

# Dynamic Graph Representation Learning for Spatio-Temporal Neuroimaging Analysis

Rui Liu<sup>ID</sup>, Member, IEEE, Yao Hu<sup>ID</sup>, Member, IEEE, Jibin Wu<sup>ID</sup>, Member, IEEE, Ka-Chun Wong<sup>ID</sup>, Zhi-An Huang<sup>ID</sup>, Member, IEEE, Yu-An Huang<sup>ID</sup>, Member, IEEE, and Kay Chen Tan<sup>ID</sup>, Fellow, IEEE

**Abstract**—Neuroimaging analysis aims to reveal the information-processing mechanisms of the human brain in a noninvasive manner. In the past, graph neural networks (GNNs) have shown promise in capturing the non-Euclidean structure of brain networks. However, existing neuroimaging studies focused primarily on spatial functional connectivity, despite temporal dynamics in complex brain networks. To address this gap, we propose a spatio-temporal interactive graph representation framework (STIGR) for dynamic neuroimaging analysis that encompasses different aspects from classification and regression tasks to interpretation tasks. STIGR leverages a dynamic adaptive-neighbor graph convolution network to capture the interrelationships between spatial and temporal dynamics. To address the limited global scope in graph convolutions, a self-attention module based on Transformers is introduced to extract long-term dependencies. Contrastive learning is used to adaptively contrast similarities between adjacent scanning windows, modeling cross-temporal correlations in dynamic graphs. Extensive experiments on six public neuroimaging datasets demonstrate the competitive performance of STIGR across different platforms, achieving state-of-the-art results in classification and regression tasks. The proposed framework enables the detection of remarkable temporal association patterns between regions of interest based on sequential neuroimaging signals, offering medical professionals a versatile and interpretable tool for exploring task-specific

Received 23 April 2024; revised 7 September 2024, 14 November 2024, and 6 January 2025; accepted 12 January 2025. Date of publication 4 February 2025; date of current version 7 March 2025. This work was supported in part by the National Nature Science Foundation of China under Grant 62202399, Grant U21A20512, and Grant 6220239; in part by the Research Grants Council of the Hong Kong, SAR, under Grant PolyU11211521, Grant PolyU15218622, Grant PolyU15215623, and Grant C5052-23G; in part by the Fundamental Research Funds for the Central Universities under Grant G2023KY05102; in part by the General Program of National Natural Science Foundation of China under Grant 62472353; and in part by the Guangdong Basic and Applied Basic Research Foundation under Grant 2024A1515011984 and Grant 2023A151514013. This article was recommended by Associate Editor J.-H. Xue. (Corresponding authors: Zhi-An Huang; Yu-An Huang.)

Rui Liu, Yao Hu, Jibin Wu, and Kay Chen Tan are with the Department of Data Science and Artificial Intelligence, The Hong Kong Polytechnic University, Hong Kong, SAR (e-mail: ruiliu@polyu.edu.hk; echo-yao.hu@polyu.edu.hk; jibin.wu@polyu.edu.hk; kaychen.tan@polyu.edu.hk).

Ka-Chun Wong is with the Department of Computer Science, City University of Hong Kong, Hong Kong, SAR (e-mail: kc.w@cityu.edu.hk).

Zhi-An Huang is with the Department of Computer Science, City University of Hong Kong (Dongguan), Dongguan 523000, China (e-mail: huang.za@cityu-dg.edu.cn).

Yu-An Huang is with the School of Computer Science, Northwestern Polytechnical University, Xi'an 710000, China (e-mail: yuanhuang@nwpu.edu.cn).

This article has supplementary material provided by the authors and color versions of one or more figures available at <https://doi.org/10.1109/TCYB.2025.3531657>.

Digital Object Identifier 10.1109/TCYB.2025.3531657

neurological patterns. Our codes and models are available at <https://github.com/77YQ77/STIGR/>.

**Index Terms**—Contrastive learning, electroencephalography (EEG), functional near-infrared spectroscopy (fNIRS), graph neural networks, interpretable visualization, magnetic resonance imaging (MRI), self-attention, spatio-temporal dynamics.

## I. INTRODUCTION

HUMAN brain is a highly intricate neurobiological system that plays a critical role in coordinating human behavior and cognition. Understanding the structure and function of the human brain network has become an enthralling endeavor with a variety of purposes, including mental disease diagnosis [1], brain-computer interface [2], and neuromorphic computing. Recent advancements have enabled researchers to measure brain activity using various neuroimaging techniques, such as magnetic resonance imaging (MRI), electroencephalography (EEG), functional near-infrared spectroscopy (fNIRS), etc. Functional magnetic resonance imaging (fMRI) detects variations in blood-oxygen-level-dependent signals using magnetic fields, EEG records spontaneous electrical activity of the brain, and fNIRS measures blood flow signals through detectors placed on the scalp [3]. Despite their ability to capture brain activity with high spatial and temporal resolution [4], there is currently no universal method to unify the analysis of complex spatio-temporal dynamics from various neuroimaging data. Especially since different types of neuroimaging data require different analysis tools, it is challenging to integrate these techniques into a cross-platform analysis system [5]. To this end, we propose a cross-platform analysis framework that can be applied to different downstream tasks, such as classification, regression, and interpretation, offering a comprehensive understanding of the human brain network.

Graph Neural Network (GNN), as shown in Fig. 1(B3), is a powerful model for transforming and analyzing graph-structured data while maintaining the graph property (including nodes, edges, and global context) and permutation invariances [7]. The intrinsic graph-structured nature of the brain has led to the widespread use of GNNs to learn the representations of the brain network from neuroimaging data. Recent studies have demonstrated the potential of GNNs in addressing heterogeneity across different neuroimaging data, making the GNN-based graph representation an ideal and general solution for analyzing such data. For example, when GNNs were applied to predict depression based on fMRI, EEG, and fNIRS data, they all focused on functional changes occurring in the prefrontal cortex [5], [8], [9]. Although these

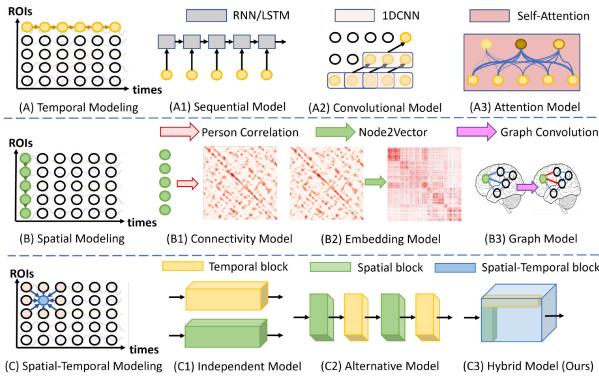


Fig. 1. Overview of different types of spatial and temporal modeling techniques (inspired by [6]). (a) Temporal modeling aims to capture dependencies between activities at different time steps. (A1)–(A3) illustrate commonly used temporal modeling techniques. (b) Spatial modeling focuses on exploring the correlations of activities at different brain regions. (B1)–(B3) illustrate commonly used spatial modeling techniques. (c) Spatio-temporal modeling captures and fuses both spatial and temporal information. Specifically, (C1) represents the independent mode and models spatial and temporal representations separately using different network blocks that work in parallel. (C2) represents the alternative mode, which models the spatial and temporal representations using separate network blocks in sequence. and (C3) represents the hybrid mode, which jointly models spatial and temporal representations using a consolidated network block (e.g., our STIGR).

studies have demonstrated the potential of GNNs in the analysis of neuroimaging data, most of them are grounded in a static manner like spatial modeling methods [Fig. 1(B1) and (B2)] [10]. Growing evidence suggests that functional brain activity can vary rapidly and involve different activated brain regions, providing important information about brain function [11]. Some recent attempts, e.g., Liu et al. [12], have combined dynamic FC analysis with static graph convolution for brain network analysis. Compared with static spatial modeling, the dynamic temporal modeling paradigm can capture the inherent temporal dynamics of brain activity [13], providing a richer and more nuanced representation of the underlying neural processes. As illustrated in Fig. 1(A1)–(A3), sequential models [14], convolutional models [15], and attention models [16] have been utilized to dynamically analyze neuroimaging and have demonstrated a remarkable capability in modeling temporal dependence. Despite notable progress in this field, the aforementioned methods primarily focus on the spatial or temporal dependencies of neuroimaging signals. This oversight neglects the potential benefits of jointly exploring their interplay, which is a critical avenue for elucidating the workings of the brain.

Spatial temporal graph neural network (STGNN) is an emerging branch of GNN to handle time-series data, which has shown great promise in capturing both spatial structure and temporal dynamics in brain activity. However, directly applying STGNN to brain network analysis has some inherent limitations that need to be addressed. First, existing STGNNs and spatio-temporal modeling methods [17], [18], [19], [20], [21] typically model spatial and temporal dependencies independently or alternatively, as illustrated in Fig. 1(C1) and (C2). In the independent mode, distinct temporal and spatial modules are utilized in parallel to capture spatial and temporal dependencies, followed by their combination to comprehend the overall spatio-temporal dependencies within

the data. Conversely, in the alternating mode, separate modules are employed to model temporal and spatial dependencies, alternately employing them to extract and incorporate spatio-temporal dependencies. However, both modes fail to comprehensively capture the intricate interplay between spatial and temporal dynamics in complex brain activity. For independent mode, the spatial and temporal module generates exchanges solely at the final stage of late fusion [22], overlooking crucial spatio-temporal interactions, such as how the connectome scales and brain networks evolve over time. For alternative mode, although there is an increased interaction between the spatial and temporal modules [23], this alternation neglects to model the relationship between nodes at previous and subsequent time points (i.e., the spatio-temporal interdependency in Fig. 3), which presents challenges in capturing the inter-relationship between spatial and temporal dynamics [24]. Thus, these models are limited in their ability to comprehensively investigate the interactions between spatial and temporal elements in the dynamic brain network. Second, STGNN typically performs temporal convolutions on the input dynamic graph. However, due to the limited kernel size of the convolution operation, these models are unable to capture long-range dependencies that are beyond the receptive fields of the convolution kernels [15]. Third, traditional STGNNs rely on a fixed adjacency matrix to represent the topological structure of the brain network [25]. However, brain activity is a dynamic process, and the connectivity between brain regions evolves over time. The static representation of brain topology using a fixed adjacency matrix makes it difficult to capture the temporal dependence of different brain regions when performing cognitive functions.

In this work, we propose a novel spatio-temporal interactive graph representation framework (STIGR) (as shown in Fig. 2), dubbed STIGR, for dynamic brain network analysis. STIGR comes with three main components for learning precise and discriminative spatio-temporal graph representations: 1) local dynamic graph convolution network (L-DGCN); 2) global spatio-temporal attentive network (G-STAN); and 3) contrastive-learning-based adjacent matrix learning (CL-AM). L-DGCN is designed to capture the complex inter-relationships between spatial and temporal dynamics by learning intrinsic local dynamic graph representations between adjacent time periods. G-STAN leverages the Transformer architecture to effectively capture long-range dependencies, resulting in enhanced global attentive spatio-temporal representations. By utilizing contrastive learning to adaptively learn the adjacency matrix of L-DGCN, CL-AM can accurately model the connection relationships of dynamic graphs. The major contributions of this work are summarized as follows.

- 1) *Pioneering Framework for Versatility Neuroimaging Analysis:* We present a promising neuroimaging analysis framework that spans multiple platforms and tasks, enabling efficient learning and representation of intricate spatio-temporal relationships within neuroimaging data. Our framework represents the first attempt to hybrid model the dynamic interactions of brain networks, as illustrated in Fig. 1(C3), addressing a critical gap in current research.

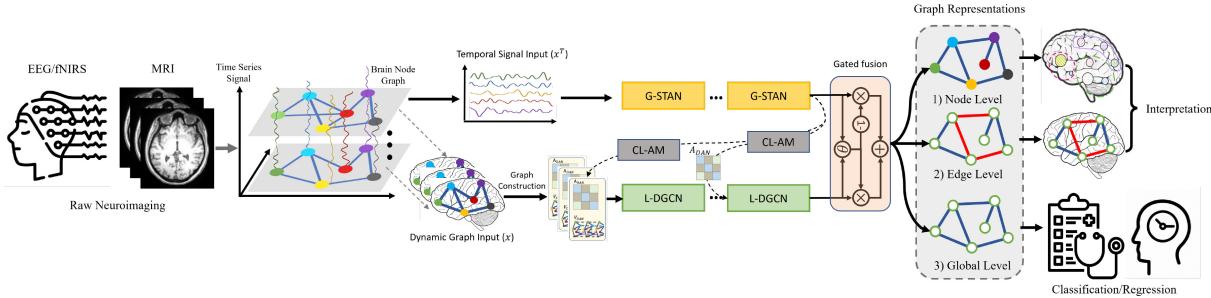


Fig. 2. STIGR framework illustration showcases its core concept as a dynamic spatio-temporal graph representation learning network based on graph convolution. By incorporating graph convolution and global attention mechanisms on the constructed dynamic spatio-temporal graphs, the network effectively learns discriminative graph representations of the neuroimaging data. Specifically, the raw neuroimaging data are first denoised to obtain temporal signal input  $x^T$  and dynamic graph input  $x$ . The temporal signal input  $x^T$  is then processed by a stack of global spatio-temporal attentive network (G-STAN), extracting global spatio-temporal representations. Simultaneously, the dynamic graph input  $x$  is divided into multiple windows, generating different brain node graphs. By utilizing contrastive learning guided by the G-STAN, the CL-AM module generates an adaptive adjacency matrix for the dynamic brain graph. Subsequently, a stack of local dynamic graph convolution networks (L-DGCNs) with the CL-AM module captures local spatio-temporal representations from the dynamic graph input  $x$ . The learned graph representations from both perspectives are fused using a gated fusion module, enabling the handling of various downstream tasks at the node, edge, and graph levels.

- 2) *Innovative Spatio-Temporal Components for Enhanced Learning:* We introduce two novel spatio-temporal modeling components, L-DGCN and G-STAN, which utilize graph convolution and attention mechanisms to learn discriminative spatio-temporal representations from local and global perspectives. In addition, our approach employs contrastive learning to adaptively refine the adjacency matrix of L-DGCN, guided by G-STAN, thereby enhancing the interaction between local and global learning processes.
- 3) *Exceptional Performance and Interpretability Across Neuroimaging Tasks:* Extensive experiments demonstrate the superior performance of our proposed STIGR in cross-platform/task neuroimaging analysis across a diverse range of modality, including MRI, EEG, and fNIRS. It outperforms State-of-the-Art (SOTA) neuroimaging analysis methods in tasks, such as mental disorder classification, motor imagery classification, and brain age prediction. Furthermore, STIGR facilitates interpretable neuroimaging analysis, allowing the exploration of discriminative patterns at various levels, including nodes, edges, and graphs.

## II. PROBLEM DEFINITION

The objective of our research is to address a key problem in neuroimaging analysis: the development of an effective graph representation framework for a set of tasks, such as classification, regression, and interpretation. In this context, we define the input data as a set,  $D = \{X, Y\}$ , where  $X$  represents a series of neuroimaging time-series sequences, and  $Y$  represents the corresponding outcomes, such as diagnostic labels or clinical scores. Each sequence in  $X$  is of a fixed dimension,  $N \times T$ , where  $N$  denotes the spatial dimension of the neuroimaging data, such as the number of ROIs for fMRI or electrodes/channels for EEG, and  $T$  represents the number of time points across all sequences. The goal of our study is to learn a function  $f : X \rightarrow Y$  that can accurately classify or regress these sequences, thus contributing to our understanding and analysis of neuroimaging data. Our objective is to apply

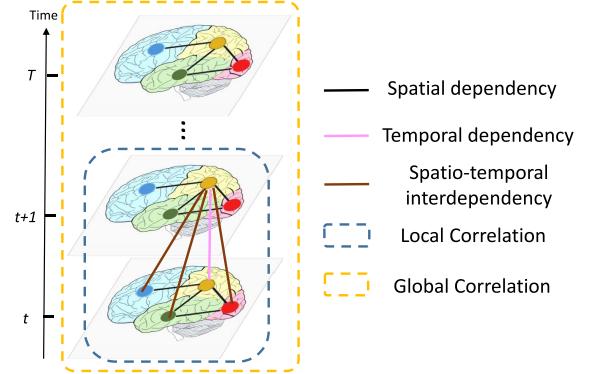


Fig. 3. Illustration of three types of dependencies and two types of correlations in spatio-temporal neuroimaging data: spatial dependency (black lines connecting different ROIs), temporal dependency (pink line representing the influence of a node on itself at the next time point), and spatio-temporal interdependency (brown lines indicating interdependencies between a node and its neighboring nodes at the next time point). Local correlation is shown by the blue dashed rectangle, representing the correlation between a dynamic graph and its temporal and spatial neighbors, encompassing spatial, temporal, and spatio-temporal dependency. Global correlation is depicted by the yellow dashed rectangle, illustrating the correlation between distant graphs.

graph-based analytical techniques to a time-series matrix  $X \in \mathbb{R}^{N \times T}$ , representing neuroimaging data. To facilitate this, we transform the time-series data into a graph representation. This involves constructing a dynamic graph network  $G(X) = \{\mathcal{G}(1), \dots, \mathcal{G}(T)\}$  where each graph  $\mathcal{G}(t)$  at time  $t$  consists of a vertex set  $V(t) = \{X_1(t), \dots, X_N(t)\}$  denoting  $N$  nodes and an edge set  $A(t)$  delineating the connectivity between nodes based on the neighborhood function  $\mathcal{N}(i)$ . This graph-based representation allows us to apply sophisticated algorithms designed for graph data to a set of tasks (such as classification, regression, and interpretation) for the original time-series datasets. Therefore, this problem can be formulated as an optimization task, which can be defined as follows:

$$f^* = \arg \min_f \mathbb{E}_{X,Y} [\mathcal{L}(f(G(X)), Y) + \mathcal{R}(f; A)] \quad (1)$$

where  $\mathbb{E}$  denotes the expected value of the loss function  $\mathcal{L}$  over the space of  $(X, Y)$ , and  $\mathcal{R}$  represents an additional

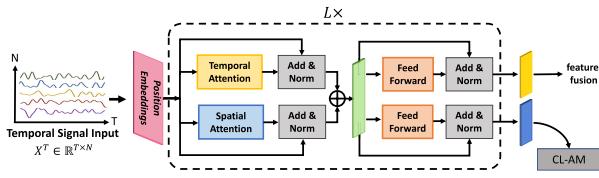


Fig. 4. Proposed G-STAN architecture comprises two parallel attention networks designed to extract global attention information from both spatial and temporal dependences in the data.

regularization term on  $A$ . It is worth noting that in this article, the superscript  $T$  in  $x^T$  denotes the transpose operation, whereas the subscript  $T$  in  $X_T$  denotes a temporal variable.

### III. MODEL DESIGN

Fig. 3 presents a visualization of the interdependencies and correlations associated with dynamic graph representation. These interdependencies encompass spatial, temporal, and spatio-temporal relationships, collectively forming our local correlations. They capture connections between a dynamic graph and its temporal and spatial neighbors, facilitating the analysis of local relationships. In addition, it is crucial to consider the global correlation, which accounts for the interdependence between distant graphs and mitigates the impact of short-term variations in neuroimaging analysis [26].

To effectively incorporate these multiple dependencies and correlations, we present the graph representation learning framework of STIGR, as depicted in Fig. 2. The STIGR comprises three main components: 1) G-STAN; 2) L-DGCN; and 3) CL-AM. The G-STAN and L-DGCN are designed to jointly learn the spatio-temporal graph representation of neuroimaging data at global and local levels, respectively. To better model the node connections between different windows, CL-AM is designed to adaptively extract the cross-window temporal connectivity in L-DGCN based on contrastive learning. The local and global spatio-temporal representations are then adaptively integrated using a gated fusion module to obtain the discriminative graph representation.

#### A. Global Spatio-Temporal Attentive Network

Capturing global dependency is essential to provide a broader context and understanding of the overall brain dynamics [16]. The Transformer [27] is highly effective for modeling global dependency, which explicitly computes the dependency of tokens at different time steps. However, the vanilla Transformer was designed for sequence data and did not account for the spatial dependencies (e.g., ROIs connectivity) in the brain network. To address this limitation, we propose the G-STAN module based on the Transformer encoder to learn global spatio-temporal representations by decoupling the attention mechanism into temporal and spatial dimensions.

As illustrated in Fig. 4, the G-STAN module takes the temporal signal sequence  $X^T \in \mathbb{R}^{T \times N}$  as its input and produces both temporal representations  $E_T^{(L)} \in \mathbb{R}^{T \times N}$  and spatial representations  $E_S^{(L)} \in \mathbb{R}^{N \times T}$  as its outputs, which serve distinct purposes for feature fusion and CL-based adjacent matrix learning, respectively. Specifically, the input sequence

$X^T$  is first processed by a position embedding layer to obtain the hidden feature representation  $E \in \mathbb{R}^{T \times N}$  that encodes position information. The position embedding is implemented based on [28] and can be formulated as follows:

$$E = \text{PositionEmbedding}(X^T) = X^T + PE \quad (2)$$

where the position embedding  $PE$  is represented as a learnable table of vectors, which assigns a unique vector representation to each ROI signal based on its position index of the brain. Then, G-STAN utilizes a dual self-attention structure with  $L$  layers to separately extract the attention focus from the temporal and spatial dimensions of  $E$ . In particular, temporal attention selectively fuses information from different time points of the same node (e.g., ROI of fMRI or channel of EEG), while spatial attention fuses information from different nodes that are collected at the same time point. In the temporal attention block, we learn three matrix representations  $Q$ ,  $K$ , and  $V$  (queries, keys, and values) given  $E$  as

$$Q = f_Q(E), K = f_K(E), V = f_V(E) \quad (3)$$

where  $f_Q$ ,  $f_K$ , and  $f_V$  are the corresponding projection functions of *queries*, *keys*, and *values*. The dot-product similarity is then used to compare *queries* with *key-value* pairs to obtain the attention distribution on the *values*. If the *queries* and the *keys* are similar (i.e., high attention weight), the corresponding *values* are considered related. The resulting weighted values matrix becomes the output of the attention block. Following [27], multihead attention (MHA) is adopted to extract temporal-level attention, projecting the *queries*, *keys*, and *values* with  $h$  different linear projection heads. The output representation of the temporal attention block  $\tilde{E}_T$  can be formulated as follows:

$$\tilde{E}_T = \text{MHA}(E) = \langle \text{head}_1, \dots, \text{head}_h \rangle W^o \quad (4)$$

$$\text{head}_j = \text{Att}_j(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{N}}\right)V \quad (5)$$

where  $W^o \in \mathbb{R}^{N \times N}$  represents the linear transformation matrices of the head concatenation. Distinct from conventional temporal convolution layers, G-STAN leverages a self-attention based on the Transformer encoder to capture long-term temporal dependencies in time-series data. Without convolution layers, this self-attention module can learn the significance of each time point in the context of the entire sequence by attending to different positions of the input sequence.

Likewise, the spatial attention block is constructed in a similar way to the temporal attention block, but with the spatial-level  $Q$ ,  $K$ , and  $V$  matrices learned over the transpose of the input embedding feature (i.e.,  $E^T \in \mathbb{R}^{N \times T}$ ) to output the spatial attention for all nodes at the same time point (denoted as  $\tilde{E}_S$ ). Since the temporal and spatial attention blocks operate independently to calculate  $\tilde{E}_T$  and  $\tilde{E}_S$ , then element-wise addition  $\oplus$  is used to capture the interdependencies between spatial and temporal representations, i.e.,  $\tilde{E} = \tilde{E}_T \oplus \tilde{E}_S$ . The output is then passed through two parallel feed-forward networks, allowing our network to extract spatio-temporal interdependencies while exploring modality-specific intradependencies. Following the design principle of

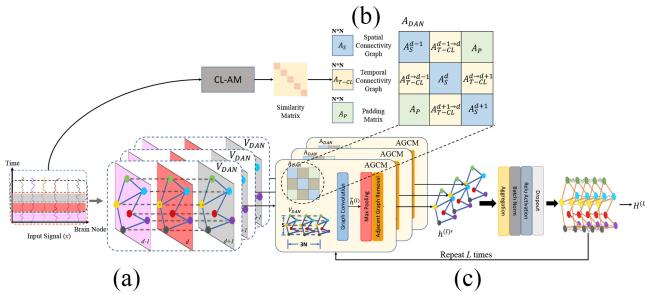


Fig. 5. (a) Input illustration for the dynamic adaptive-neighbor graph generated along the time axis. (b) Adjacency matrix of the dynamic adaptive-neighbor graph in (a). (c) Workflow of L-DGCN and architecture of AGCM.

Transformer, we apply residual connection with layer normalization, abbreviated as *Add&Norm*, after the dual attention and feedforward layers. We stack  $L$  spatio-temporal attention layers to successively update the fused embeddings. Finally, G-STAN outputs temporal and spatial representations  $E_T^{(L)}$  and  $E_S^{(L)}$  for further global-local feature fusion and CL-based adjacent matrix learning, respectively.

### B. Local Dynamic Graph Convolution Network

To capture the full range of spatial, temporal, and spatio-temporal interdependencies inherent in dynamic brain activity, we introduce a novel dynamic adaptive-neighbor graph construction ( $\mathcal{G}_{DAN} \in \{V_{DAN}, A_{DAN}\}$ ) to be used in our L-DGCN, as shown in Fig. 5(a) and (b). Specifically, we divide the  $T$  time points into  $D$  scanning windows with a length of  $S$  each ( $D = \lfloor T/S \rfloor$ ). Our graph convolution approach distinguishes itself from standard spatial convolutions by forming connections between nodes across adjacent temporal windows, generating localized spatio-temporal subgraphs adept at capturing short-term dynamics in brain networks. Such a dynamic adaptive-neighbor graph construction extends beyond individual spatial graphs, linking them across time to map intricate spatio-temporal interdependencies.

The scale of  $\mathcal{G}_{DAN}$  is determined by the number of adjacent scanning windows. For example, we set it to 1 as shown in Fig. 5, resulting in  $V_{DAN} = [V[d-1], V[d], V[d+1]] \in \mathbb{R}^{3N \times S}$  [as shown in Fig. 5(a)] and  $A_{DAN}$  with the size of  $3N \times 3N$  [as shown in Fig. 5(b)]. Inspired by the recent advance in [24], we leverage  $A_{DAN}$  to record three types of adjacency matrix  $\in N \times N$ : 1) spatial connectivity graph matrix; 2) temporal connectivity graph matrix; and 3) padding matrix, as follows.

- 1) The spatial connectivity graph matrix ( $A_S$ ) represents the node connections within dynamic subgraphs in each adjacent scanning window and is positioned on the diagonal of  $A_{DAN}$ . The construction of  $A_S$  is based on the calculation of the Pearson correlation coefficient (PCC). Specifically, the  $A_S$  for the  $d$ th window is computed as follows:

$$A_S^d(i, j) = \frac{\sum_{k=1}^S (v_d(i, k) - \bar{v}_d(i)) (v_d(j, k) - \bar{v}_d(j))}{\sqrt{\sum_{k=1}^S (v_d(i, k) - \bar{v}_d(i))^2 (v_d(j, k) - \bar{v}_d(j))^2}} \quad (6)$$

where  $i$  and  $j$  represent the  $i$ th and  $j$ th nodes, respectively. Notably, we selected the top 10% of the spatial correlation in  $A_S$  to reduce redundant information. In this

approach, we define significant FCs as those with a PCC exceeding the threshold  $\tau = 0.1$ , which are assigned a value of 1, while all other correlations are assigned a value of 0.

- 2) The temporal connectivity graph matrix ( $A_{T-CL}$ ), which works in conjunction with the spatial graph matrix  $A_S$ , is designed to capture the temporal relationships between nodes across different subgraphs. The CL-AM module plays a pivotal role in adaptively learning the construction of  $A_{T-CL}$ , a process that we detail in Section III-B3) for further clarity.
- 3) The padding matrix ( $A_P$ ) is included to fill the empty space left after dividing the temporal signal into adjacent scanning windows. It is constructed by setting all its entries to 0 and is used for padding operations.

In this way,  $\mathcal{G}_{DAN}$  can extract the local spatio-temporal correlations of dynamic graph matrices by the following adjacent graph convolution module (AGCM).

We further extend our analysis to the entire time course and employ AGCM to capture hybrid spatio-temporal relations among the entire dynamic brain network, as shown in Fig. 5(c). The L-DCGN module accepts a graph input signal  $X \in \mathbb{R}^{N \times T}$  as its input and outputs a learned graph representation  $H^{(L)} \in \mathbb{R}^{N \times T}$ . Utilizing the proposed dynamic graph construction method, the input graph signal  $X$  is segmented into  $D$  subgraphs, each denoted as  $\mathcal{G}_{DAN} \in \{V_{DAN}, A_{DAN}\}$ . Subsequently, for each subgraph  $\mathcal{G}_{DAN}$ , a graph convolution operation aggregates the central node's features with those of its neighboring nodes across adjacent windows to form a comprehensive hybrid graph representation. For the  $l$ th input graph representation of  $d$  subgraph, denoted as  $h_d^{(l)} \in \mathbb{R}^{3N \times S}$ , the graph convolution operation is applied or performed as

$$\tilde{h}_d^{(l)} = \sigma(A_{DAN}(h_d^{(l)} W + b)) \quad (7)$$

where  $W \in \mathbb{R}^{S \times S}$  and  $b \in \mathbb{R}^S$  are learnable parameters, and  $\sigma$  represents the *ReLU* activation function. After graph convolution, max pooling is performed over the time dimension to aggregate the representation from adjacent windows, resulting in a graph representation  $h_d^{(l)'} \in \mathbb{R}^{N \times S}$ . Each scanning window of the spatio-temporal graph input  $H^{(l)}$  is processed in parallel by AGCM, generating the local graph representations  $h_d^{(l)'}$  (where  $d \in [1, D]$ ) as shown in Fig. 5(c). It should be noted that zero-padding is applied when processing the first and last scanning windows. Then, an aggregation layer is employed to concatenate  $h_d^{(l)'}$  as follows:

$$H^{(l)'} = \langle h_1^{(l)'}, \dots, h_D^{(l)'} \rangle \in \mathbb{R}^{D \times N \times S} \quad (8)$$

where  $\langle . \rangle$  denotes the concatenation operator. Furthermore, a batch normalization layer, *ReLU* function, and a dropout layer are cascaded to generate the final local spatio-temporal graph representation  $H^{(l+1)}$  of layer  $l$ . The final output after  $L$  layers of L-DGCN is represented as  $H^{(L)}$ .

### C. Contrastive-Learning-Based Adjacent Matrix

To effectively capture the dynamic hybrid graph representation of each central node and its neighboring nodes across

adjacent windows, it is essential to model the node connections within and between different windows simultaneously. Generally, the spatial connectivity graph matrix  $A_S$  can be calculated using static coherency-based measurements (e.g., PCC) to model the degree of similarity between connectome nodes within the same window. However, modeling node connections between different windows  $A_{T-CL}$  is rather challenging due to their complex cross-temporal connectivity characteristics. Therefore, we propose a CL-AM module based on contrastive learning, which contrasts the similarity of learned representations between all scanning windows to model cross-temporal connectivity without label information.

The intuition behind contrastive learning is to pull similar samples closer and push away dissimilar samples, based on their representation. Here, we classify adjacent scanning windows at the same node as “positive” data points, indicating they belong to the same category. In contrast, scanning windows at different nodes are labeled as “negative” data points, indicating they belong to different categories. The representation of scanning windows can be extracted by G-STAN. As a result, G-STAN acts as a shared encoder for the CL-AM module to encode the underlying shared information between different scanning windows.

In addition to the perspective of similarity of features, contrastive learning in our work can also be interpreted from the view of contrastive predictive coding which involves learning representations that predict future data points from past observations [29]. In our CL-AM module, we leverage this principle by using embeddings from one session as context to predict the embeddings of the next session, assuming temporal continuity in data. Our goal is to make the embeddings of the corresponding nodes in adjacent sessions more similar to each other (values close to 1 on the diagonal of the similarity correlation matrix, SC) than to the embeddings of other nodes at different time points (values close to 0 for off-diagonal elements). This encourages the model to capture both intrasession node distinctions and the temporal evolution of nodal features. The contrastive loss function plays a key role by minimizing the distance between positive pairs (same node across sessions) and maximizing the distance between negative pairs (different or noncorresponding nodes across sessions). Consequently, the SC matrix, derived from these trained embeddings, reflects the degree of similarity between nodes across sessions. Higher similarity indicates a strong spatiotemporal connection, implying predictability of node states over time. Consequently, the final SC serves as a proxy for  $A_{T-CL}$ , enabling the construction of a dynamic graph that reflects the underlying spatiotemporal connectivity patterns.

As depicted in Fig. 6, the CL-AM module processes the spatial-level representation  $E_S^{(l)} \in \mathbb{R}^{N \times T}$  as input and produces a learned similarity matrix SC as output based on contrastive learning. Specifically, the spatial-level representation  $E_S^{(l)}$  is divided into  $D$  sequences of length  $S$  to learn the pair-wise temporal connectivity. For example, to learn  $A_{T-CL}^{d-1 \rightarrow d}$ , we first encode the feature representations of the previous window  $d - 1$  ( $E_S^{(d-1)} \in \mathbb{R}^{N \times S}$ ) and the current window  $d$  ( $E_S^{(d)} \in \mathbb{R}^{N \times S}$ ) using the spatial-level attention network of G-STAN. We then feed the representations of the two scanning windows into two

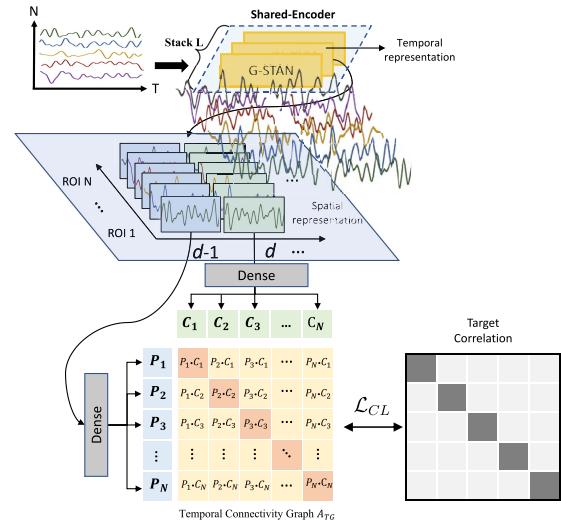


Fig. 6. Structure of the CL-AM, where the spatial network in G-STAN serves as a shared encoder.

different linear projection heads to map the representations to a latent space for similarity calculation. Therefore, we obtain two embedding vectors of  $N$  dimensions that represent the regional features of the window  $d - 1$  and the window  $d$ , denoted as  $P$  and  $C$ , respectively.

The CL-AM module is trained to predict which of the  $N \times N$  possible pairings actually occurred by taking the embeddings of  $N$  nodes in two windows as samples for contrastive learning. The CL-AM maximizes the similarity of the previous and current embeddings of the  $N$  positive pairings while minimizing the similarity of the embeddings of the rest  $\{N^2 - N\}$  negative pairings. The total number of  $A_{T-CL}$  that needs to be calculated in the learning process is denoted as  $M$ . Given the  $i$ th and  $j$ th nodes, the contrastive loss function of  $m$ th  $A_{T-CL}$  is calculated as follows:

$$\mathcal{L}_{CL}^m = \sum_i (1 - SC_{ii})^2 + \lambda \sum_i \sum_{i \neq j} SC_{ij}^2 \quad (9)$$

where  $\lambda$  is a positive constant that trades off the importance of the positive pairs (first term) and negative pairs (second term). The similarity coefficient  $SC_{ij}$  is computed between  $P_i$  and  $C_j$  as follows:

$$SC_{ij} = \frac{P_i C_j}{\|P_i\|_2 \cdot \|C_j\|_2} \quad (10)$$

where  $\|\cdot\|_2$  is  $l^2$ -norm. By optimizing  $\mathcal{L}_{CL}^m$ , the temporal connectivity graph matrix  $A_{T-CL}^{d-1 \rightarrow d}$  can be well represented by the learned SC for dynamic graph convolution. The total contrastive loss is defined as  $\mathcal{L}_{CL} = \sum_{m=1}^M \mathcal{L}_{CL}^m$  for the processing of L-DGCN, where  $m$  indicates the total number of  $A_{T-CL}$  needed to calculate in the learning process. The overall contrastive loss  $\mathcal{L}_{CL}$  can be regarded as a regularization term in (1) that serves to optimize  $A_{T-CL}$ .

#### D. Overall Architecture

Algorithm 1 outlines the program flow of STIGR. First, we reshape the input dynamic graph signal  $X \in \mathbb{R}^{N \times T}$

**Algorithm 1** Pseudocode of the STIGR

**Require:** Preprocessed data  $X$ , and label  $Y$

**Ensure:** Predicted probabilities of testing set  $Y_{pred}^{te}$

- 1:  $[X^{tr}, X^{te}, Y^{tr}, Y^{te}] \leftarrow \text{Split}(X, Y)$
- 2: Initialize STIGR.
- 3: **for**  $e = 1, \dots, epochs$  **do**
- // Global Spatio-Temporal Attentive Network
- 4:    $E = \text{PositionEmbedding}(X^{tr\top})$
- 5:    $E_T^{(0)} = E; E_S^{(0)} = E^\top$
- 6:   **for**  $l = 0, \dots, L-1$  **do** // l: # of attention layers
- 7:      $\tilde{E}^{(l)} = \text{MHA}(E_S^{(l)}) \oplus \text{MHA}(E_T^{(l)})$
- 8:      $E_T^{(l+1)} = \text{AddNorm}(\text{FFN}(\tilde{E}^{(l)}), \tilde{E}^{(l)})$
- 9:      $E_S^{(l+1)} = \text{AddNorm}(\text{FFN}(\tilde{E}^{(l)\top}), \tilde{E}^{(l)\top})$
- 10:  **end for**
- // Local Dynamic Graph Convolution Network
- 11:   $H^{(0)} = X^{tr}$
- 12:  **for**  $l = 0, \dots, L-1$  **do** // l: # of L-DGCN layers
- 13:    $A_{T-CL}^{(l)} = \text{CL-AM}(E_S^{(l)})$
- 14:    $A_{DAN}^{(l)} = \langle A_S, A_{T-CL}^{(l)}, A_P \rangle$
- 15:    $H^{(l+1)} = \text{L-DGCN}(H^{(l)}, A_{DAN}^{(l)})$
- 16:  **end for**
- 17:   $R = (1 - \theta)E_T^{(L)} \oplus \theta H^{(L)}$  // Fused representation
- 18:  // Tower networks for specific tasks
- 19:   $output \leftarrow \text{Tower}(R)$
- 20:   $Loss \leftarrow \alpha \mathcal{L}(output, Y^{tr}) + (1 - \alpha) \mathcal{L}_{CL}$
- 21:  STIGR.update( $Loss$ )
- 22: **end for**
- 23:  $Y_{pred}^{te} \leftarrow \text{STIGR.predict}(X^{te})$

Note: The  $\langle \cdot \rangle$  represents the construction of the adjacency matrix for the DAN graph.

into  $X^\top \in \mathbb{R}^{T \times N}$  and pass it through the  $L$  layer G-STAN module. This module learns the global spatial and temporal representations  $E_S^{(L)}$  and  $E_T^{(L)}$ , respectively. Based on the learned  $E_S^{(L)}$ , we then employ the CL-AM module to adaptively learn the temporal connectivity matrix  $A_{T-CL}$  to construct the adjacent matrix  $A_{DAN}$ . The dynamic graph signal  $X$  and learned adjacent matrix  $A_{DAN}$  are subsequently fed into the  $L$  layer L-DGCN module to derive the local spatio-temporal representation  $H^{(L)}$ . Upon obtaining the final outputs  $E_T^{(L)}$  and  $H^{(L)}$ , the gated fusion module is used to combine the local and global spatio-temporal representations as  $R = (1 - \theta)E_T^{(L)} \oplus \theta H^{(L)}$ , where  $\theta$  is a learnable parameter. The fused feature embeddings  $R$  are then directed to different output networks for various downstream tasks. For classification, the output tower network is composed of a global average pooling layer that aggregates node-level features to represent the entire graph at a global level and two dense layers that perform binary classification. The contrastive loss acts as a regularization effect on the representation learning process, resulting in the overall optimization function in (1)

$$f^* = \arg \min_f \mathbb{E}_{X,Y} [\alpha \mathcal{L}(f(G(X)), Y) + (1 - \alpha) \mathcal{L}_{CL}] \quad (11)$$

where  $\alpha$  is a tunable parameter that balances the weights of the losses.

**IV. EXPERIMENTS AND RESULTS****A. Data Acquisition**

**ABIDE I & II**<sup>1</sup>: ABIDE I comprises 1035 valid resting-state fMRI samples aggregated from 17 different brain imaging sites, including 505 autism spectrum disorder (ASD) subjects and 530 typical controls (TCs). ABIDE II dataset is also a multisite dataset containing 1113 resting-state fMRI from 19 different sites, comprising 521 ASD participants and 592 TCs.

**ADHD-200**<sup>2</sup>: ADHD-200 includes 939 preprocessed resting-state fMRI samples from 8 imaging sites, involving 358 children and adolescents with ADHD and 581 TCs.

**COBRE**<sup>3</sup>: The COBRE dataset consists of 146 resting-state fMRI samples involving 72 schizophrenia (SCZ) and 74 TCs.

**BCICIV\_2a**<sup>4</sup>: BCICIV\_2a is a motor imagery dataset with training and testing sets based on EEG recordings collected from nine subjects. The dataset includes four different motor imagery tasks, using a 22-channel EEG.

**BCIC2015**<sup>5</sup>: BCIC2015 is an EEG dataset obtained from 26 healthy participants who completed an error detection task based on the P300-Speller paradigm. The objective of this experiment is to utilize the analysis of 55 channels of EEG signals following the receipt of feedback to identify instances.

**fNIRS-BCI**<sup>6</sup>: fNIRS-BCI consists of fNIRS recordings obtained from a total of eight individuals, including three males and five females. During the experiments, a 52-channel fNIRS system was used to record changes in oxygenated hemoglobin and deoxygenated hemoglobin in the prefrontal cortex.

The performance of the experiments was evaluated in terms of accuracy (ACC), sensitivity (SEN), and specificity (SPE). A more detailed description of the model and parameters settings can be found in the supplementary Section I.

**B. Model Comparison**

In this section, we evaluate the performance of STIGR on fMRI, EEG, and fNIRS datasets for preliminary classification tasks, and compare it with SOTA methods, as shown in Fig. 7. Since most SOTA methods do not provide open-source code, we ensured fairness by directly comparing the performance of our models with the results reported in the respective papers. For fMRI datasets, we achieve competitive accuracy of  $72.73(\pm 3.37)\%$  (SEN: 76.08%, SPE: 68.84%, and AUC: 79.03%) on ABIDE I,  $72.53(\pm 3.59)\%$  (SEN: 73.00%, SPE: 71.56%, and AUC: 77.59%) on ABIDE II,  $76.15(\pm 3.74)\%$  (SEN: 79.47%, SPE: 72.35%, and AUC: 81.75%) on ADHD-200, and  $88.42(\pm 4.9)\%$  (SEN: 89.92%, SPE: 87.29%, and AUC: 92.44%) on COBRE, respectively. Notably, these datasets are cross-site and collected from different medical institutions with varying scanning equipment and protocols, resulting in significant data heterogeneity that can significantly affect the accuracy of the analysis [36].

<sup>1</sup>[http://fcon\\_1000.projects.nitrc.org/indi/abide](http://fcon_1000.projects.nitrc.org/indi/abide)

<sup>2</sup>[http://fcon\\_1000.projects.nitrc.org/indi/adhd200](http://fcon_1000.projects.nitrc.org/indi/adhd200)

<sup>3</sup>[http://fcon\\_1000.projects.nitrc.org/indi/retro/cobre.html](http://fcon_1000.projects.nitrc.org/indi/retro/cobre.html)

<sup>4</sup><https://www.bbci.de/competition/iv/>

<sup>5</sup><https://www.kaggle.com/c/inria-bci-challenge>

<sup>6</sup><http://bnci-horizon-2020.eu/database/data-sets>





**TABLE V**  
EFFECTIVENESS ANALYSIS OF DIFFERENT ATLAS, LEARNABLE COEFFICIENTS, AND NEIGHBOR STEP

| Section A: Learned and tunable coefficient values for models trained on different datasets    |                       |                   |             |                   |             |
|---|-----------------------|-------------------|-------------|-------------------|-------------|
| Coefficient   | ADHD-200              | ABIDE I           | COBRE       | BCICIV_2a         | fNIRS-BCI   |
| $\theta$ (in gated fusion)  | 0.43                  | 0.57              | 0.53        | 0.62              | 0.51        |
| $\alpha$ (in final loss function)   | 0.72                  | 0.68              | 0.57        | 0.74              | 0.62        |
| Section B: Effectiveness analysis of threshold $\tau$ in adjacent matrix construction.        |                       |                   |             |                   |             |
| Model   | Threshold ( $\tau$ )  | ABIDE I           |             | ADHD-200          |             |
| STIGR   | 0                     | 72.1( $\pm 4.0$ ) | 78.2        | 75.3( $\pm 3.7$ ) | 81.0        |
|   | 0.1                   | 72.7( $\pm 3.4$ ) | <b>79.0</b> | 76.2( $\pm 3.6$ ) | <b>81.8</b> |
|   | 0.2                   | 71.8( $\pm 3.2$ ) | 78.8        | 75.5( $\pm 3.9$ ) | 80.9        |
|   | 0.3                   | 71.9( $\pm 3.9$ ) | 77.5        | 74.8( $\pm 3.6$ ) | 81.2        |
|   | 0.5                   | 71.2( $\pm 3.4$ ) | 76.5        | 74.5( $\pm 3.8$ ) | 81.1        |
| Section C: Effectiveness analysis of different neighbor step on ABIDE I and ADHD-200 dataset. |                       |                   |             |                   |             |
| Model   | Neighbor step ( $k$ ) | ABIDE I           |             | ADHD-200          |             |
| STIGR   | 1                     | 72.7( $\pm 3.4$ ) | <b>79.0</b> | 76.2( $\pm 3.6$ ) | <b>81.8</b> |
|   | 2                     | 71.9( $\pm 3.5$ ) | 78.0        | 75.1( $\pm 4.0$ ) | 80.7        |
|   | 3                     | 71.2( $\pm 3.2$ ) | 77.5        | 74.5( $\pm 3.6$ ) | 80.5        |

analysis, while G-STAN extracts long-term spatio-temporal dependencies (global attention). Since the CL-AM module is proposed based on L-DGCN and G-STAN, we cannot evaluate its effectiveness by eliminating the other two components. However, we can assess the impact of the generated  $A_{T-CL}$  by directly comparing the performance of Model (d) and Model (e). The introduction of the CL-AM module results in significant accuracy improvements of 1.2%, 1.3%, and 2.1% for the ABIDE I, ADHD-200, and BCICIV\_2a datasets, respectively. This finding underscores the importance of adjacent matrix construction in graph convolution and highlights the potential of contrastive learning in modeling representation similarity. Overall, Model (e), which integrates all three components (L-DGCN, G-STAN, and CL-AM), consistently outperforms the baseline model, achieving accuracy improvements of 4.5%, 4.8%, and 8.8% on the ABIDE I, ADHD-200, and BCICIV\_2a datasets, respectively. Consequently, we conclude that all components clearly contribute to the performance enhancement of STIGR.

#### D. Parameter Sensitivity Analysis

To evaluate the effectiveness of the proposed components and novel loss function, we have tabulated the values of the learned coefficients  $\theta$  (in gated fusion) and  $\alpha$  (in the final loss function) in Section A of Table V. As anticipated, it is evident that neither  $\theta$  nor  $\alpha$  converge to the extremes of 0.0 or 1.0 across different datasets. Instead, it tends to remain around the mid-value of 0.5, which suggests that both the L-DGCN and G-STAN components jointly play an equitable role in augmenting the model performance across various datasets. In addition to the analysis of  $\theta$  across datasets, we conduct specific experiments on the ADHD-200 dataset to further examine the behavior of  $\theta$  within a unified dataset. The results are clearly presented in the accompanying Fig. 8. From the findings, it can be observed that the mean value of  $\theta$  initially starts around 0.5 and gradually decreases to 0.33, eventually converging at 0.42. This trajectory of  $\theta$  during training offers valuable insight into the parameter's adaptation within a unified dataset. Besides, it is also worth mentioning that the coefficient of the loss function,  $\alpha$ , consistently remains above 0.5, suggesting that the binary cross-entropy classification loss plays a more crucial role than the contrastive loss. We perform visualizations by varying the parameter  $\alpha$  from 0 to 1 [Fig. 9(a)], showcasing the performance differences observed

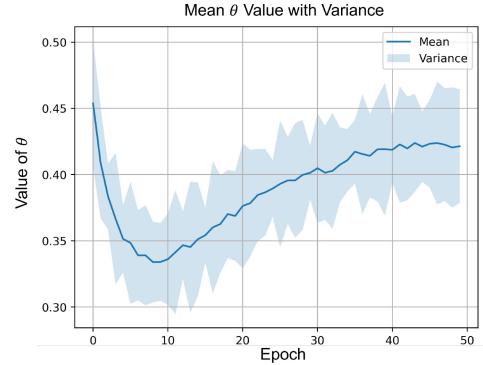


Fig. 8. Variation of learnable coefficient  $\theta$  during the training process on ADHD-200.

with the different  $L_{CL}$  loss weights. The accuracy remains relatively stable as the weight of  $L_{CL}$  is increased, with a sharp trend toward decreasing performance as alpha approaches from 0.2 to 0 (i.e., full  $L_{CL}$  contribution). The model performance across different datasets shows a peak in accuracy when the weight of  $L_{CL}$  is around 0.3. This suggests that a certain amount of  $L_{CL}$  is beneficial for enhancing accuracy, but too much weight on  $L_{CL}$  could lead to performance degradation. These insights are critical for fine-tuning the balance between the contrastive loss and other components of the loss function.

We also perform a parameter sensitivity analysis on two specific aspects: 1) the number of layers  $L$  and 2) the length of scanning windows  $S$ , specifically on the ABIDE I and ADHD-200. To evaluate the impact of the number of layers  $L$ , we explore settings to vary  $L$  from 1 to 4. Regarding the length of scanning windows  $S$ , we test the model under three different configurations:  $S = 15$ ,  $S = 30$  TR, and a hybrid setting,  $S = 15/30$  TR. In the case of  $S = 15/30$  TR, where we use 15 TR for the first layer and 30 TR for the second. The experimental results are presented in Fig. 9. With respect to the number of layers, we observe an enhancement in performance with an increase in  $L$  up to 2 layers, after which the performance noted a marginal decline. Regarding the length of scanning windows, the variations in performance across the different session lengths were minimal, suggesting that the model performance is not heavily dependent on this hyperparameter. Interestingly, the hybrid setting of  $S = 15/30$  TR exhibited a performance improvement over the singular settings. This can be attributed to its alignment with the hierarchical nature of GNNs' message-passing algorithms [43], which progressively aggregate information from closer to farther neighbors.

In addition, to examine the impact of the adjacency matrix construction in graph modeling, we conduct a series of experiments within our L-DGCN module. By conducting experiments with  $\tau$  values ranging from 0 to 0.5, we have gathered the experimental results that clarify the relationship between threshold settings and model performance, as shown in Section B of Table V. It is worth noting that when  $\tau = 0$ , it indicates that all FC values were utilized for matrix construction. Upon analyzing the result, we identify that increasing  $\tau$  initially improved the performance of our model and eventually reached a plateau with a slight decline.

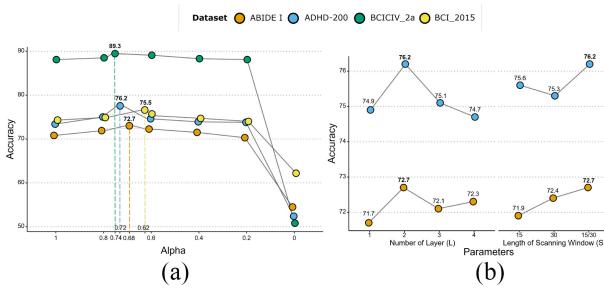


Fig. 9. Performance comparison with different hyperparameters. (a) Performance comparison with different  $\alpha$ . (b) Performance comparison with different  $L$  and  $S$ .

However, peak performance for both the ABIDE I and ADHD-200 datasets was observed at  $\tau = 0.1$ . These findings validate our proposition that a moderate  $\tau$  effectively balances the elimination of redundant information with the retention of pertinent data, thereby enhancing the efficiency of our model. Generally, variations in performance across different  $\tau$  levels are minimal, demonstrating the model's capability to effectively handle both scenarios of information redundancy at  $\tau = 0$  and information scarcity at  $\tau = 0.5$ .

In this section, we delve into how varying neighbor steps affect our framework's performance. By adjusting the neighbor steps, denoted as  $k$ , when constructing the dynamic adaptive-neighbor graph ( $V_{\text{DAN}}$ ), we explore different construction scenarios.  $V_{\text{DAN}}$  includes nodes within a range of  $k$  steps from the central node  $d$ , represented as  $V_{\text{DAN}} = [V[d - k], V[d], V[d + k]]$ . We assess  $k$  from 1 to 3 in simulation experiments, tabulated in Section C of Table V. Results show that  $k = 1$  yields the best performance on ABIDE I and ADHD-200 datasets. Increasing  $k$  leads to a slight decrease in Acc and Auc on both datasets, indicating that prioritizing immediate neighbors effectively extracts relevant information. This is because finer time granularity at  $k = 1$  allows for a precise representation of temporal dynamics and stronger correlations. However, as  $k$  increases, granularity decreases, limiting the capacity to capture fine-grained temporal relationships. Therefore, a smaller  $k$  is preferable for achieving higher performance on neuroimaging datasets. We also assess the efficacy of employing various atlases. Detailed results are provided in supplementary Section IV.

#### E. Brain Age Prediction

Accurately predicting brain age is crucial in neuroimaging analysis, as it adds in the early detection of vital indicators of an individual's brain health [44]. Many recent neuroimaging studies utilize deep learning perform age prediction and classification tasks concurrently to improve model generalization across various applications [45]. However, these studies typically require a complete separation of the two tasks during model training, necessitating training each task independently with specific modifications to the model architecture. Leveraging the powerful representation capabilities of STIGR, we aim to overcome this limitation by directly utilizing graph-level representations learned from previous classification tasks. STIGR is trained on classification tasks

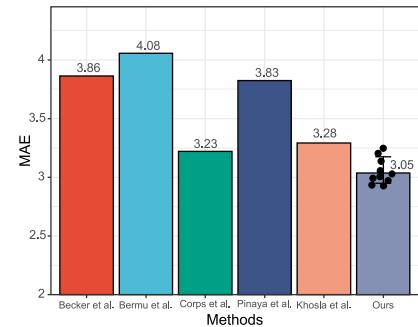


Fig. 10. Performance comparison of brain age prediction on the ABIDE I.

as a feature extractor to obtain high-quality graph representations. These representations are then fed into a three-layer fully connected network for training and subsequently used for age prediction. This approach enables training a simple regression head network instead of the entire model, as done in previous studies. As depicted in Fig. 10, our method achieves competitive performance compared to recent research on brain age prediction in the ABIDE I dataset [44], [45], [46], with a mean absolute error (MAE) of  $3.05(\pm 0.28)$  years. It is worth noting that this comparison is not conducted under identical conditions, as these methods utilize different data volumes and validation approaches. Our method outperforms the runner-up with a decrease in MAE of 0.18 years, showcasing the effectiveness of leveraging learned graph representations to reduce training time for regression tasks while achieving superior results. We also conduct brain age prediction on ABIDE II and ADHD-200 with MAEs of  $3.69(\pm 0.39)$  and  $2.62(\pm 0.11)$ , respectively. In summary, STIGR is a versatile framework that can effectively learn graph representations for different tasks.

#### F. Interpretation Analysis

In addition to the global property of graph subjects explored in previous experiments, the rich information encoded in the node and edge features of the learned graph representation remains crucial for a comprehensive neuroimaging analysis. As shown in Fig. 2, an ideal graph representation learning framework should support a variety of graph tasks, spanning node-level, edge-level, and global-level analyses. Therefore, we delve into STIGR's capacity to learn informative node-level and edge-level representations for model interpretation. This approach allows us to directly probe the model's focus on specific brain regions and their connections, offering a biologically interpretable and scalable window into its decision-making process. The fused graph representation  $R$  (feature map), generated by the gated fusion layer through forward propagation, serves as a representation that preserves essential input patterns [47]. The magnitude of the activations in  $R$  corresponding to each ROI/electrode can be directly interpreted as a measure of its importance. Higher magnitudes suggest greater relevance to the downstream task. This is because graph convolutions and attention mechanisms amplify the signal from important nodes [48], leading to higher activations

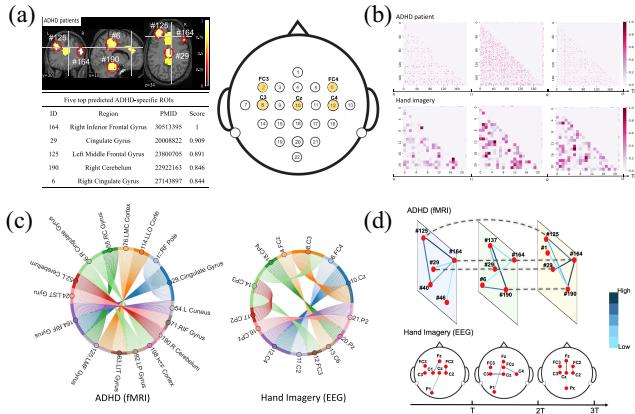


Fig. 11. (a) and (b) visualize the discriminative ROIs (fMRI)/channels (EEG) obtained from static and dynamic node-level representations, respectively. (c) and (d) illustrate the visualization of discriminative connections obtained from static and dynamic edge-level representations.

in  $R$ . Mathematically, this process can be formulated as

$$\text{Score}_{j,c} = \frac{1}{N_c} \sum_{i=1}^{N_c} R_{i,j}^c \quad (12)$$

where  $N_c$  is the number of subjects in class  $c$  and  $R_{i,j}^c$  represents the activation of the  $j$ th ROI/electrode for the  $i$ th subject in class  $c$ , obtained as the output of our well-trained model without the tower network. To illustrate the practical potential of the learned representations for model interpretation, we conduct case studies on two datasets: 1) ADHD-200 for fMRI and 2) BCICIV\_2a for EEG. For detailed interpretability analysis of the other datasets, please refer to the supplementary Section VI.

In static node-level analysis, attention scores for each ROI are computed by aggregating their saliency values from feature maps. Using the weighted sum, we determine the relative importance of each ROI and visualize the top-5 ROIs/channels with the highest attention scores [Fig. 11(a)]. The left panel displays ADHD-related ROIs inferred by the learned static node-level representation. Manual investigation validated the five predicted ADHD-specific ROIs: *Right Inferior Frontal Cortex* (164), *Cingulate Gyrus* (29), *Left Middle Frontal Gyrus* (125), *Right Cerebellum* (190), and *Right Cingulate Gyrus* (6), associated with neurological manifestations of ADHD. For example, Rubia et al. [49] reported dysfunction in *right inferior frontal cortex* (164) of ADHD subjects, a key center of cognitive control. Similarly, the right panel shows the top five channels essential for detecting changes in EEG signal during the motor imagery task (hand versus feet):  $C3$ ,  $C4$ ,  $FC3$ ,  $FC4$ , and  $Cz$ , aligning with prior studies [50]. In the dynamic node-level analysis, we partition the feature map into three time segments. The mean FC matrices are calculated for each segment to measure the correlations between nodes in the neuroimaging data. These matrices are then visualized to track dynamic changes in node-level significance over time, as illustrated in Fig. 11(b). In particular, this visualization demonstrates that the significant ROIs and channels identified

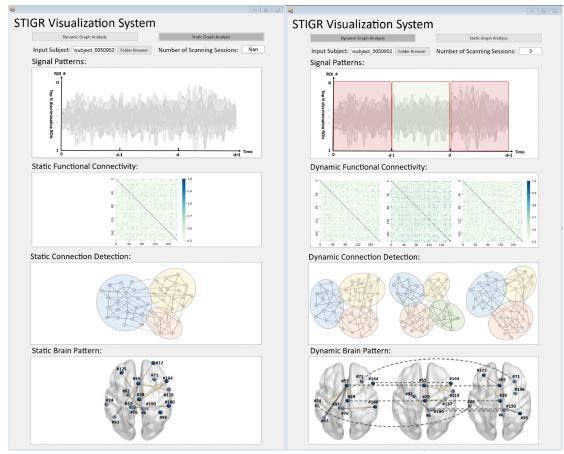


Fig. 12. Example demonstration of the STIGR visualization system, showcasing its potential in aiding the analysis of individual-level neurological patterns for ADHD.

in the static analysis retain their significance in the dynamic analysis. This consistency reinforces the reliability of our findings, validating the credibility of results in both static and dynamic node-level analyses.

In addition, in the static edge-level analysis, we derive the mean FC using the obtained feature maps, quantifying the overall connectivity between different ROIs. The top-10 FC connections with the highest values are visualized, shown in Fig. 11(c). We find that the strongest connectivity was observed between the right inferior frontal cortex region (164) and the cingulate gyrus (29) in participants with ADHD, which is consistent with a recent study [49] that reported increased positive functional connectivity between these two regions in adolescents with ADHD. For the EEG task of hand imaging, the connections between  $C3$  and  $Cz$ , as well as between  $C4$  and  $Cz$ , were found to be the most significant, in agreement with the results reported in [51]. Interestingly, all five of the top ADHD-specific ROIs and hand-specific channels in the static node-level representation were also identified in the top 10 discriminative FC-connected regions in the static edge-level representation, indicating strong consistency.

Similarly to dynamic node-level analysis, we also partition the obtained feature map into three distinct time segments. The corresponding FC matrices are calculated and specific top-5 FC connections are identified that exhibited significance in the ADHD/hand imagery group but not in the control (TC/foot imagery). In Fig. 11(d), we present the dynamic edge connection patterns and demonstrate the efficiency of exploring the evolution of these dynamic interaction brain networks over time to identify reproducible dynamic edge connection patterns as potential neuroimaging biomarkers. Specifically, our statistical analysis revealed that only for this particular pattern [the upper panel of Fig. 11(d)], 35.6% (129/362) of ADHD subjects in the ADHD-200 dataset exhibited a similar dynamic FC pattern. Furthermore, in the BCICIV\_2a dataset, 40.3% (58/144) of the samples showed a similar dynamic connection pattern, as illustrated in the lower panel of Fig. 11(d). These findings suggest that the dynamic activity patterns identified by STIGR have great potential to advance

understanding of brain activity mechanisms. To validate the effectiveness of our graph representation-based interpretation analysis, we also compare the results with the traditional layer-wise relevance propagation method in the supplementary material.

Finally, we have developed a visualization system for subject-level diagnosis, integrating the analysis approaches detailed above. Illustrated in Fig. 12, our trained STIGR effectively captures significant nodes and edges from static and dynamic viewpoints, allowing exploration of ADHD-associated neurological patterns when selecting an input subject (an ADHD subject sample is shown). Specifically, we explore FC connectome at the node-level, community detection, and edge-level connection patterns within this system. In particular, these individual-level results align substantially with the ADHD-specific findings highlighted in Fig. 11. Our visualization system is designed to improve the interpretation of learned graph representations, offering medical professionals also an interpretable tool to explore task-specific neurological patterns at the individual level. In essence, this framework supports informed and data-driven medical decision-making.

## V. CONCLUSION

In this article, we have proposed STIGR, an effective graph representation learning framework for capturing discriminative spatio-temporal graph representations of neuroimaging data. STIGR combines both local and global spatio-temporal dynamics through two major components: 1) L-DGCN and 2) G-STAN. L-DGCN performs graph convolution operations while simultaneously exploring spatio-temporal dependencies by considering the temporal connection of dynamic graphs; G-STAN is developed to learn the missing global spatio-temporal representation and fuse it with L-DGCN to enhance the overall dynamic graph representation ability. In addition, we have designed a CL-AM module based on contrastive learning to adaptively learn the adjacent node matrix for L-DGCN.

Extensive experiments across three neuroimaging datasets (fMRI, EEG, and fNIRS) demonstrate STIGR's superiority over current SOTA methods on benchmark datasets. Ablation studies confirm the contributions of the proposed components. The learned graph representations are versatile and applicable to various downstream tasks, including classification, regression, and interpretation at node, edge, and graph levels. Overall, our framework offers a promising solution for effective graph representation learning in multiple spatio-temporal neuroimaging contexts. In the future, we expect to combine STIGR with multimodality learning to explore the feasibility of multimodal graph representations in the context of neuroimaging.

## REFERENCES

- [1] C. Lian, M. Liu, Y. Pan, and D. Shen, "Attention-guided hybrid network for dementia diagnosis with structural MR images," *IEEE Trans. Cybern.*, vol. 52, no. 4, pp. 1992–2003, Apr. 2022.
- [2] P. Zhong, D. Wang, and C. Miao, "EEG-based emotion recognition using regularized graph neural networks," *IEEE Trans. Affect. Comput.*, vol. 13, no. 3, pp. 1290–1301, Jul.–Sep. 2020.
- [3] R. Liu, Z.-A. Huang, Y. Hu, Z. Zhu, K.-C. Wong, and K. C. Tan, "Attention-like multimodality fusion with data augmentation for diagnosis of mental disorders using MRI," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 35, no. 6, pp. 7627–7641, Jun. 2024.
- [4] D. Ahmedt-Aristizabal, M. A. Armin, S. Denman, C. Fookes, and L. Petersson, "Graph-based deep learning for medical diagnosis and analysis: Past, present and future," *Sensors*, vol. 21, no. 14, p. 4758, 2021.
- [5] Y. Kong et al., "Multi-stage graph fusion networks for major depressive disorder diagnosis," *IEEE Trans. Affect. Comput.*, vol. 13, no. 4, pp. 1917–1928, Oct.–Dec. 2022.
- [6] S. Wein, A. Schüller, A. M. Tomé, W. M. Malloni, M. W. Greenlee, and E. W. Lang, "Forecasting brain activity based on models of spatiotemporal brain dynamics: A comparison of graph neural network architectures," *Netw. Neurosci.*, vol. 6, no. 3, pp. 665–701, 2022.
- [7] X. Cao et al., "scPriorGraph: Constructing biosemantic cell–cell graphs with prior gene set selection for cell type identification from scRNA-seq data," *Genome Biol.*, vol. 25, no. 1, p. 207, 2024.
- [8] C. Dai et al., "Brain EEG time-series clustering using maximum-weight clique," *IEEE Trans. Cybern.*, vol. 52, no. 1, pp. 357–371, Jan. 2020.
- [9] Q. Yu, R. Wang, J. Liu, L. Hu, M. Chen, and Z. Liu, "GNN-based depression recognition using spatio-temporal information: A fNIRS study," *IEEE J. Biomed. Health Inform.*, vol. 26, no. 10, pp. 4925–4935, Oct. 2022.
- [10] J. Wang, Q. Wang, H. Zhang, J. Chen, S. Wang, and D. Shen, "Sparse multiview task-centralized ensemble learning for ASD diagnosis based on age-and sex-related functional connectivity patterns," *IEEE Trans. Cybern.*, vol. 49, no. 8, pp. 3141–3154, Aug. 2019.
- [11] Z.-A. Huang, R. Liu, Z. Zhu, and K. C. Tan, "Multitask learning for joint diagnosis of multiple mental disorders in resting-state fMRI," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 35, no. 6, pp. 8161–8175, Jun. 2024.
- [12] L. Liu et al., "BrainTGL: A dynamic graph representation learning model for brain network analysis," *Comput. Biol. Med.*, vol. 153, Feb. 2023, Art. no. 106521.
- [13] P. Yang et al., "Fused sparse network learning for longitudinal analysis of mild cognitive impairment," *IEEE Trans. Cybern.*, vol. 51, no. 1, pp. 233–246, Jan. 2021.
- [14] R. Liu, Z.-A. Huang, M. Jiang, and K. C. Tan, "Multi-LSTM networks for accurate classification of attention deficit hyperactivity disorder from resting-state fMRI data," in *Proc. 2nd Int. Conf. Ind. Artif. Intell. (IAI)*, 2020, pp. 1–6.
- [15] G. Bertasius, H. Wang, and L. Torresani, "Is space-time attention all you need for video understanding?" in *Proc. 38th Int. Conf. Mach. Learn.*, vol. 2, 2021, p. 4.
- [16] R. Liu, Z.-A. Huang, Y. Hu, Z. Zhu, K.-C. Wong, and K. C. Tan, "Spatial–temporal co-attention learning for diagnosis of mental disorders from resting-state fMRI data," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 35, no. 8, pp. 10591–10605, Aug. 2024.
- [17] B.-H. Kim, J. C. Ye, and J.-J. Kim, "Learning dynamic graph representation of brain connectome with spatio-temporal attention," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 34, 2021, pp. 4314–4327.
- [18] S. Gadgil, Q. Zhao, A. Pfefferbaum, E. V. Sullivan, E. Adeli, and K. M. Pohl, "Spatio-temporal graph convolution for resting-state fMRI analysis," in *Proc. Int. Conf. Med. Image Comput. Comput.-Assist. Intervent.*, 2020, pp. 528–538.
- [19] Y. Li, Y. Liu, Y.-Z. Guo, X.-F. Liao, B. Hu, and T. Yu, "Spatio-temporal-spectral hierarchical graph convolutional network with semisupervised active learning for patient-specific seizure prediction," *IEEE Trans. Cybern.*, vol. 52, no. 11, pp. 12189–12204, Nov. 2022.
- [20] Q. Zhu et al., "Spatio-temporal graph hubness propagation model for dynamic brain network classification," *IEEE Trans. Med. Imag.*, vol. 43, no. 6, pp. 2381–2394, Jun. 2024.
- [21] D. Yang et al., "A novel Spatio-temporal hub identification in brain networks by learning dynamic graph embedding on Grassmannian manifolds," *IEEE Trans. Med. Imag.*, early access, Nov. 19, 2024, doi: [10.1109/TMI.2024.3502545](https://doi.org/10.1109/TMI.2024.3502545).
- [22] F. Hu, L. Zhang, X. Yang, and W.-A. Zhang, "EEG-based driver fatigue detection using spatio-temporal fusion network with brain region partitioning strategy," *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 8, pp. 9618–9630, Aug. 2024.
- [23] L. Chen et al., "An explainable spatio-temporal graph convolutional network for the biomarkers identification of ADHD," *Biomed. Signal Process. Control*, vol. 99, Jan. 2025, Jan. 106913.

- [24] C. Song, Y. Lin, S. Guo, and H. Wan, "Spatial-temporal synchronous graph convolutional networks: A new framework for spatial-temporal network data forecasting," in *Proc. AAAI Conf. Artif. Intell.*, vol. 34, 2020, pp. 914–921.
- [25] R. Liu, Z.-A. Huang, Y. Hu, L. Huang, K.-C. Wong, and K. C. Tan, "Spatio-temporal hybrid attentive graph network for diagnosis of mental disorders on fMRI time-series data," *IEEE Trans. Emerg. Topics Comput. Intell.*, vol. 8, no. 6, pp. 4046–4058, Dec. 2024.
- [26] M. Gadgil, E. Peterson, J. Tregellas, S. Hepburn, and D. C. Rojas, "Differences in global and local level information processing in autism: An fMRI investigation," *Psychiatry Res.*, vol. 213, no. 2, pp. 115–121, 2013.
- [27] A. Vaswani et al., "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 30, 2017, pp. 1–15.
- [28] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," 2018, *arXiv:1810.04805*.
- [29] A. Van Den Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," 2018, *arXiv:1807.03748*.
- [30] S. U. Amin, M. Alsulaiman, G. Muhammad, M. A. Bencherif, and M. S. Hossain, "Multilevel weighted feature fusion using convolutional neural networks for EEG motor imagery classification," *IEEE Access*, vol. 7, pp. 18940–18950, 2019.
- [31] J.-S. Bang, M.-H. Lee, S. Fazli, C. Guan, and S.-W. Lee, "Spatio-spectral feature representation for motor imagery classification using convolutional neural networks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 33, no. 7, pp. 3038–3049, Jul. 2022.
- [32] P. Gaur, H. Gupta, A. Chowdhury, K. McCreadie, R. B. Pachori, and H. Wang, "A sliding window common spatial pattern for enhancing motor imagery classification in EEG-BCI," *IEEE Trans. Instrum. Meas.*, vol. 70, pp. 1–9, 2021.
- [33] K. Das and R. B. Pachori, "Electroencephalogram based motor imagery brain computer interface using multivariate iterative filtering and spatial filtering," *IEEE Trans. Cogn. Develop. Syst.*, vol. 15, no. 3, pp. 1408–1418, Sep. 2023.
- [34] Y. Hou, T. Chen, X. Lun, and F. Wang, "A novel method for classification of multi-class motor imagery tasks based on feature fusion," *Neurosci. Res.*, vol. 176, pp. 40–48, Mar. 2022.
- [35] I. Siviero, G. Menegaz, and S. F. Storti, "Functional connectivity and feature fusion enhance multiclass motor-imagery brain-computer interface performance," *Sensors*, vol. 23, no. 17, p. 7520, 2023.
- [36] P. Khan, P. Ranjan, and S. Kumar, "Data heterogeneity mitigation in healthcare robotic systems leveraging the Nelder-Mead method," in *Artificial Intelligence for Future Generation Robotics*. Amsterdam, The Netherlands: Elsevier, 2021, pp. 71–82.
- [37] P. C. Petrantonakis and I. Kompatsiaris, "Single-trial NIRS data classification for brain-computer interfaces using graph signal processing," *IEEE Trans. Neural Syst. Rehabil. Eng.*, vol. 26, no. 9, pp. 1700–1709, Sep. 2018.
- [38] J. Lu, H. Yan, C. Chang, and N. Wang, "Comparison of machine learning and deep learning approaches for decoding brain computer interface: An fNIRS study," in *Proc. Int. Conf. Intell. Inf. Process.*, 2020, pp. 192–201.
- [39] Z. Wang, J. Zhang, X. Zhang, P. Chen, and B. Wang, "Transformer model for functional near-infrared spectroscopy classification," *IEEE J. Biomed. Health Inform.*, vol. 26, no. 6, pp. 2559–2569, Jun. 2022.
- [40] J. Han, J. Lu, J. Lin, S. Zhang, and N. Yu, "A functional region decomposition method to enhance fNIRS classification of mental states," *IEEE J. Biomed. Health Inform.*, vol. 26, no. 11, pp. 5674–5683, Nov. 2022.
- [41] U. Mahmood, Z. Fu, V. Calhoun, and S. Plis, "Attend to connect: End-to-end brain functional connectivity estimation," in *Proc. Workshop Geom. Topol. Represent. Learn.*, 2021, pp. 1–8.
- [42] G. Wen et al., "Graph self-supervised learning with application to brain networks analysis," *IEEE J. Biomed. Health Inform.*, vol. 27, no. 8, pp. 4154–4165, Aug. 2023.
- [43] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 1263–1272.
- [44] B. G. Becker, T. Klein, C. Wachinger, and A. D. N. Initiative, "Gaussian process uncertainty in age estimation as a measure of brain abnormality," *NeuroImage*, vol. 175, pp. 246–258, Jul. 2018.
- [45] C. Bermudez et al., "Anatomical context improves deep learning on the brain age estimation task," *Magn. Reson. Imag.*, vol. 62, pp. 70–77, Oct. 2019.
- [46] J. Corps and I. Rekik, "Morphological brain age prediction using multi-view brain networks derived from cortical morphology in healthy and disordered participants," *Sci. Rep.*, vol. 9, no. 1, pp. 1–10, 2019.
- [47] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 2921–2929.
- [48] P. E. Pope, S. Kolouri, M. Rostami, C. E. Martin, and H. Hoffmann, "Explainability methods for graph convolutional neural networks," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 10772–10781.
- [49] K. Rubia et al., "Functional connectivity changes associated with fMRI neurofeedback of right inferior frontal cortex in adolescents with ADHD," *NeuroImage*, vol. 188, pp. 43–58, Mar. 2019.
- [50] X. Ma, S. Qiu, and H. He, "Multi-channel EEG recording during motor imagery of different joints from the same limb," *Sci. Data*, vol. 7, no. 1, p. 191, 2020.
- [51] K. K. Ang, Z. Y. Chin, C. Wang, C. Guan, and H. Zhang, "Filter bank common spatial pattern algorithm on BCI competition IV datasets 2a and 2b," *Front. Neurosci.*, vol. 6, p. 39, Mar. 2012.