in human neuroscience, vol. 16, p. 918969, 2022.
- [203] F. Noman, S.-Y. Yap, R. C.-W. Phan, H. Ombao, and C.-M. Ting, 'Graph autoencoder-based embedded learning in dynamic brain networks for autism spectrum disorder identification', in 2022 IEEE International Conference on Image Processing (ICIP), 2022, pp. 2891–2895.
- [204] M. Liu et al., 'Multiscale functional connectome abnormality predicts cognitive outcomes in subcortical ischemic vascular disease', Cerebral Cortex, vol. 32, no. 21, pp. 4641–4656, 2022.
- [205] Z. Dong, J. S. X. Chong, B. C. Kok, and J. H. Zhou, 'COOP-DHGNN: a Framework for Joint Classification and Prediction of Brain Functional Connectivity Using Sparse Trajectory Dataset with Application to Early Dementia', in 2022 IEEE International Conference on Big Data (Big Data), 2022, pp. 4972–4978.
- [206] Y. Zhang, X. He, Y. H. Chan, Q. Teng, and J. C. Rajapakse, 'Multi-modal graph neural network for early diagnosis of Alzheimer's disease from sMRI and PET scans', Computers in Biology and Medicine, vol. 164, p. 107328, 2023.

SUPPLEMENTARY MATERIALS

TABLE S1

MAPPING OF COMMON ABBREVIATIONS USED TO THEIR FULL NAMES.

Abbreviation	Full name
I-WL	I-dimensional Weisfeiler-Leman
AD	Alzheimer's Disease
ADHD	Attention Deficit Hyperactivity Disorder
ASD	Autism Spectrum Disorder
BG	Brain Graph
CAM	Class Activation Mapping
CNN	Convolutional Neural Network
dFC	dynamic Functional Connectivity
DFS	Deep Feature Selection
DNN	Deep Neural Network
FC	Functional Connectivity
fMRI	functional MRI
GAT	Graph Attention Network
GCN	Graph Convolutional Network
GIN	Graph Isomorphism Network
GNN	Graph Neural Network
GO	Gene Ontology
HOFC	High-Order Functional Connectivity
HSIC	Hilbert-Schmidt Independence Criterion
ICA	Independent Component Analysis
IG	Integrated Gradients
k-NN	k-Nearest Neighbours
LSTM	Long Short-Term Memory
MCI	Mild Cognitive Impairment
MDD	Major Depressive Disorder
MI	Mutual Information
ML	Machine Learning
NC	Normal Controls
NLP	Natural Language Processing
PCA	Principal Component Analysis
PD	Parkinson's Disease
PG	Population Graph
RF	Random Forest
RFE	Recursive Feature Elimination
RL	Reinforcement Learning
ROI	Region of Interest
sFC	static Functional Connectivity
SNR	Signal-to-Noise Ratio
SVM	Support Vector Machine
SZ	Schizophrenia
TDC	Typically Developing Children

TABLE S2

KEY CHARACTERISTICS OF EXPLAINERS. METHODS ABOVE THE LINE ARE SELF-INTERPRETABLE, WHILE THOSE BELOW ARE POST-HOC METHODS. E=EDGE, N=NODE, NF=NODE FEATURES, S=SUBGRAPH.

Explainer	Approach	Granularity	Target	Trainable
KER-GNN	Structural	Instance	NF	Yes
FGNN	Structural	Instance	NF	Yes
LRI	Informational	Instance	S	Yes
GIB	Informational	Instance	S	Yes
BrainGNN	Pooling	Instance	N	Yes
GAT	Attention	Instance	E	Yes
IG	Gradients	Instance	NF	No
LRP	Decomposition	Instance	NF	No
GNN-LRP	Decomposition	Instance	E	No
Excitation BP	Perturbation	Instance	NF	No
Occlusion	Perturbation	Instance	NF	No
SHAP	Perturbation	Instance	NF	No
GNNExplainer	Perturbation	Instance	S, NF	Yes
LIME	Surrogate	Instance	NF	Yes
GraphLIME	Surrogate	Instance	NF	Yes
GNNInterpreter	Graph generation	Model	S	Yes
XGNN	Graph generation	Model	S	Yes

Model accuracy vs Size of dataset

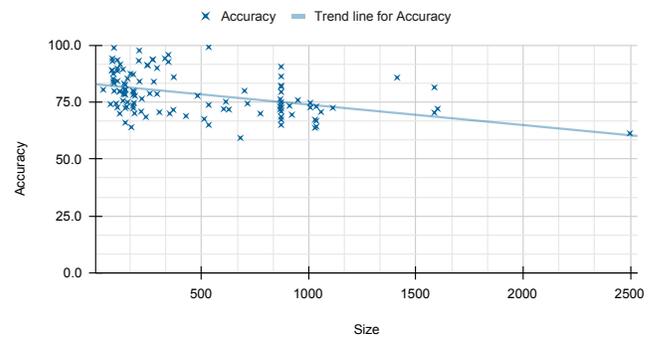


Fig. S1. Plot of model accuracy against dataset size. It is evident that majority of the studies are rather small (below 500) and classification accuracies generally decrease as the dataset size increases.

SD vs Size

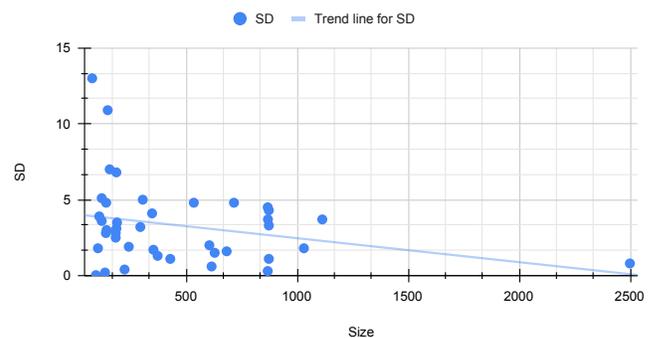


Fig. S2. Plot of the standard deviation of model accuracy against dataset size. Small datasets tend to have larger standard deviations. Standard deviations generally decrease as dataset size increases.

STUDIES WITHOUT BIOMARKER DISCOVERY

In this section, we summarise studies that did not perform biomarker discovery.

A. ADHD

4 sFC studies did not highlight salient features. Saboksayr *et al.* [159] applied the GCN framework proposed by Parisot *et al.* [24] on the ADHD-200 dataset. Spectral GCNs (Defferrard) were used in a PG approach with the graph constructed based on similarity of FC features as well as age, gender and site information. For node features, RFE is done using a ridge classifier on the FC features. Using the AAL116 atlas on the ADHD-200 dataset with 921 subjects (362 ADHD, 559 TDC), they obtained a mean accuracy of 69.5%. In contrast, Rakhimberdina *et al.* [160] used simple graph convolution (SGC) [161] (i.e. a linear graph convolutional model where non-linear activations are removed) as the graph convolution layer. PG was constructed using the hamming distance of the subject's phenotypic features (gender, handedness and site). Node features were the vectorised FC matrix. Using the HO110 atlas, their SGC model obtained an accuracy of 74.4% (ChebGCN: 71.6%) on a subset of ADHD-200 with 714 subjects (357 ADHD, 357 TDC).

Yao *et al.* [162] proposed a framework called multi-scale triplet GCN (MTGCN), where multiple templates were used to parcellate brain regions and construct corresponding FC matrices. For each template, a triplet GCN (TGCN) inputs a triplet of three graphs (anchor, positive, negative) and outputs the similarity among the triplet. A positive sample is one which belongs to the same class as the anchor sample, while the negative sample belongs to a different class from the anchor sample. Next, a weighted fusion scheme is used to combine outputs from multi-scale TGCNs and produce a prediction. Spectral GCN (Defferrard) was used, with a k-NN graph was used for the adjacency matrix and node features are based on the FC between pairs of ROIs. Using features from the AAL116 and BNT273 atlas on site PKU from the ADHD-200 dataset with 189 subjects (112 ADHD, 77 TDC), MTGCN achieved an accuracy of 77.8% (TGCN: 74.1%). This was extended in [163] to Mutual MTGCN (MMTGCN), where a template mutual learning strategy was proposed to combine the outputs of each TGCN (instead of weighted fusion). Motivated by the assumption that the topology of each template contains complementary information, the Kullback-Leibler divergence between the outputs of each pair of TGCNs is computed and used to compute a mimicry loss. By minimising this loss, MMTGCN achieved an accuracy of 71.8% (as compared to 70.6% in MTGCN) on a larger subset of ADHD-200 with 627 subjects (276 ADHD, 351 TDC) and applying a wider range of atlases (AAL116, BN273, CC200, BASC325).

2 multimodal studies were performed on ADHD datasets. Liu *et al.* [164] used a vanilla 2-layer GCN architecture (Defferrard), but they proposed a novel multimodal fusion approach. Important features from SC were extracted via XGBoost and used to adjust the corresponding edge weights in the functional graph. Instead of using FC directly, the functional graph was built based on Mahalanobis distance

with non-linear projection via a Gaussian kernel. Using the AAL90 atlas on the Consortium for Neuropsychiatric Phenomics (CNP) dataset, their model achieves an accuracy of 93.7% (notably, alternative models using a single modality obtained 80 – 90% accuracy). However, their disease class is formed by combining ADHD, SZ and bipolar disorder patient population. Thus, even though they identified salient features, it is impossible to determine whether they are ADHD-specific.

Yao *et al.* [165] proposed a heterogeneous graph attention convolutional network (HGACN) for modelling multi-modal dataset. FC values from fMRI and fractional anisotropy values from diffusion MRI (dMRI) were combined to form a single heterogeneous graph that has two types of nodes (fMRI or dMRI) and three types of edges (functionally connected, structurally connected, location-related). This graph is used to generate an attention matrix by using attention scores to learn relationships between nodes of the same and different types. Subsequently, these matrices are used as input to a heterogeneous graph convolution layer (BG) where each node type has its own graph convolution weights. This is followed by a fully connected layer for classification. Using the AAL116 atlas on a private multimodal dataset with 187 subjects (110 ADHD, 77 TDC), HGACN achieves an accuracy of 70.1% (GCN with early fusion: 67.9%).

B. ASD

Wang *et al.* [166] experimented with 5 different node features for GCNs (Kipf, BG) to study how they influence model performance. This includes (i) the original fMRI signals, (ii) one-hot encoding of the node position, (iii) 8 node statistics (local efficiency, 4 centralities measures, and 3 local clustering coefficient), (iv) node correlation, and (v) a combination of all 4 options. The adjacency matrix used is a binarised matrix, but full details on the binarisation process do not seem to be provided. The GCN is followed by a readout layer that comprises a concatenation of the outputs of max pooling and average pooling layers. Using the AAL116 atlas on a subset of the ABIDE dataset with 184 subjects (79 ASD, 105 TDC), they obtained an accuracy of 74.7% when node statistics and node correlation were used as node features (62.6% when only node correlation was used).

Yin *et al.* [59] used an architecture with 3 GATv2 layers (Kipf, BG) followed by global max pooling and a fully connected layer for ASD classification. The graph used by GAT is a thresholded FC matrix, while the node vector contains handcrafted features: 3 graph centralities (Graph degree, betweenness, eccentricity), 4 statistics of the fMRI time series (mean, variance, skewness and kurtosis). Using the Power264 atlas on ABIDE dataset with 871 subjects (403 ASD, 468 TDC), they obtained an accuracy of 82.3% (autoencoder + DNN: 79.2%).

Zhang *et al.* [52] proposed a GCN based on low-rank subspace (LRGCN) to address the issue of inter-site heterogeneity. Low-Rank Representation (LRR) with group sparse regularisation is used to learn a harmonised low rank subspace from multiple sites. ChebGCN (BG) was used with k-NN graph and node vector using the features produced from LRR. Using the

TABLE S3

SUMMARY OF FINDINGS FROM ADHD STUDIES, INCLUDING THOSE WITH BIOMARKER DISCOVERY PERFORMED (PLACED ABOVE MIDLINE). ‘SIZE’ REFERS TO THE SIZE OF THE DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL).

Reference	Dataset (Size)	Dataset distribution	Atlas	GNN	Graph	Result	Baseline
[94]	Private (240)	120 ADHD, 120 TDC	CC200	Kipf	BG	68.5%	GAT 58.6%
[114]	ADHD-200 (603)	260 ADHD, 343 TDC	AAL116	New	BG	72.0%	GAT 68.0%
[159]	ADHD-200 (921)	362 ADHD, 559 TDC	AAL116	Defferrard	PG	69.5%	-
[160]	ADHD-200 (714)	357 ADHD, 357 TDC	HO110	SGC	PG	74.4%	Cheb 71.6%
[162]	PKU (189)	112 ADHD, 77 TDC	AAL116, BN273	Defferrard	BG	77.8%	TGCN 74.1%
[163]	ADHD-200 (627)	276 ADHD, 351 TDC	Multiple	Defferrard	BG	71.8%	MTGCN 70.6%
[164]	UCLA (272) [M]	41 ADHD, 115 TDC	AAL90	Defferrard	BG	93.7%	SVM 90.2%
[165]	Private (187) [M]	110 ADHD, 77 TDC	AAL116	New	BG	70.1%	GCN 67.9%

HO110 atlas on a subset of ABIDE (CPAC, only 5 sites used: NYU, UM, UCLA, USM, Leuven) with 474 subjects (221 ASD, 253 TDC), LRGCN obtained an accuracy of 73.5% on site NYU (Denoising autoencoder: 66.0%) amongst numerous results reported on individual sites.

Yang *et al.* [47] used 2 layers of GAT (8 heads with exponential linear unit activation) via BG and passed the learnt representation to an MLP to perform ASD classification. Node features are based on feature selection of FC features via two-sample t-test, followed by LASSO. The key novelty introduced is the use of Pearson’s correlation-based Spatial Constraints Representation (PSCR) to construct spatial-constrained sparse functional brain networks for the BG. They argue that Pearson’s correlation (PC) does not remove potential confounding effects of other brain regions, while partial correlation is inappropriate when the number of time points is fewer than the number of brain regions. Sparse representation (SR) provides regularisation effects that suppress weak connections (due to noise) but have lower effect size, poorer robustness and limited scalability than PC. Thus, they proposed PSCR which balances between PC and Spatial Constraints Representation. Using the HO110 atlas on the ABIDE I dataset with 871 subjects (468 TDC, 403 ASD, C-PAC pipeline), they achieved an accuracy of 72.4%. Their experiments further showed that BG construction influences model performance. PSCR does better than PC, but PC outperforms SR. They did not explicitly perform model explainability, but provided a visualisation of features selected via t-test and LASSO.

Chu *et al.* [167] hypothesized that existing works ignore potential complementary functional topology information at different spatial scales. Therefore, they combined information from multiple scales (by using two atlases: AAL116 and CC200). For each scale, they formed the adjacency matrix by using the absolute value of the Pearson’s Correlation and the node features by using the connection profile. GCN (Defferrard, BG) embeddings obtained from the different scales are concatenated to train the model. On a subset of the ABIDE

dataset with 184 subjects (79 ASD, 105 TDC), they obtained an accuracy of 79.5% (GCN: 75.8% with AAL116, 75.3% with CC200).

Mai *et al.* [168] proposed the use of variational graph autoencoder (GVAE) that takes in an FC matrix and produces a modified graph to improve robustness of classifiers on noisy fMRI datasets. The encoder is made up of 2 GCN layers (Kipf, BG) and the decoder produces an edge probability matrix. This matrix goes through further sampling without changing the number of edges connected to each ROI. This updated matrix is then used for downstream classification tasks. Using the CC200 atlas on the ABIDE dataset with 1035 subjects, they achieved an accuracy of 66.8% (BrainGNN: 57.3%).

Yang *et al.* [169] used vanilla GCNs (Kipf, PG) but they proposed to incorporate graphlet topological counting in the construction of the PG. Node features comprise the vectorised FC matrix after RFE (5000 features retained) along with a graphlet counting vector. The latter is computed by counting the number of graphlets (i.e. subgraphs/motifs) with 2 to 5 nodes in the FC matrix (with self-loops removed and thresholding). Edge weights are computed by using demographic features and imaging data. However, instead of computing the similarity based on the FC matrix, they used the graphlet counting vectors. GCN layers were followed by pooling and DropEdge layers. Using the CC200 atlas on the ABIDE dataset with 1035 subjects (539 ASD, 573 TDC), they obtained an accuracy of 64.3%. (ASD-DiagNet [170]: 63.2%).

Cao *et al.* [171] utilised DeepGCN, a 16 layer GCN (Defferrard, PG) to perform ASD classification. To avoid over-smoothing, residual connections and DropEdge were introduced. DropEdge randomly sets the non-zero elements of the adjacency matrix to 0 doing training and the authors view it as a way to augment the dataset. The node features contained the vectorised FC matrix with 2000 features after RFE (via ridge regression). Using the HO111 atlas on the ABIDE dataset with 871 subjects (403 ASD, 468 TDC), they obtained an accuracy of 73.7% (GCN: 70.4%).

TABLE S4

SUMMARY OF FINDINGS FROM ASD STUDIES, INCLUDING THOSE WITH BIOMARKER DISCOVERY PERFORMED (TOP SECTION), WITHOUT BIOMARKER DISCOVERY PERFORMED (MIDDLE SECTION), AND SMALLER DATASETS (BOTTOM SECTION). 'SIZE' REFERS TO THE SIZE OF THE DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL).

Reference	Dataset (Size)	Dataset distribution	Atlas	GNN	Graph	Result	Baseline
[115]	ABIDE (866)	402 ASD, 464 TDC	MODL128	Defferrard	BG	71.8%	k-NN 65.8%
[58]	ABIDE (871)	403 ASD, 468 TDC	HO112	Defferrard	Mix	81.8%	pGCN 71.4%
[116]	ABIDE (949)	419 ASD, 530 TDC	Multiple	Kipf	PG	75.9%	MLP 75.2%
[117]	ABIDE (1112)	-	AAL116	GraphSAGE	BG	72.5%	pGCN 69.7%
[118]	ABIDE (871)	403 ASD, 468 TDC	HO111	Kipf	PG	79.5%	MLP 78.1%
[51]	ABIDE (1057)	525 ASD, 532 TDC	CC200	Edge	BG	70.7%	DNN 70.0%
[119]	ABIDE (871)	403 ASD, 468 TDC	Multiple	Defferrard	Mix	86.3%	pGCN 69.5%
[120]	ABIDE (613) [d]	286 ASD, 327 TDC	HO110	Mix	PG	75.2%	MVGCN 72.0%
[121]	ABIDE (1035) [d]	505 ASD, 530 TDC	CC200	Edge	Mix	73.1%	SVM 64.0%
[122]	ABIDE (867) [d]	416 ASD, 451 TDC	HO110	Defferrard	BG	76.3%	GCN 73.2%
[123]	ABIDE (1007) [M]	481 ASD, 526 TDC	AAL116	New	BG	72.7%	GCN 70.4%
[81]	ABIDE (1007) [M]	481 ASD, 526 TDC	AAL116	New	BG	74.7%	GCN 70.4%
[166]	ABIDE (184)	79 ASD, 105 TDC	AAL116	Kipf	BG	74.7%	-
[59]	ABIDE (871)	403 ASD, 468 TDC	Power264	Kipf	BG	82.3%	DNN 79.2%
[52]	ABIDE NYU (180)	72 ASD, 98 NC	HO110	Cheb	BG	73.5%	DAE 66.0%
[47]	ABIDE (871)	403 ASD, 468 TDC	HO110	GAT	BG	72.4%	PC 65.5%
[167]	ABIDE (184)	79 ASD, 105 TDC	AAL116, CC200	Defferrard	BG	79.5%	GCN 75.8%
[168]	ABIDE (1035)	505 ASD, 530 TDC	CC200	Kipf	BG	66.8%	[31] 57.3%
[169]	ABIDE (1035)	539 ASD, 573 TDC	CC200	Kipf	PG	64.3%	[170] 63.2%
[171]	ABIDE (871)	403 ASD, 468 TDC	HO111	Defferrard	PG	73.7%	GCN 70.4%
[172]	ABIDE (871)	403 ASD, 468 TDC	HO111	Defferrard	PG	70.9%	pGCN 69.5%
[124]	ABIDE (871)	403 ASD, 468 TDC	AAL116	Defferrard	PG	90.6%	EV-GCN 85.9%
[173]	ABIDE (871)	403 ASD, 468 TDC	AAL116, HO112	Snow	PG	67.3%	pGCN 76.4%
[174]	ABIDE (871)	403 ASD, 468 TDC	HO110	Defferrard	BG	65.0%	-
[175]	ABIDE (172)	70 ASD, 102 TDC	Power264	Defferrard	BG	64.0%	[174] 55.0%
[176]	ABIDE (871)	403 ASD, 468 TDC	HO111	Defferrard	PG	75.0%	GCNN 68.0%
[177]	ABIDE (774)	379 ASD, 395 TDC	CC392	Kipf	BG	70.0%	FCN 71.0%
[56]	ABIDE (866)	402 ASD, 464 TDC	AAL116	Defferrard	Mix	73.1%	pGCN 66.4%
[160]	ABIDE (871)	403 ASD, 468 TDC	HO110	SGC	PG	68.6%	Cheb 67.5%
[162]	ABIDE (1029)	485 ASD, 544 TDC	AAL116, BN273	Defferrard	BG	67.3%	TGCN 65.2%
[178]	ABIDE (1029) [d]	485 ASD, 544 TDC	AAL116	Kipf	PG	63.7%	GCN 59.6%
[179]	ABIDE (512) [d]	-	AAL90	Kipf	BG	67.6%	GCN 63.5%
[197]	Biopoint (118)	75 ASD, 43 TDC	DX148	GraphSAGE	BG	70.0%	-
[198]	Biopoint (118)	75 ASD, 43 TDC	DK84	GAT	BG	79.7%	CNN 78.1%
[31]	Biopoint (118)	75 ASD, 43 TDC	DK84	New	BG	79.8%	GAT 77.4%
[200]	ABIDE3 (303)	130 ASD, 173 TDC	SF200	GIN	BG	70.6%	SVM 67.4%
[201]	ABIDE (351)	155 ASD, 196 TDC	AAL116	Kipf	BG	70.0%	GCN 62.0%
[202]	ABIDE NYU (92) [d]	45 ASD, 47 TDC	AAL116	Kipf	BG	79.9%	FCN 72.6%
[203]	ABIDE (144) [d]	70 ASD, 74 TDC	Power264	Kipf	BG	66.0%	SVM 63.8%

Anirudh *et al.* [172] proposed a bootstrapped version of GCN to address the lack of techniques to choose an appropriate graph. This involves using an ensemble of GCNs, each branch having a different percentage of edge features being dropped. Each branch comprises 3 GCN layers (Defferrard, PG) For node features, RFE was used to reduce feature dimensionality of the vectorised FC matrix to 2000. The output (logits) of each GCN is combined via averaging. They experimented with 3 types of graphs: (i) feature-based graph weighted by gender and site, (ii) noisy version of (i) with

edges dropped (set at 30%), (iii) identity graph, i.e. using the imaging features only, without edge features. Using the HO atlas on the ABIDE dataset with 871 subjects, they obtained an accuracy of 70.9% (population GCN: 69.5%) and observed a similar slight increase in performance when the other 2 types of graphs were used.

Mao *et al.* [124] proposed a GCN (Defferrard, PG) based on a relational attention mechanism to address the issue of using a static graph. Attention weights between node i and j are computed by dividing the concatenation of the linearly

projected node embeddings at those nodes with that of the sum of the neighbours of node i . This is followed by a 2-layer MLP. The authors believe that their method generates relations adaptively for different individuals in the population, such that they can learn the unique information of individuals better. Experimenting with 7 different atlases (EZ116, TT97, AAL116, HO111, CC200, CC400, DOS160) on the ABIDE dataset with 871 subjects, they obtained an accuracy of 90.6% (GCN: 69.8%, EV-GCN: 85.9%) with the AAL116 atlas. Performance across atlases only varied by about 1%.

Pan *et al.* [173] proposed MAMF-GCN which combines fMRI and non-imaging data in a PG framework via multiple channels and uses the attention mechanism to fuse information together in an adaptive manner. Node features are generated via Pearson's correlation of the FC matrix. Adjacency matrix for each data modality is formed separately: for imaging data graph, adjacency matrix is based on the cosine similarity function (via k-NN) whereas a pairwise association encoder (PAE) is used in non-imaging (phenotypes such as age, gender, education) data graph. Snowball GCN (akin to DenseNet for GCN) was used to learn embeddings of the original graphs, overcoming the over-smoothing problem that occurs when using multiple layers of GCN. Each modality has its own snowball GCN (modality-specific channel), but a channel common convolution module was also introduced where weights sharing across GCNs are enforced. The attention mechanism is then used to combine all 4 learnt representations. Using features from the AAL116 and HO112 atlas (features from both atlases are used together to form the functional k-NN graph) on the ABIDE dataset with 871 subjects (468 ASD, 403 TDC), MAMF-GCN obtained an accuracy of 97.7% (GCN(PG), 76.4%). Notably, they used DSM information to build the PG. Without it (i.e. using site information, gender, age), accuracy drops to 67.3% but varies greatly depending on the choice of non-imaging data used.

Ktena *et al.* [174] proposed to learn graph similarity metrics using Siamese GCN (s-GCN). The architecture takes in a pair of inputs which goes through the s-GCN (pair of GCN with shared weights, implemented via Defferrard). Dot product is computed between the outputs of s-GCN, which is then passed through a fully connected layer that predicts the similarity between the input pair. s-GCN involves a constrained variance loss that limits the variance for each class to be below a certain threshold. This helps to prevent the issue of the similarity estimates collapsing to the class mean when a variance minimisation approach is adopted. By using a k-NN classifier on the estimated pairwise similarities, classification can be performed. Using the HO110 atlas on the ABIDE dataset with 871 subjects (403 ASD, 468 TDC), they obtained an accuracy of around 65% (barchart presented without exact numbers).

Ma *et al.* [175] proposed an end-to-end similarity learning framework called Higher-order Siamese GCN (HS-GCN) for metric learning of graph datasets. The framework is made of Siamese neural networks and uses two GCNs (Defferrard) as the twin networks. HS-GCN performs higher-order convolutions by incorporating higher-order proximity using random walks into GCNs. In GCNs, filters are defined in the graph spectral domain. FC matrices with only positive values

were used as the graph for the GCN (BG) while connection profile was used as the node feature. Using the Power264 atlas on a subset of the ABIDE dataset with 172 subjects (70 ASD, 102 TDC), they obtained an accuracy of around 64% (S-GCN: 55%). For model explainability, they performed k-means clustering on the brain network embedding by the higher-order GCN.

Jiang *et al.* [56] proposed hierarchical GCN (HI-GCN), which involves a multi-scale scheme and a multi-graph clustering operation to remove noisy and irrelevant connections. Node features representing the connection profiles are passed through a GCN called f-GCN (Defferrard, BG) which also contains EigenPooling layers. and the resulting embedding is used as node features of a PG for another GCN called p-GCN. Notably, they also split positive and negative correlations and go through clustering and graph convolution independently. These are only combined before the fully connected layer. Using the AAL dataset on the ABIDE dataset with 866 subjects (402 ASD, 464 TDC), Hi-GCN obtained an accuracy of 73.1% (GCN: 66.4%).

Studies elaborated in earlier sections include one from Rakhimberdina *et al.* [160] (see Section IV-A for architectural details), which obtained an accuracy of 68.6% (ChebGCN: 67.5%) using the HO110 atlas on the ABIDE dataset with 871 subjects (403 ASD, 468 TDC) via their architecture based on simplified graph convolution (phenotypic information used for PG construction includes gender, age group and site). Yao *et al.*'s MTGCN (see Section IV-A) [162] obtained 67.3% accuracy on the ABIDE I and II dataset with 1029 subjects (485 ASD, 544 TDC).

Other studies include one from Masood *et al.* [176], who proposed a GCNN-LSTM model to model population graphs. The construction of GCNN is based on Parisot *et al.* [24]. Using the ABIDE dataset with 871 subjects (403 ASD, 468 TDC, CPAC, atlas used not mentioned but presumably HO111 following [24]), GCNN-LSTM obtained an accuracy of 75% (GCNN: around 68%, estimated from barplot). Felouat *et al.* [177] used topological metrics and GCNs (Kipf, BG) to create new relevant features of graphs that could be helpful in the classification tasks. Node features included betweenness centrality, clustering coefficient, squares clustering coefficient, closeness centrality, eigenvector centrality and degree centrality. Edges were constructed using edge betweenness centrality, with proportional thresholding performed. Classification is done via a 'voting classifier model', which is essentially an ensemble of models that is combined via majority voting. Using CC400 atlas (392 distinct ROIs) on the ABIDE dataset with 774 subjects (379 ASD, 395 TDC), they obtained an accuracy of 70.0% (FCN: 71.0%).

A dFC study by Peng *et al.* [178] proposed Graph canonical correlation Analysis (CCA) for Temporal Self-supervised Learning (GATE), which involves a two-step GCN learning procedure: (i) similarity-based self-supervised learning on an unlabelled fMRI PG and (ii) fine-tuning on a small labelled fMRI dataset for a classification task. The self-supervised learning strategy involves two new graph augmentation strategies from dFC (step-window augmentation (S-A), multi-scale window augmentation (M-A)) and leveraging CCA to fuse

the different temporal embeddings of the GCN encoder (from original and augmented data), maximising their mutual information. S-A uses neighbouring sliding windows as related views, while M-A uses sliding windows of different sizes. These views are encoded by the same GCN (Kipf, PG) and the weights are optimised via their proposed input-consistency regularisation loss based on CCA. For downstream prediction, the encoder is finetuned (but the adjacency matrix is replaced by an identity matrix) along with a linear layer. Using the AAL116 atlas on ABIDE with 1029 subjects (485 ASD, 544 TDC), GATE obtained an accuracy of 63.7% (GCN: 59.6%)

Liu *et al.* [179] proposed a temporal graph learning framework for brain networks (BrainTGL). dFC matrices from the various time windows are given as input to an attention-based pooling layer that learns a score of all nodes, for each supernode. This score is also used to compute the weight of the superedges, producing a coarsened graph. These outputs are passed to a dual temporal graph learning module that comprises a signal representation learning (S-RL) module and a temporal graph representation learning module (TG-RL). The S-RL module is formed by a stack of convolutional layers and it reduces the dynamic FC data (within a time window) to a vector that is used as node features for the graph. Then, the TG-RL takes in the coarsened graph from each time window and passes them through a modified LSTM with a multi-skip combination that allows it to capture temporal dependencies of varying lengths. The LSTM also contains GCN layers (Kipf, BG). Finally, it is proposed that an ensemble of BrainTGL modules, each with different window lengths, is used to address the issue of finding the optimal window length. Using the AAL90 atlas on a subset of ABIDE (512 subjects with a specific range of sequence length, class distribution unclear), BrainTGL-ensemble obtained an accuracy of 67.6% (GCN: 63.5%, BrainTGL: 65.3%).

C. MDD

Zhu *et al.* [180] utilised a deep graph convolutional neural network (DGCNN) which entails 4 GCN layers (Kipf), a Sort-Pooling layer, a 1D convolution layer with max pooling and dense layers. Binarised adjacency matrices (with a threshold of 0.3) were used as the graph for the GCN (BG) and node features contained the connection profile. Using the DOS160 atlas on the REST-meta-MDD dataset with 1601 subjects (830 MDD, 771 NC), they obtained an accuracy of 72.1% (baseline GCN: 67.4%) Venkatapathy *et al.* [181] proposes an ensemble of GNNs (BG) made up of GCN, GAT and GraphSAGE. A k-NN graph is used as the adjacency matrix and the node features are connection profiles. Using the DOS160 on the REST-meta-MDD dataset with 1586 subjects (821 MDD, 765 NC), the ensemble of GNNs achieved 70.4% accuracy (baseline GCN: 65.2%). Notably, they classified between first-episode and recurrent MDDs patients too, achieving similar accuracies.

Other studies using static FC include MAMF-GCN [173] (see Section IV-B), where their snowball GCN and multi-channel attention architecture, achieved a very high accuracy of 99.2% (GCN: 62.6%, SGC [161]: 79.1%, EV-GCN: 95.9%) on a subset of rest-meta-MDD dataset (site: Southwest University) with 533 subjects (282 MDD, 251 NC). A short paper by

Pitsik *et al.* [182] uses a GCN (BG) with a single GraphConv layer. FC matrix was used as the graph while the node features are one-hot vectors. Using the AAL166 atlas on a private dataset with 91 subjects (41 MDD, 50 NC), they achieved a high accuracy of 98.9% as well as perfect generalisation to unseen dataset. This was further developed in an extended work [183] where more extensive hyperparameter tuning (e.g. thresholding of adjacency matrix) and topological analysis were done. Using the AAL166 atlas on a private dataset with 84 subjects (49 NC, 35 SZ) achieving an accuracy of 93.0% with an architecture of 2 GraphConv layers, with adjacency matrix thresholded to a sparsity of 15% and only positive correlations kept. They also noted that the optimal number of GNN layers is significantly correlated with the length of the shortest path.

A dFC study by Yao *et al.* [184] proposed Temporal-Adaptive GCN (TAGCN) which used both static and dynamic FC for MDD classification. The framework contains two main building blocks: (i) adaptive GCN layer (BG) for constructing flexible FC topology structure that is subject-specific, (ii) temporal convolutional layer to extract dynamic information. A new adjacency matrix is also proposed. It is defined as the sum of three matrices - A , R , and S . A is constructed using k-NN with Pearson's correlation as the similarity metric. R is learnable and optimisable. S is formed from a normalised embedding Gaussian function, wherein it learns topology information of FC in each time-series block. These matrices are then added together and multiplied with the inputs and a weight matrix, producing the output of the GCN. Using the HO112 on a subset of the REST-meta-MDD dataset (site: Southwest University) with 533 subjects (282 MDD, 251 NC), they reported an accuracy of 73.8% (GAT: 70.1%). Ablation studies revealed that R is most important, with accuracy reducing to 71.3% when it was removed.

Other dFC studies include BrainTGL [179] (see Section IV-B) which used the AAL90 atlas on the NMU-MDD dataset with 427 subjects (181 MDD, 246 NC). BrainTGL-ensemble achieved an accuracy of 68.9% (BrainTGL: 66.2%, Brain-NetCNN [19]: 65.8%).

A multi-modal MDD study by Wang *et al.* [185] proposed an adaptive multimodal neuroimage integration (AMNI) framework utilising T1w and fMRI data simultaneously. Intermediate fusion was performed, combining representations of T1w learnt by a 3D CNN with representations of FC learnt from a GCN (Kipf, BG). Both sets of features are then fused via concatenation, but it is noted that feature adaptation is performed by introducing a maximum mean discrepancy loss to the overall loss function. Using the HO112 atlas on a subset of the REST-meta-MDD dataset (site: Southwest University) with 533 subjects (282 MDD, 251 NC), they reported an accuracy of 65%.

D. SZ

2 sFC studies did not perform biomarker discovery. Chen *et al.* [132] adopted a GCN framework proposed by Lee *et al.* [186] which involves 3 blocks of GCN layer (Kipf) with top-k pooling as well as a multi-scale readout layer that takes in the

TABLE S5

SUMMARY OF FINDINGS FROM MDD STUDIES, INCLUDING THOSE WITH BIOMARKER DISCOVERY PERFORMED (PLACED ABOVE MIDLINE). 'SIZE' REFERS TO SIZE OF DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL).

Reference	Dataset (Size)	Dataset distribution	Atlas	GNN	Graph	Result	Baseline
[82]	REST-meta-MDD (1586)	821 MDD, 765 NC	DOS160	Defferrard	BG	81.5%	-
[30]	REST-meta-MDD, psymri (2498)	1249 MDD, 1249 NC	HO112, CC200	Kipf	BG	61.3%	SVM-RBF 61.2%
[125]	Private (218)	129 MDD, 89 NC	BM82, JHU81	Kipf	BG	70.9%	GAT 68.2%
[126]	Private (75)	29 MDD, 44 NC	Yeo114	Defferrard	PG	74.1%	SVM 69.8%
[127]	Private (277) [d]	180 MDD, 97 NC	Unclear	GAT	BG	84.0%	-
[128]	REST-meta-MDD (681) [d]	356 MDD, 325 NC	AAL116	Kipf	BG	59.3%	STNet 52.0%
[180]	REST-meta-MDD (1601)	830 MDD, 771 NC	DOS160	Kipf	BG	72.1%	GCN 67.4%
[181]	REST-meta-MDD (1586)	821 MDD, 765 NC	DOS160	Mix	BG	70.4%	GCN 65.2%
[173]	REST-meta-MDD SWU (533)	282 MDD, 251 NC	AAL116, HO112	Snowball	PG	99.2%	SGC 79.1%
[182]	Private (91)	41 MDD, 50 NC	AAL166	GraphConv	BG	98.9%	-
[183]	Private (84)	35 MDD, 49 NC	AAL166	GraphConv	BG	93.0%	-
[184]	REST-meta-MDD SWU (533) [d]	282 MDD, 251 NC	HO112	New	BG	73.8%	GAT 70.1%
[179]	NMU (427) [d]	181 MDD, 246 NC	AAL90	Kipp	BG	68.9%	CNN 65.8%
[185]	REST-meta-MDD SWU (533) [M]	282 MDD, 251 NC	HO112	Kipf	BG	65.0%	-

TABLE S6

SUMMARY OF FINDINGS FROM SZ STUDIES, INCLUDING THOSE WITH BIOMARKER DISCOVERY PERFORMED (PLACED ABOVE MIDLINE). 'SIZE' REFERS TO THE SIZE OF THE DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL).

Reference	Dataset (Size)	Dataset distribution	Atlas	GNN	Graph	Result	Baseline
[129]	Private (1412)	505 SZ, 907 NC	AAL116, DOS160	Defferrard	BG	85.8%	SVM 80.9%
[131]	Private (345) [M]	140 SZ, 205 NC	AAL90, BNA246	Kipf	BG	95.8%	SVM 81.2%
[133]	COBRE (154) [M]	67 SZ, 87 NC	DK293	Kipf	BG	75.0%	SVM 71.0%
[132]	Private (345)	140 SZ, 205 NC	AAL90	Kipf	BG	92.7%	MLP 75.0%
[160]	COBRE (145)	71 SZ, 74 NC	HO110	SGC	PG	80.6%	Cheb 76.5%
[187]	COBRE (120) [d]	-	AAL90	Kipf	BG	91.6%	ST-CRN 88.3%
[164]	UCLA (272) [M]	50 SZ, 115 NC	AAL90	Defferrard	BG	93.7%	SVM 90.2%

representations from each pooling layer. In the readout layer at each scale, max-pooling and mean-pooling are performed and concatenated. These concatenated representations are then summed up before passing it to an MLP for classification. The graph for GCN (BG) was a thresholded FC matrix (edge is set as 1 if above the threshold, else 0). Node features included connectivity features (ALFF, ReHo and DC) and network properties (nodal efficiency, betweenness centrality). Using features from AAL90 on a private dataset with 345 subjects (140 SZ, 205 NC), they found that connectivity

features (92.7%) did better than centrality measures (73.3%) (MLP: 75.0%).

Other studies applied on multiple neurological disorders (where SZ is one of them) include Rakhimberdina *et al.* [160] (see Section IV-A for architectural details), which obtained an accuracy of 80.6% (ChebGCN: 76.5%) using the HO110 atlas on the COBRE dataset with 147 subjects (71 SZ, 74 NC) via their architecture based on simplified graph convolution (phenotypic information used for PG construction includes gender, age group and handedness).

1 study looked into the use of dFC for SZ prediction. Huang *et al.* [187] proposed a framework that consists of two main parts: graph construction and a hierarchical representation learning (HARL) module. In graph construction, they use a 1D convolution operator to learn the weight of the edges between states of the dynamic brain network, assuming that each state may have potential connections with other states. HARL consists of three levels, each level containing a graph convolutional layer (Kipf), a graph pooling layer (top-k based on self-attention scores) and a readout layer (where the node representation is concatenated with the maximum out of all neighbours, then summed up across all nodes). The output of the readout layer from each level is combined via concatenation and passed to a fully connected layer to produce the final prediction. Graph used for the GNN is constructed based on node features with the 1D convolution described above. Node features are based on the FC matrix from each sliding window. Using the AAL90 atlas on the COBRE dataset with 120 subjects (class distribution not mentioned), they obtained an accuracy of 91.6% (ST-CRN: 88.3).

1 multimodal study on SZ was conducted by Liu *et al.* [164] (refer to Section IV-A for results). However, their disease class is formed by combining ADHD, SZ and bipolar disorder patient population and no SZ-specific insights were reported.

E. Dementia

Zhao *et al.* [188] used a GCN (Cheb, PG) for MCI prediction. They constructed the graph using different scanner information and gender. These were then fed into a GCN network, which contains graph convolution and pooling. Using the AAL90 atlas on the ADNI-2 and ADNI-3 dataset with 184 subjects (40 LMCI, 77 EMCI, 67 NC), they obtained an accuracy of 78.4% for EMCI vs NC classification (GCN: 72.6%) and 84.3% for LMCI vs NC classification (GCN: 81.6%).

Inspired by U-Net, Qin *et al.* [189] introduced a U-shaped hierarchical GCN framework (U-GCN, based on BG), which includes down-sampling, up-sampling and skip connection operators for graph data. Node features are assigned through graph theoretical metrics such as clustering coefficient and nodal efficiency. An adjacency matrix thresholded at 0.5 was used as the graph. Using an unspecified atlas on the ADNI dataset with 91 subjects (44 AD, 47 NC), U-GCN obtained an accuracy of 83.3% (Higher-order GNN: 72.2%).

Other studies include Jiang *et al.*'s [56] Hi-GCN (see Section IV-B) achieved an accuracy of 75.6% (pGCN: 73.7%) for AD vs MCI classification on the ADNI dataset with 133 subjects (99 MCI, 34 AD). Li *et al.*'s [119] TE-HI-GCN (see Section IV-B) which obtained an accuracy of 89.4% (GCN: 76.5%) for AD vs MCI classification on the same dataset. Unlike their study on ASD, model explainability was not performed for AD classification. A study by Wang *et al.* [166] which experimented with various node vectors while using the AAL116 atlas on a subset of the ADNI dataset with 137 subjects (69 NC, 68 MCI), obtained an accuracy of 78.5% when one-hot encoding and node correlation were used as node features (76.7% when only node correlation was used).

An *et al.* [190] used a GCN (Kipf, PG) where node vectors are features obtained from feature selection. They used features from ALFF, ReHo, FC and dFC (sliding window with thresholding). 4 ALFF and 5 ReHo features that survived feature selection via t-test were kept. For dFC, feature selection from RF to obtain 10 features then further performed RFE-SVM, which kept 9 features. In total, 18 handcrafted features were used. Graph construction is not clearly specified. Using AAL116 atlas on a private dataset with 246 (98 AD, 148 NC) scans from 204 patients (containing 94 AD and 110 NC subjects), they obtained an accuracy of 91.3%.

Song *et al.* [191] introduced a novel method of constructing the graph using both static and dynamic FC. A high-order dFC matrix is constructed by multiplying the dFC matrix with its transpose. Then, all dFC matrices are combined with the sFC matrix by averaging across all dFC matrices and taking a weighted average between the sFC and mean dFC matrix. Their proposed GCN model (Defferrard, PG) consists of two graph convolution layers with ReLU activations and one softmax layer. Node features are extracted from the combined sFC/dFC matrix via RFE and the edges are constructed using demographic information. Using the AAL90 atlas on the ADNI dataset with 184 subjects (40 LMCI, 77 EMCI, 67 NC), their model achieves an accuracy of 82.7% for EMCI vs NC classification (GCN: 65.5%) and 88.7% for LMCI vs NC classification (GCN: 87.8%).

A notable study from Li *et al.* [153] took a different approach from other papers: instead of using diagnostic labels, they looked across the spectrum of AD by deriving groups of subjects containing similar patterns of $A\beta$ deposition. They defined graphs with both static and dynamic FC whereby the A was denoted by the static FC between the regions and the X was formed as the one-hot representation for the ROI index. By using their obtained clusters (which correspond to AD, MCI and NC) for different AD phenotypes, the authors train GCN models [66] on both static and dynamic FC (separately and together) and show that models trained over both static and dynamic FC outperformed the baselines. Using the Shen268 atlas on a subset of the OASIS-3 dataset with 258 samples (227 NC, 14 MCI, 17 AD; note that these diagnostic labels are not used), they achieved an accuracy of 78.8% (sFC only: 70.2%, dFC only: 72.7%). This suggests that the diagnostic labels might not be as reliable as thought.

Liu *et al.* [192] proposed Multiscale Atlas-based GCN (MAGCN), which integrated information contained in the multiple spatial-scale dFCNs (Defferrard). Each scale has a GCN which takes in the outputs of the large scale (after an atlas mapping operation). These outputs are concatenated and passed to an LSTM, whose outputs are then subjected to an attention-based multiple instance learning pooling layer, followed by a fully-connected layer to classify the disease status of the subject. Using the SC200 atlas on ADNI-2 and ADNI-GO datasets with 481 subjects (117 EMCI, 364 NC), MAGCN achieved an accuracy of 77.8% for EMCI classification.

Hu *et al.* [193] proposed an architecture that simultaneously performs node and graph classification for dFC data (based on sliding windows) on datasets involving AD and frontotem-

TABLE S7

SUMMARY OF FINDINGS FROM DEMENTIA STUDIES, INCLUDING THOSE WITH BIOMARKER DISCOVERY PERFORMED (TOP SECTION), WITHOUT BIOMARKER DISCOVERY PERFORMED (MIDDLE SECTION), AND DEMENTIA SUBTYPES/STAGES WITH FEW STUDIES (BOTTOM SECTION). ‘SIZE’ REFERS TO THE SIZE OF THE DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL).

Reference	Dataset (Size)	Dataset distribution	Atlas	GNN	Graph	Result	Baseline
[134]	ADNI3 (168) [d]	82 SMC, 86 NC	AAL90	Kipf	BG	87.5%	SVM 87.5%
[141]	ADNI (138) [d]	44 SMC, 94 NC	AAL90	Kipf	PG	83.2%	GCN 81.1%
[143]	ADNI & Private (207) [M]	44 SMC, 163 NC	AAL90	Kipf	PG	93.2%	GCN 86.5%
[26]	ADNI (88) [M]	44 SMC, 44 NC	AAL90	Defferrard	PG	84.1%	GCN 76.1%
[61]	ADNI (910)	345 EMCI, 565 NC	SC200	Kipf	BG	73.4%	GCN 66.9%
[135]	ADNI (101) [d]	53 EMCI, 48 NC	YEO114	Kipf	BG	74.4%	CNN 67.6%
[140]	ADNI (88) [d]	44 EMCI, 44 NC	AAL90	Defferrard	PG	87.5%	GCN 84.2%
[141]	ADNI (180) [d]	86 EMCI, 94 NC	AAL90	Kipf	PG	80.0%	GCN 75.6%
[142]	ADNI (154) [M]	77 EMCI, 67 NC	AAL90	Defferrard	PG	85.4%	GCN 71.5%
[143]	ADNI & Private (249) [M]	86 EMCI, 163 NC	AAL90	Kipf	PG	91.2%	GCN 85.5%
[26]	ADNI (88) [M]	44 EMCI, 44 NC	AAL90	Defferrard	PG	85.2%	GCN 75.0%
[140]	ADNI (82) [d]	38 LMCI, 44 NC	AAL90	Defferrard	PG	89.0%	GCN 79.6%
[142]	ADNI (107) [M]	40 LMCI, 67 NC	AAL90	Defferrard	PG	93.5%	GCN 73.8%
[143]	ADNI & Private (329) [M]	166 LMCI, 163 NC	AAL90	Kipf	PG	94.2%	GCN 87.5%
[26]	ADNI (82) [M]	38 LMCI, 44 NC	AAL90	Defferrard	PG	89.0%	GCN 80.7%
[136]	ADNI (83)	33 AD, 50 NC	AAL90	New	BG	94.2%	-
[138]	ADNI2 (292) [d]	115 AD, 177 NC	AAL116	Defferrard	BG	90.0%	GCN 83.8%
[139]	ADNI (107) [d]	59 AD, 48 NC	AAL90	Kipf	BG	89.8%	GCN 87.5%
[188]	ADNI2, ADNI3 (144)	77 EMCI, 67 NC	AAL90	Defferrard	PG	78.4%	GCN 72.6%
[188]	ADNI2, ADNI3 (107)	40 LMCI, 67 NC	AAL90	Defferrard	PG	84.3%	GCN 81.6%
[189]	ADNI (91)	44 AD, 47 NC	-	New	BG	83.3%	HO-GNN 72.2%
[56]	ADNI (133)	99 MCI, 34 AD	AAL116	Defferrard	Mix	75.6%	pGCN 73.7%
[119]	ADNI (133)	99 MCI, 34 AD	Multiple	Defferrard	Mix	89.4%	GCN 76.5%
[166]	ADNI (137)	68 MCI, 69 NC	AAL116	Kipf	BG	78.5%	Ablation 76.7%
[190]	Private (246) [d]	98 AD, 148 NC	AAL116	Kipf	PG	91.3%	-
[191]	ADNI (144) [d]	77 EMCI, 67 NC	AAL90	Defferrard	PG	82.7%	GCN 65.5%
[191]	ADNI (107) [d]	40 LMCI, 67 NC	AAL90	Defferrard	PG	88.7%	GCN 87.8%
[153]	OASIS3 (258) [d]	-	SHEN268	Kipf	BG	78.8%	Ablation 72.7%
[192]	ADNI2, GO (481) [d]	117 EMCI, 364 NC	SC300	Defferrard	BG	77.8%	-
[193]	FTD (181) [d]	95 FTD, 86 NC	AAL90	Kipf	Mix	78.2%	ST-GCN 73.8%
[193]	ADNI (107) [d]	59 AD, 48 NC	AAL90	Kipf	Mix	72.7%	ST-GCN 71.8%
[163]	ADNI (370) [d]	191 MCI, 179 NC	Multiple	Defferrard	BG	86.0%	TGCN 83.8%
[178]	FTD (181) [d]	95 FTD, 86 NC	AAL116	Kipf	PG	72.4%	GCN 64.5%
[138]	ADNI2 (368) [d]	191 EMCI, 177 NC	AAL116	Defferrard	BG	71.6%	GCN 70.5%
[194]	ADNI (210) [M]	105 EMCI, 105 NC	AAL90	Defferrard	PG	84.1%	Ablation 78.5%
[196]	ADNI3 (209) [M]	93 MCI, 116 NC	DX146	Kipf	BG	97.7%	-
[204]	Private (140)	59 aMCI, 81 NCI	SC500	Kipf	BG	80.4%	GCN 73.6%
[204]	Private (138)	57 naMCI, 81 NCI	SC500	Kipf	BG	79.0%	GCN 72.0%
[139]	FTD (181) [d]	95 FTD, 86 NC	AAL90	Kipf	BG	87.2%	GCN 83.1%
[58]	ADNI (221)	121 MCI, 100 NC	HO112	Defferrard	Mix	76.5%	HI-GCN 71.1%
[205]	ADNI (292)	100 MCI, 192 NC	AAL90	New	BG	78.6%	GCN 72.8%
[93]	Private (144) [d]	69 SCD, 75 NC	AAL90	Kipf	BG	72.4%	-

poral dementia (FTD). The first branch (node classification) vectorises the FC matrix to use as input data to a GCN with 2 graph convolution layers (Kipf, PG). Adjacency matrix is formed via k-NN based on vectorised FC matrix only instead of using demographic data. However, they also introduce the attention mechanism to improve the adjacency matrix at every training iteration. The second branch (graph classification) uses the functional profile as input to the GCN, also with 2 graph convolution layers (Kipf, BG) and the adjacency matrix is formed via k-NN based on the connection profile. This

is followed by a readout layer that computes the average representation across all sliding windows. This representation is passed through another graph convolution layer but a key difference in this layer is that the adjacency matrix is formed from the average of adjacency matrices obtained in the first branch. Finally, the representations from both branches are projected by another learnable matrix before being fused via concatenation. Using the AAL90 atlas on two datasets (FTD with 181 subjects: 95 FTD, 86 NC; ADNI with 107 subjects: 59 AD, 48 NC), they obtained an accuracy of 78.2% for FTD

classification (ST-GCN: 73.8%, Hi-GCN: 76.7%) and 72.7% for AD classification (ST-GCN: 71.8%, Hi-GCN: 70.8%). Ablation studies reveal that removing the temporal information results in a drop in accuracy of around 3% for both tasks.

Other dFC studies without biomarker discovery include Yao *et al.*'s MMTGCN [163] (see Section IV-A) which achieves an accuracy of 86.0% for MCI classification on ADNI dataset with 370 (191 MCI, 179 NC) while using multiple atlases. GATE proposed by Peng *et al.* [178] used the AAL116 atlas on the FTD dataset of 181 subjects (95 FTD, 86 NC) and obtained an accuracy of 72.4% (GCN: 64.5%). DS-GCN proposed by Xing *et al.* [138] was discussed in the main manuscript for their experiments on AD, but they also performed EMCI experiments without highlighting biomarkers. Using the AAL90 atlas on the ADNI dataset with 368 subjects (191 EMCI, 177 NC), they obtained an accuracy of 71.6% (GCN 70.5%).

2 studies with biomarker discovery using multi-modal datasets of dementia were found. Liu *et al.* [194] utilised a spectral GCN (Defferrard, PG) where node features comprise GM volumes from T1w images and nodal shortest path length from fMRI (not FC). Feature selection was performed using group LASSO, specifically MTFs-gLASSO [195] (done individually for each modality, then concatenated). The graph was built based on correlation of node vectors, along with age, gender and MMSE scores. Using the AAL90 atlas on the ADNI dataset with 210 subjects (105 EMCI, 105 NC), they achieved an accuracy of 84.1%. Ablation studies showed that the best performance was obtained by combining imaging and non-imaging phenotypic features and group LASSO did better than vanilla LASSO or feature selection via t-test. Other multimodal studies include one from Zhang *et al.* [196], who proposed a Graph Convolutional Recurrent Neural Network (GCRNN), which incorporated DTI and dFC into the analysis. SC is used as the graph for GCN (Kipf) to capture spatial relationships. These embeddings are passed to an RNN (they experimented with both LSTM and GRU) to capture temporal patterns of fMRI data. Using the DX148 (146 ROIs after removing 2) on the ADNI-3 dataset with 209 subjects (93 MCI, 116 NC), they obtained an accuracy of 97.7% for NC/MCI classification.

F. PD

Chan *et al.* [148] proposed an architecture called JOIN-GCLA that takes in both multimodal neuroimaging and multi-omics datasets. Since many data modalities are considered, they proposed a scalable PG approach involving multiple GCNs. Multimodal neuroimaging features are first concatenated and used as inputs to a GCN. The graph is built based on the similarity of the node features, as measured by Pearson's correlation. The output of the GCN is passed to multiple branches of GCNs, each branch using a different adjacency matrix that is built based on similarity of the omics features (also measured via Pearson's correlation). Finally, an attention mechanism is used to combine the intermediate prediction and produce the final prediction. Using features from the AAL116 atlas on the PPMI dataset (376 original scans - 351 PD, 25 NC - with significant imbalance addressed via CycleGAN),

their JOIN-GCLA architecture produces Matthew correlation coefficients above 0.8 over several different combinations of imaging and omics modalities. However, they noted that data imbalance and the small size of the test dataset are limiting factors in the study and further studies on larger datasets should be conducted. While connectivity features were not analysed, attention scores revealed that DNA methylation and single nucleotide polymorphism were the most important omics modalities.

VII. STUDIES WITH BIOMARKER DISCOVERY

A. ASD

In this section, we summarise ASD studies that performed biomarker discovery, but were not included in the main manuscript due to space constraints. These 7 studies have smaller datasets and they are summarised in Table S9.

Li *et al.* proposed a series of techniques to identify salient features from the Biopoint dataset (task-fMRI) containing 118 subjects (75 ASD, 43 TDC). In [197], they proposed the use of Infomax loss to embed informative and robust brain regional fMRI representations for both graph-level classification and region-level functional difference detection tasks. The pipeline involves a classifier and discriminator that take in representations from a GNN encoder (GraphSAGE, BG). The classifier and discriminator are optimised simultaneously. Node features are handcrafted from features such as degree, GLM coefficients, mean and standard deviation of task-fMRI, and ROI centre coordinates. Using the DX148 atlas, their model achieved the highest F-score of 0.70. For model explainability, they extracted the representations learnt by the last graph convolutional layer and used t-SNE to visualise the node presentations in 2D space. These gave discriminatory regions, from which regions with **Silhouette score** greater than 0.1 included prefrontal cortex, cingulate cortex, visual regions, and other social, emotional and execution related brain regions. In [198], they proposed Pooling Regularized-GNN (PR-GNN), which comprises multiple blocks of node convolutional (with edge features) and node pooling (ranking-based, e.g. TopK and SAGE) layer. Graph convolutions are performed via GAT. Various losses (distance-based loss, e.g. maximum mean discrepancy; group-level consistency loss) were introduced to influence the pooling layer such that ranking scores have sufficient variability to separate the nodes well, and to ensure that nodes linked to the same class are ranked similarly. Using the DK84 atlas, PR-GNN achieved an accuracy of 79.7% (BrainNetCNN: 78.1%). Salient ROIs extracted from the proposed **graph pooling** layers include dorsal striatum, thalamus and frontal gyrus.

Finally, brain functional network is known for its modular organisation and the modulations to this modular organisation during different diseases/tasks [199]. In [31], they proposed to include modular information by introducing an ROI-aware graph convolutional layer that treats the brain nodes as communities and trains different weights for nodes belonging to different communities. They then perform node-pooling for reducing the network dimensions and finally feature concatenation. Node features were the connection profile and the

TABLE S8

SUMMARY OF FINDINGS FROM PD STUDIES, INCLUDING THOSE WITH BIOMARKER DISCOVERY PERFORMED (PLACED ABOVE MIDLINE). ‘SIZE’ REFERS TO THE SIZE OF THE DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL).

Reference	Dataset (Size)	Dataset distribution	Atlas	GNN	Graph	Result	Baseline
[144]	Private (150) [M]	75 PD, 75 NC	DK86	GAT	BG	73.0%	MLP 66.0%
[146]	PPMI (41) [M]	22 PD, 19 NC	BNA246	Kipf	BG	80.5%	MLP 58.5%
[148]	PPMI (376) [M]	351 PD, 25 NC	AAL116	Kipf	PG	>0.8 MCC	-

TABLE S9

SUMMARY OF FINDINGS FROM STUDIES THAT IDENTIFIED POTENTIAL BIOMARKERS. ‘SIZE’ REFERS TO THE SIZE OF THE DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL). ‘TYPE’ REFERS TO THE TYPE OF EXPLANATION.

Reference	Dataset (Size)	Atlas	Explainer (Type)	Salient features
ASD				
[197]	Biopoint (118)	DX148	Clustering (ROI)	Regions: PFC, cingulate cortex
[198]	Biopoint (118)	DK84	Pooling (ROI)	Regions: dorsal striatum, thalamus, frontal gyrus
[31]	Biopoint (118)	DK84	Pooling (ROI)	Regions: frontal gyrus, temporal lobe, cingulate gyrus, occipital pole, angular gyrus
[200]	ABIDE (303)	SF200	Pooling (ROI)	Regions: right parietal, left visual, right lateral PFC, left PFC, left cingulate
[201]	ABIDE (351)	AAL116	Attention (ROI)	Regions: hippocampus, PARAH, putamen, thalamus
[202]	ABIDE (92) [d]	AAL116	Pooling (ROI)	Regions (Lo): bilateral cerebellum and right hippocampus ; Regions (Ho): left insula, left putamen, medial aspect of right SFG
[203]	ABIDE (144) [d]	Power264	Clustering (module)	Modules: ASD stronger FC in visual, DMN and SN ; TDC higher FC in sensory and auditory networks.

graph was a thresholded FC matrix computed via partial correlations. Using the DK84 atlas, they achieved an accuracy of 79.8% (GAT: 77.4%). Model interpretability via **graph pooling** operation revealed group-level biomarkers such as regions in the frontal gyrus, temporal lobe, cingulate gyrus, occipital pole, and angular gyrus. Notably, BrainGNN is able to produce both individual-level and group-level biomarkers.

Yang et al. [200] utilised a 4-layer GIN architecture (with attention-based graph pooling performed at each layer) for ASD classification. The graph of the GNN was the thresholded FC matrix (BG) while node features are one-hot encodings. Using the SF200 atlas on a subset of ABIDE I and II with 303 subjects (130 ASD, 173 TDC), they obtained an accuracy of 70.6% (SVM: 67.4%). Model explainability via the **pooling** layer revealed salient ROIs (above 75th percentile across more than half of the experiments) were present in all functional networks (Yeo-7) but dominated by the visual and frontoparietal control networks.

Chu et al. [201] focused on the adaptation of images between different sites via an unsupervised domain adaptation module called A²GCN. They define a labelled source and an unlabelled target domain with the data coming from the train and test sets, respectively. The loss function contains terms to reduce the loss between the embeddings obtained from the source and the target domains, and the label cross entropy. The X and A are defined by Pearson’s Correlation and the GNN is the spectral GCN defined in [66]. They also define an attention module (before the domain adaptation

module) which they use to get salient nodes in the brain. Using the AAL116 atlas on a subset of the ABIDE dataset (NYU, UM, UCLA) with 351 subjects (155 ASD, 196 TDC), they achieved accuracies ranging from 68.7% to 72.3% (GCN: 61.1% to 63.1%) depending on the choice of source and target site. Ablation studies demonstrated that the adaption module contributed most greatly to model performance. For model explainability, they focused on site UM and extracted the representations produced by the attention module. 19 features with ‘strong correlations’ were chosen, but it was unclear how these were computed. Regions included the hippocampus, parahippocampal gyrus, putamen and thalamus.

Zhao et al. [202] proposed self-attention mechanism GCN (SA-GCN) which uses SAGPool in between GCN layers (Kipf, BG) to reduce parameters. They proposed using both low-order and high-order functional graph networks (LO-FGN, HO-FGN). The former refers to the usual FC, but with thresholding performed. The latter involves computing the correlation of the FC matrix, followed by thresholding. Node vector used for LO-FGN is the mean and variance of the (subset of) time series at that ROI, while HO-FGN uses the FC matrix. The graph embeddings of each FGN are combined after the SA-GCN layer, but the fusion technique is unclear. They used sliding windows to generate FC matrices, which they view as a form of data augmentation as they did not model the temporal relationship between them. Using the AAL116 atlas on a subset of ABIDE (site NYU) with 92 subjects (45 ASD, 47 TDC), they obtained an accuracy of 79.9% (GCN.Lo:

68.0%, GCN_Ho: 74.5%, FCN 72.6%). However, they noted that the performance was rather sensitive to the choice of parameters such as window width, step size and threshold. Based on their **pooling** approach, the probability of occurrence of the ROI after pooling was used as feature scores. Lo-FGN and Ho-FGN have different biomarkers: the former highlighted the bilateral cerebellum (index 8, 103rd and 104th ROI in AAL116) and right hippocampus, while the latter highlighted left insula, left putamen and the medial aspect of the right superior frontal gyrus. The latter is largely linked to the DMN.

Noman et al. [203] proposed an architecture based on dynamic graph autoencoders (DyGAE) to model dFC data produced via sliding windows. DyGAE involves a GAE with 2 GCN layers (Kipf, BG) in the encoder to extract representations of dFC data. GAE reproduces the adjacency matrix A for each timepoint. This matrix is a thresholded and binarised version of the original dFC matrix. DyGAE is followed by a readout layer which takes in the latent embeddings from DyGAE at each timepoint. They proposed 2 options for downstream classification: (i) fully connected layers that do weighted voting for all windows for each subject, (ii) DyGAE-LSTM. Using Power264 atlas on a subset of ABIDE I (3 sites: NYU, UCLA, Utah) with 144 subjects (70 ASD, 74 TDC), they obtained an accuracy of 66.0% for the weighted voting version of DyGAE (DyGAE-LSTM: 54.8%, BrainNetCNN: 54.5%, SVM: 63.8%, Hi-GCN (sFC): 66.6%). For model explainability, they constructed higher order dFC and applied K-means **clustering**. They obtained 3 distinct states and showed that ASD has stronger FC connections in visual, DMN and salience networks while TDC has higher connectivity in sensory and auditory networks.

B. Dementia

In this section, we summarise dementia studies that performed biomarker discovery, but were not included in the main manuscript due to space constraints. These 6 studies cover subtypes/stages of dementia that do not have sufficient studies. They are summarised in Table S10.

Liu et al. [204] proposed a multiscale connectome for vascular dementia classification. This entails multiple layers of GCNs (Kipf, BG) with residual connections (akin to DenseNet). Each GCN uses a different parcellation, starting from finer ones (500 ROIs) to coarser ones (100 ROIs). Between each pair of GCN (with atlases of different scales), an atlas-pooling layer combines ROIs based on an overlapping ratio threshold (set at 0). Learnt representations from each GCN are fused via concatenation before passing them to fully connected layers for classification. Using Schaefer's parcellations with a range of 100 to 500 ROIs, on a private dataset with 197 subjects (59 amnesic MCI, 57 non-amnesic MCI, 81 NCI; NCI refers to presence of subcortical vascular disease without cognitive impairment), they achieved an accuracy of 80.4% for NC vs aMCI classification (GCN, 200 ROI: 73.6%), 79.0% for NC vs naMCI classification (GCN, 100 ROI: 72.0%). **Grad-CAM** was used for model explainability, revealing that the limbic and DMN networks are consistently most crucial for aMCI vs NCI classification.

naMCI vs NCI classification has considerably more variations in crucial networks across ROI scales.

Wang *et al.* [139] proposed a GNN architecture that incorporates self-attention to adaptively learn the adjacency matrix used by the GCN (details in main manuscript). Using the AAL90 atlas on the Frontotemporal Dementia (FTD) dataset (from the Frontotemporal Lobar Degeneration Neuroimaging Initiative (NIFDI)) with 181 subjects (95 FTD, 86 NC), they obtained an accuracy of 87.2% (SGC [161]: 83.1%). Model explainability was achieved via **attention** scores. Regions in the frontal and temporal regions of the brain were highlighted.

Zhang et al.'s [58] LG-GNN (see Section IV-B for details) was also used to analyse AD/MCI subjects. Using the HO112 atlas on the ADNI dataset with 221 subjects (121 MCI, 100 AD), they obtained an accuracy of 76.5% (Hi-GCN: 71.1%) They noted that brain regions with the largest **attention** weights are mainly in the hippocampus, superior parietal lobule, temporal gyrus, inferior frontal gyrus and insular cortex. However, it is unclear which experiment these insights were derived from.

Dong et al. [205] proposed a Co-opetition Hypergraph Graph Neural Network (COOP-DHGNN) for MCI classification and prediction of FC trajectory that works even with a sparse trajectory dataset (i.e. timepoints with missing FC data). Drawing inspiration from generative adversarial networks, COOP-DHGNN comprises a conditional discriminator, conditional generator and a classifier (based on hypergraphs) trained in a collaborative and competitive manner. The classifier is first pre-trained using cross entropy loss. They highlighted that existing GCNs do not make effective use of edge features and proposed a Dual Hypergraph Module (DHM) to combine node and edge feature embeddings. This involves a hypergraph convolution, which necessitates a Dual Hypergraph Transformation (DHT) where the incidence matrix (where each row shows whether the corresponding node is incident to the hyperedge) is transposed. HyperDrop is akin to a top-k pooling layer but performs edge-wide pooling instead of nodes. Subsequently, the generator and discriminator are trained and the pre-trained classifier is only trained after a certain number of epochs. The generator and discriminator are trained with 3 losses: (i) a GAN loss, (ii) collaborative loss and (iii) auxiliary loss. GAN loss is similar to the adversarial loss in GAN, while collaborative loss encourages the classifier to classify synthetic data from the generator correctly. The auxiliary loss helps the generator learn better by incorporating a mean squared error loss and a Pearson's correlation loss. The trained generator is used to reconstruct FC matrices at timepoints where data is missing. Using the AAL90 atlas on the ADNI dataset with 292 subjects (100 MCI, 192 NC) with varying number of timepoints, they achieved an accuracy of 78.6% (GCN: 72.8%). Model interpretability is enabled via the **HyperDrop** layer, revealing that connections between the orbital part of the left inferior frontal gyrus with the right superior temporal pole and the left middle temporal gyrus to be the most salient connections.

Zhang et al. [93] used a GNN-based architecture involving GCNConv (BG, Kipf) and SAGPool layers for the classification of patients with subjective cognitive decline (SCD).

TABLE S10

SUMMARY OF FINDINGS FROM STUDIES THAT IDENTIFIED POTENTIAL BIOMARKERS. ‘SIZE’ REFERS TO THE SIZE OF THE DATASET. WHEN MODALITIES BEYOND SFC ARE USED, THEY ARE MARKED WITH [D] (DFC) OR [M] (MULTIMODAL). ‘TYPE’ REFERS TO THE TYPE OF EXPLANATION.

Reference	Dataset (Size)	Atlas	Explainer (Type)	Salient features
Dementia (amnesic MCI)				
[204]	Private (140) (amnesic)	SC500	GradCAM (module)	Modules: limbic network and DMN
[204]	Private (138) (non-amnesic)	SC500	GradCAM (module)	Modules: more variations in crucial networks
Dementia (Frontotemporal)				
[139]	FTD (181) [d]	AAL90	Attention (ROI)	Regions: amygdala, precentral gyrus, PARAH
Dementia (MCI)				
[58]	ADNI (221)	HO112	Attention (ROI)	Regions: hippocampus, SPL, temporal gyrus, IFG, insular cortex
[205]	ADNI (292)	AAL90	Pooling (connection)	Connections: between the orbital part of the left IFG with the right STP and the left MTG
Dementia (subjective cognitive decline)				
[93]	Private (144) [d]	AAL90	Multiple (ROI)	Regions: all: left thalamus; SVM and GCN: left angular, left supramarginal and bilateral middle cingulum regions

The average value of fMRI time series was used as the node vector, while FC was used as the graph. Using the AAL90 atlas on a private dataset with 144 subjects (69 SCD, 75 NC), they achieved an accuracy of 72.4%. Detailed biomarker analysis was conducted via 3 methods: statistical tests, SVM and GCN (**via attention**). The left thalamus was identified as a salient region across all 3 methods, while the left angular, left supramarginal and bilateral middle cingulum regions were salient in both GCN and SVM.

ADDITIONAL DISCUSSION - CHALLENGES

Dataset composition

With the increasing use of datasets aggregated from multiple sites, more care should be taken when choosing the sites to be used for the analysis especially when the goal is to perform biomarker discovery. For instance, such datasets are often gathered from different countries and poor generalisation performance does not immediately mean that fMRI is not useful for identifying patients with the disorder. Rather, it could indicate the presence of heterogeneity in the disorder. This could require alternative modelling techniques, or more thorough evaluation (e.g. removing sites that are very different).

Preprocessing pipelines

One aspect not covered in our analysis is the effect of variability in pre-processing pipelines (e.g. whether to perform global signal regression, Fisher transform, etc.) on model performance and biomarker robustness. These issues are complex and should be analysed in a separate study. One way to reduce the variability of pre-processing pipelines is to release pre-processed connectome datasets. We note from our analysis that the availability of pre-processed datasets drives research: about 1/3 of the studies are on ASD and almost all (33/36) use the ABIDE I dataset, especially the preprocessed connectomes

with 871 subjects (12/36). While the ABIDE II dataset is also available, few studies utilise it due to the lack of pre-processed data. Having pre-processed data made available helps to remove one source of variability and also allows results to be compared across studies. Thus, we encourage future data releases to also share pre-processed fMRI datasets, or even connectomes.

Graph construction

Across the papers that have been reviewed, we found that Pearson’s correlation was the most widely adopted approach for analysing FC in BG setups. However, several papers have highlighted the limitations of Pearson’s correlation and demonstrated how alternatives such as partial correlation [31], [114] and sparse representation could be less affected by noise and perform better [47]. Thus, future fMRI studies could consider looking beyond Pearson’s correlation and adopting alternative ways of constructing functional connectomes.

For PG, it is well known that features used for graph construction are often chosen arbitrarily. We note that studies using adaptive approaches of learning the PG [124], [139] tend to do better than other approaches using the same datasets. Thus, future studies could consider such approaches instead of handcrafting them.

Classification accuracies

Many papers on disorder prediction are centered on the novelty of their proposed architecture and demonstrate their utility by achieving a better prediction performance than baseline and state-of-the-art models. However, if the goal is to discover biomarkers that are generalisable across the population, then sizable datasets should be expected to be used to justify such claims. Furthermore, demonstrating out-of-distribution generalisation (e.g. another dataset not used during training) with limited performance degradation should

become a requirement. Few existing papers go to such extents of validating their model performance.

Additionally, most studies do not carefully account for confounds such as age, gender and scanner variability. It is unclear whether simply incorporating them in the population graph would act in a similar way as regressing them out, or whether techniques such as ComBat would be sufficient. Future studies should delve deeper into these factors and study their impact on biomarker analysis.

Effect of GNN construction on explanations

Greater care has to be taken when designing the GNN architecture for biomarker discovery applications. For instance, using BG gives the option of explaining node features or the (sub-)graph if the FC matrix is used. However, if thresholding or binarisation of the adjacency matrix was performed, subgraph-based explainers would not be recommended in these cases.

In the case of PG, only the node features would typically be used for generating explanations - it might be possible to subgraph-based explainers such as GNNE explainer to identify sub-populations from PG, but it would only provide a subjective and very indirect means of linking back to the imaging features. However, if feature selection was performed for the node features, then it might not be very ideal for biomarker discovery since the salient features are compared against an incomplete subset of features. Instead of using feature selection, we recommend that such studies adopt a BG + PG approach where the BG learns an encoding of the (multi-modal) imaging data. Although such a model seems more complex, it allows for a more direct form of model explainability.

In both BG and PG, using graphs based on k-NN makes it unsuited for explainers that are based on subgraphs. Such considerations have to be made before using the model for biomarker discovery applications.

Choice of Explainers

While a formal benchmarking study should be conducted to determine the best combination of predictors and explainers to use, pre-existing knowledge about the limitations of each explainer can be used to make more appropriate choices. For instance, numerous studies computed importance scores via gradients. However, gradients suffer from the problem of saturation and more advanced forms of IG discussed in the main manuscript (e.g. Guided IG) should be used instead of vanilla gradients/saliency.

Studies that aim to produce clinically useful results should use explainers that can produce individualised explanations, such as IG. Global explainers such as XGNN would not be useful in such a situation. However, much work is still needed to study the consistency of these explainers before they can be used for individual insights.

ADDITIONAL DISCUSSION - OUTLOOK

Multimodal biomarkers

The scope of this review paper is limited to resting-state fMRI and it is important to keep in mind that it only provides

information about brain function at rest. Whenever possible, other modalities (structural, molecular, task-fMRI, etc.) should be studied in tandem to produce a complete picture. In this regard, GNN provides a very flexible architecture to model multi-modal datasets: (i) modalities that are traditionally image-based (e.g. T1w, PET) can still be converted into graphs [206] such as morphometric similarity networks [55], [133]. (ii) GNN can be used to integrate connectome datasets with multi-omics data [157] in a parameter-efficient manner [148], (iii) PG make it possible to integrate non-imaging modalities [24]. Multi-modal biomarkers will be another emerging area that would give a more holistic view than what is currently possible with studies only using fMRI [156].

Intrinsically-interpretable models

Pooling and attention were found to be the most popular means of introducing model explainability into GNNs. However, these models are still not completely interpretable. Future studies could look into ways to develop models that are more interpretable. Several promising research directions include invertible GNN [122] and varying-coefficient model (in particular, contextual explanation network) [37].