
Pre-Training Graph Contrastive Masked Autoencoders are Strong Distillers for EEG

Xinxu Wei¹ Kanhao Zhao² Yong Jiao² Hua Xie³ Lifang He⁴ Yu Zhang^{2,1}

Abstract

Effectively utilizing extensive unlabeled high-density EEG data to improve performance in scenarios with limited labeled low-density EEG data presents a significant challenge. In this paper, we address this challenge by formulating it as a graph transfer learning and knowledge distillation problem. We propose a Unified Pre-trained Graph Contrastive Masked Autoencoder Distiller, named EEG-DisGCMAC, to bridge the gap between unlabeled and labeled as well as high- and low-density EEG data. Our approach introduces a novel unified graph self-supervised pre-training paradigm, which seamlessly integrates the graph contrastive pre-training with the graph masked autoencoder pre-training. Furthermore, we propose a graph topology distillation loss function, allowing a lightweight student model trained on low-density data to learn from a teacher model trained on high-density data during pre-training and fine-tuning. This method effectively handles missing electrodes through contrastive distillation. We validate the effectiveness of EEG-DisGCMAC across four classification tasks using two clinical EEG datasets with abundant data. The source code is available at https://github.com/weixinxu666/EEG_DisGCMAC.

1. Introduction

Electroencephalography (EEG) is a pivotal tool for elucidating neural dysfunctions, making it indispensable for the

clinical diagnosis of brain disorders (Sanei & Chambers, 2013). Manual analysis of resting-state EEG (rs-EEG) signals often suffers from low accuracy due to their inherent complexity. In contrast, computer-aided diagnostic methods offer substantial improvements in diagnostic performance. Traditional methods typically involve extracting temporal and spatial features from EEG signals and applying machine learning techniques to develop effective classifiers (Trivedi et al., 2016). Recent advances have seen deep graph learning revolutionize EEG signal analysis. Instead of treating EEG data as conventional numerical inputs, researchers now represent it as non-Euclidean graph data. Graph Neural Networks (GNNs) (Kipf & Welling, 2016) are employed to capture the intricate features and topological structures inherent in these graphs. This innovative approach has markedly enhanced the accuracy and reliability of EEG-based diagnostics, showcasing the potential of GNNs in advancing applications (Song et al., 2018).

Despite these advancements, several critical issues remain unresolved. Firstly, acquiring a substantial amount of accurately labeled clinical rs-EEG data for supervised training on a specific task is challenging due to the complexities involved in data collection (Siuly et al., 2016). Models trained on these limited labeled datasets often exhibit poor accuracy and generalization (Lashgari et al., 2020). Thus, a significant but underexplored research problem is how to effectively utilize this vast amount of unlabeled data to enhance model performance and robustness (Tang et al., 2021). Secondly, the performance of EEG devices varies markedly with the precision of the data they capture. High-density (HD) EEG devices, with their extensive array of electrodes, record high-resolution brain signals, greatly improving the accuracy of diagnostic tasks (Stoyell et al., 2021). However, these devices are often prohibitively expensive and cumbersome, limiting their practical deployment. Conversely, low-density (LD) EEG devices, which are more affordable and easier to deploy (Justesen et al., 2019), capture lower-resolution signals, thus reducing diagnostic accuracy (Cataldo et al., 2022). Addressing how to leverage rich information from HD EEG to enhance diagnostic performance with LD EEG, which is more portable, is crucial for making LD EEG-based diagnostics more accessible and practical (Kuang et al., 2021).

¹Department of Electrical and Computer Engineering, Lehigh University, Bethlehem, PA, USA ²Department of Bioengineering, Lehigh University, Bethlehem, PA, USA ³Center for Neuroscience Research, Children’s National Hospital, Washington, DC, USA ⁴Department of Computer Science and Engineering, Lehigh University, Bethlehem, PA, USA. Correspondence to: Yu Zhang and Lifang He (co-corresponding authors) <{yuzi20, lih319}@lehigh.edu>.

In this paper, we address these challenges through a series of innovative methods. We construct graphs from EEG data and apply GNNs to extract topological features and train the model effectively. To leverage unlabeled EEG data to enhance performance on limited labeled data, we frame this as a Graph Transfer Learning (GTL) problem. We propose a graph self-supervised pre-training approach (Xie et al., 2022) on a large volume of heterogeneous unlabeled EEG graphs. This pre-trained model is subsequently fine-tuned on the scarce labeled data, allowing knowledge acquired from the extensive unlabeled dataset to improve performance on the labeled data. We introduce a novel unified graph self-supervised pre-training paradigm, GCMAP-PT, which combines Graph Contrastive Pre-training (GCL-PT) (Qiu et al., 2020) with Graph Masked Autoencoder Pre-training (GMAE-PT) (Hou et al., 2022). This approach integrates contrastive and generative pre-training by reconstructing contrastive samples and contrasting the reconstructed samples, enabling them to jointly supervise and optimize each other, thereby enhancing overall model performance. To improve model performance with HD EEG data when training on LD EEG data, we address this as a Graph Knowledge Distillation (GKD) problem (Yang et al., 2020) and design a Graph Topology Distillation (GTD) loss function. This allows a student model trained on LD EEG to learn from a teacher model with HD EEG by accounting for missing electrodes through contrastive distillation, while simultaneously compressing model parameters. Moreover, to ensure that models pre-trained with GTL excel as distillers in downstream GKD tasks, we integrate GTL and GKD by contrasting the queries and keys of the teacher and student models during the GTL pre-training process. This integration demonstrates that our unified pre-trained graph contrastive masked autoencoders serve as effective distillers, providing a robust solution for EEG analysis.

2. Related Works

2.1. Graph Neural Networks for EEG Modeling

Recent advancements in Graph Neural Networks (GNNs) have demonstrated their potential in enhancing the modeling and interpretation of EEG data. Notably, the Dynamical Graph Convolutional Neural Network (DGCNN) (Song et al., 2018) was introduced to improve emotion recognition by dynamically learning the interrelationships among EEG channels. Similarly, the Regularized Graph Neural Network (RGNN) (Zhong et al., 2020) applied a regularization strategy to advance emotion recognition from EEG data. Liu et al. (Liu et al., 2023) tackled a similar problem by developing a novel method for emotion recognition from few-channel EEG signals, integrating deep feature aggregation with transfer learning. For medical EEG field, Tang et al. (Tang et al., 2021) employed self-supervised GNNs

to advance seizure detection and classification, achieving significant improvements in identifying rare seizure types.

2.2. Self-Supervised Graph Pre-Training

Self-supervised learning (SSL) pre-training (Zhang et al., 2022b) has proven effective in harnessing extensive unlabeled datasets. Two predominant SSL methods are contrastive learning-based (CL-PT) pre-training, originating from computer vision, and generative-based masked autoencoders (MAE-PT) pre-training, adapted from natural language processing (NLP). These pre-training techniques have been extended to graph models. For instance, GCC (Qiu et al., 2020), GraphCL (You et al., 2020) and GRACE (Zhu et al., 2020) were among the pioneers in applying contrastive learning to graphs by leveraging graph augmentation to generate sample pairs and construct contrastive losses. Concurrently, GraphMAE (Hou et al., 2022), GraphMAE2 (Hou et al., 2023) and GPT-GNN (Hu et al., 2020) adapted the generative masked pre-training approach from NLP (Devlin et al., 2018) to graphs. These methods involve masking nodes and edges, followed by reconstruction, enabling graphs to capture and refine local topological features. Although these methods have been widely applied in CV, NLP, and graphs domains, only CMAE (Huang et al., 2023) has combined CL-based and MAE-based methods, and it has only been applied in the CV field. To the best of our knowledge, no method has elegantly integrated GCL-PT and GMAE-PT methods in the graph domain.

2.3. Graph Knowledge Distillation

Graph knowledge distillation focuses on transferring knowledge from a complex, large-scale model (teacher) to a more streamlined and efficient model (student), thus preserving performance while reducing computational demands. G-CRD (Joshi et al., 2022) introduced a distillation loss function for GNN-to-GNN transfer, employing a contrastive learning strategy to enhance similarity among nodes of the same class and increase separation between different classes. MSKD (Zhang et al., 2022a) proposed a multi-teacher distillation approach, integrating various teacher GNN models of different scales into a single student GNN model. Approaches such as Graph-MLP (Hu et al., 2021), and VQ-Graph (Yang et al., 2024) focused on transferring knowledge from structure-aware teacher GNNs to structure-agnostic student MLPs.

3. Methodology

3.1. EEG Graph Construction

An EEG graph can be formally represented as $\mathcal{G} = (\mathcal{V}, \mathcal{A}, \mathcal{X})$. The matrix $\mathcal{X} \in \mathbb{R}^{n \times d}$ represents the node features, with n indicating the number of nodes (or electrodes)

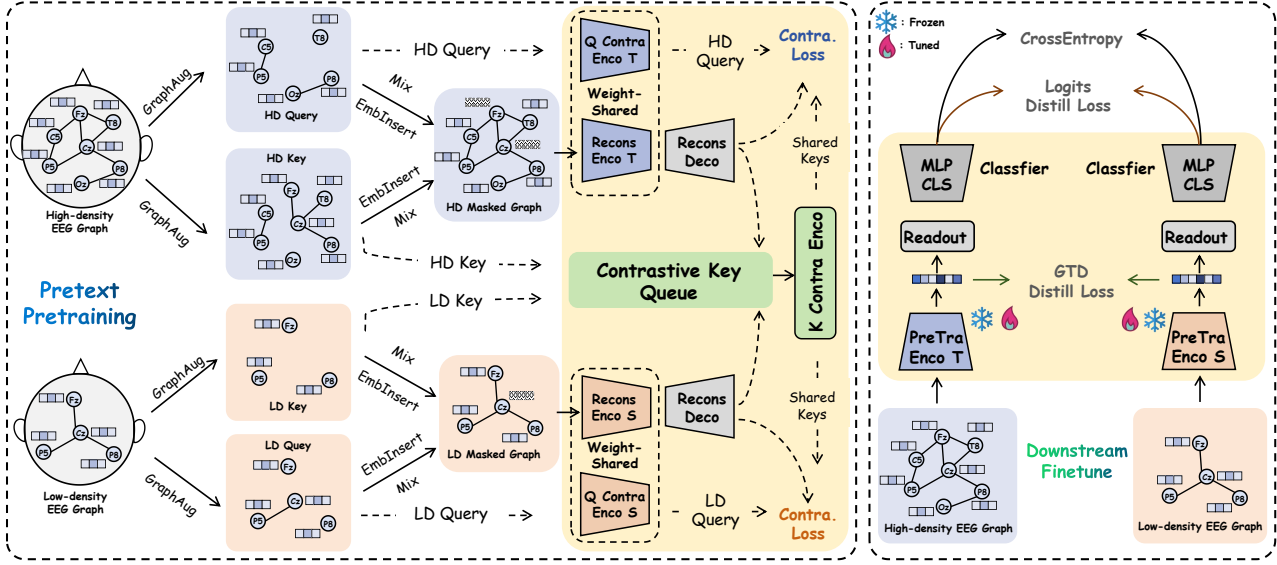


Figure 1: The proposed EEG-DisGCMC framework consists of two main stages: a pretext pre-training (PT) stage and a downstream fine-tuning (FT) stage. Note that we can perform two types of fine-tuning: ‘Tuned’ refers to fine-tuning all the parameters of the model, while ‘Frozen’ means freezing most layers of the model and only fine-tuning the parameters of the top-level layers. Note that the encoders (Enco) of our model can be adopted *Graph Transformer* or *vanilla GCNs*.

and d specifying the dimensionality of the feature vector associated with each node. $\mathcal{A} \in \mathbb{R}^{n \times n}$ denotes the adjacency matrix. The EEG graph is derived from the original EEG time series signals recorded by EEG caps, where n represents the number of channels or electrodes. To convert resting-state EEG (rs-EEG) time series into graph representations, we first apply band-pass filtering to extract EEG signals within the following frequency bands: θ (4-8 Hz), α (8-14 Hz), β (14-30 Hz), and γ (30-50 Hz). Subsequently, we compute the power spectral density (PSD) features for each band, selecting the α band for this study. These PSD features are utilized as node features for the EEG graph. The Pearson correlation is then computed between nodes to construct the adjacency matrix, which represents the edge connectivity.

3.2. Unified Graph Pre-Training for Distillation

To fully leverage the extensive amount of unlabeled EEG data, we propose a graph self-supervised pre-training approach to pre-train EEG models from the graph-based perspective. Our *motivation* stems from the observation that prior research has predominantly focused on either contrastive-based or generative-based pre-training methods for EEG time series, with limited studies addressing these techniques within the context of EEG graph models. To address this gap, we introduce a unified graph self-supervised pre-training paradigm, termed GCMC-PT, based on the following assumptions:

Assumption 1: (Combining GCL and GMAE for Enhanced

Distillation) Hybridizing contrastive-based and generative-based pre-training by simultaneously reconstructing contrastive pairs and contrasting the reconstructed samples provides a more robust distiller, rather than applying these methods separately or in sequence.

Assumption 2: (Joint Pre-Training of Teacher and Student Models) Both the teacher and student models benefit from joint pre-training through the contrasting of each other’s positive and negative pairs, leading to improved distillation performance.

Consider two types of EEG graph inputs: the high-density EEG graph $\mathcal{G}^h = (\mathcal{V}^h, \mathcal{A}^h, \mathcal{X}^h) \in \mathbb{R}^{m \times d}$ and the low-density EEG graph $\mathcal{G}^l = (\mathcal{V}^l, \mathcal{A}^l, \mathcal{X}^l) \in \mathbb{R}^{n \times d}$, where m and n represent the number of nodes (or electrodes) in \mathcal{G}^h and \mathcal{G}^l , respectively, and $m \geq n$. Note that \mathcal{G}^l can be regarded as a subgraph of \mathcal{G}^h . Additionally, two graph encoders are employed: a teacher graph encoder with extensive parameters and robust feature extraction capabilities, and a lightweight student graph encoder with fewer parameters and comparatively lower learning capacity. It is noteworthy that \mathcal{A}^h and \mathcal{A}^l can be dynamically learned and adjusted throughout the training process. The teacher and student models are adaptable to different types of GNNs, such as transductive spectral-based traditional GCNs (like DGCNN(Song et al., 2018)) or spatial-based graph transformers (Yun et al., 2019).

Since \mathcal{V}^l is derived from \mathcal{V}^h , that is $\mathcal{V}^l \subseteq \mathcal{V}^h$, we partition the complete node set \mathcal{V}^h (the Complete/HD Set) in \mathcal{G}^h into

the Deleted Set \mathcal{V}^d and the Remaining/LD Set \mathcal{V}^l . The set \mathcal{V}^d comprises $(m - n)$ nodes present in \mathcal{G}^h but absent from \mathcal{G}^l , representing the removed electrodes/nodes. Conversely, \mathcal{V}^l includes the n nodes retained in \mathcal{G}^l . The relationships among these sets can be expressed as $\mathcal{V}^l = \mathcal{V}^h - \mathcal{V}^d$, where $\mathcal{V}^l \subseteq \mathcal{V}^h$, $\mathcal{V}^d \subseteq \mathcal{V}^h$, $\mathcal{V}^d \cap \mathcal{V}^l = \emptyset$, and $\mathcal{V}^d \cup \mathcal{V}^l = \mathcal{V}^h$. Thus, the complete set is composed of the deleted set and the remaining LD set.

As illustrated in Fig. 1, to construct the contrastive-based pre-training paradigm, graph augmentation techniques (You et al., 2020) are initially applied to HD and LD graphs by randomly dropping nodes and removing edges. This process yields *Query graphs* ($\mathcal{Q}^h, \mathcal{Q}^l$) and *Key graphs* ($\mathcal{K}^h, \mathcal{K}^l$). Finally, the total augmented graphs are denoted as $\hat{\mathcal{G}}^h$ and $\hat{\mathcal{G}}^l$. This can be formulated as:

$$\begin{aligned}\hat{\mathcal{G}}^h &= \text{Mix}(\mathcal{Q}^h, \mathcal{K}^h) = \text{Aug}(\mathcal{G}^h) \\ \hat{\mathcal{G}}^l &= \text{Mix}(\mathcal{Q}^l, \mathcal{K}^l) = \text{Aug}(\mathcal{G}^l)\end{aligned}\quad (1)$$

Where $\text{Aug}(\cdot)$ means graph augmentation and $\text{Mix}(\cdot)$ represents the integration of two graph sets.

To achieve the goal of *reconstructing the contrastive pairs* as outlined in Assumption 1, the *masked graphs* for GMAE-PT are constructed from the mixed contrastive augmented samples by substituting the dropped nodes with learnable embeddings. Subsequently, both the teacher and student encoders are employed to encode the masked graphs into the graph embeddings. To accomplish GMAE-PT, graph decoders for both teacher and student encoders are utilized to reconstruct the masked graph embeddings into the original input graphs by applying the *MSE Loss* as the reconstruction loss \mathcal{L}_{Rec} on both the reconstructed node features $\tilde{\mathcal{X}}$ and graph structures $\tilde{\mathcal{A}} = \tilde{\mathcal{X}} \cdot \tilde{\mathcal{X}}^{\text{tr}}$ (Yang et al., 2024).

$$\mathcal{L}_{\text{Rec}} = \left\| \mathcal{X} - \tilde{\mathcal{X}} \right\|_2^2 + \left\| \mathcal{A} - \tilde{\mathcal{X}} \cdot \tilde{\mathcal{X}}^{\text{tr}} \right\|_2^2 \quad (2)$$

where $\tilde{\mathcal{X}}^{\text{tr}}$ means the transpose of $\tilde{\mathcal{X}}$. Then the reconstructed HD and LD query ($\tilde{\mathcal{Q}}^h, \tilde{\mathcal{Q}}^l$) and key ($\tilde{\mathcal{K}}^h, \tilde{\mathcal{K}}^l$) graphs are split out from the reconstructed $\hat{\mathcal{G}}^h$ and $\hat{\mathcal{G}}^l$.

To achieving the goal of *contrasting the reconstructed samples* in Assumption 1, the reconstructed HD and LD query and key graphs are mixed with the original contrastive samples generated by augmentation as additional contrastive HD and LD query and key samples to form the *extended* contrastive HD and LD query ($\mathcal{Q}_{\text{ex}}^h, \mathcal{Q}_{\text{ex}}^l$) and key ($\mathcal{K}_{\text{ex}}^h, \mathcal{K}_{\text{ex}}^l$) samples.

$$\begin{aligned}\mathcal{Q}_{\text{ex}}^h &= \{\mathcal{Q}^h, \tilde{\mathcal{Q}}^h\} & \mathcal{Q}_{\text{ex}}^l &= \{\mathcal{Q}^l, \tilde{\mathcal{Q}}^l\} \\ \mathcal{K}_{\text{ex}}^h &= \{\mathcal{K}^h, \tilde{\mathcal{K}}^h\} & \mathcal{K}_{\text{ex}}^l &= \{\mathcal{K}^l, \tilde{\mathcal{K}}^l\}\end{aligned}\quad (3)$$

To achieve the goal of *joint the teacher and student pre-training via contrasting the reconstructed samples* in Assumption 2, the extended key samples of both are mixed to

form a larger *Key Samples Pool* $\mathcal{K}_{\text{ex}}^{\text{hl}}$.

$$\mathcal{K}_{\text{ex}}^{\text{hl}} = \{\mathcal{K}^h, \tilde{\mathcal{K}}^h, \mathcal{K}^l, \tilde{\mathcal{K}}^l\} = KQ(\{\mathcal{K}_{\text{ex}}^{\text{hl}+}, \mathcal{K}_{\text{ex}}^{\text{hl}-}\}) \quad (4)$$

Following (He et al., 2020; Qiu et al., 2020), we adopt a *Key Queue*, denoted as $KQ(\cdot)$, to store a large number of mixed extended *key samples pool* and *key encoders* for both teacher and student to convert $\mathcal{K}_{\text{ex}}^{\text{hl}}$ to be key embeddings for jointly pre-training the teacher and student encoders via a joint contrastive loss function (Qiu et al., 2020) as follows:

$$\begin{aligned}\mathcal{L}_{\text{cl}}^T &= -\log \left(\frac{\exp(\mathcal{Q}_{\text{ex}}^h \cdot \mathcal{K}_{\text{ex}}^{\text{hl}+} / \tau)}{\sum_{i=0}^K \exp(\mathcal{Q}_{\text{ex}}^h \cdot \mathcal{K}_{\text{ex}}^{\text{hl}+} / \tau)} \right) \\ \mathcal{L}_{\text{cl}}^S &= -\log \left(\frac{\exp(\mathcal{Q}_{\text{ex}}^l \cdot \mathcal{K}_{\text{ex}}^{\text{hl}+} / \tau)}{\sum_{i=0}^K \exp(\mathcal{Q}_{\text{ex}}^l \cdot \mathcal{K}_{\text{ex}}^{\text{hl}+} / \tau)} \right)\end{aligned}\quad (5)$$

where τ represents the temperature coefficient.

Then, as shown in Fig. 6 (d) in the supplementary, we simultaneously contrast the extended HD queries and LD queries with all the mixed extended keys in the queue, which consists of both HD and LD keys with the corresponding HD and LD reconstructed keys, to construct the positive and negative pairs with their corresponding positive and negative keys $\{\mathcal{K}_{\text{ex}}^{\text{hl}+}, \mathcal{K}_{\text{ex}}^{\text{hl}-}\}$ in the queue, computing contrastive loss to jointly optimize the query and key encoders of both teacher and student models. Therefore, the joint contrastive loss function for GCL-PT $\mathcal{L}_{\text{cl}}^{\text{Joint}}$ is composed of the teacher contrastive loss $\mathcal{L}_{\text{cl}}^T$ and the student contrastive loss $\mathcal{L}_{\text{cl}}^S$. And the joint reconstruction loss function $\mathcal{L}_{\text{Rec}}^{\text{Joint}}$ for GMAE-PT consists of the teacher reconstruction loss $\mathcal{L}_{\text{Rec}}^T$ and the student reconstruction loss $\mathcal{L}_{\text{Rec}}^S$ as follows:

$$\mathcal{L}_{\text{cl}}^{\text{Joint}} = \mathcal{L}_{\text{cl}}^T + \mathcal{L}_{\text{cl}}^S \quad \mathcal{L}_{\text{Rec}}^{\text{Joint}} = \mathcal{L}_{\text{Rec}}^T + \mathcal{L}_{\text{Rec}}^S \quad (6)$$

The overall loss $\mathcal{L}_{\text{Pretrain}}$ for both the teacher and student encoders pre-training is composed of the contrastive-based loss $\mathcal{L}_{\text{cl}}^{\text{Joint}}$ and the generative-based loss $\mathcal{L}_{\text{Rec}}^{\text{Joint}}$.

$$\mathcal{L}_{\text{Pretrain}} = \mathcal{L}_{\text{cl}}^{\text{Joint}} + \mathcal{L}_{\text{Rec}}^{\text{Joint}} \quad (7)$$

3.3. Graph Topology Distillation for HD-LD EEG

In the downstream stage, the pre-trained models are fine-tuned for specific classification tasks using limited labeled EEG data. We employ the Cross-Entropy loss for classification. To transfer logit-based knowledge, we adopt the classic logit distillation loss $\mathcal{L}_{\text{Dis}}^{\text{logits}}$ (Hinton et al., 2015), using KL divergence to align the predicted logit distributions, allowing the pre-trained student model to mimic the logits of the pre-trained teacher model. Moreover, since \mathcal{G}^h contains more nodes than \mathcal{G}^l , the topological information \mathcal{A}^h learned by the pre-trained teacher model from the high-density graph is more precise and discriminative than \mathcal{A}^l , learned by the pre-trained student model from the low-density graph. These topological features capture the spatial

connectivity of EEG electrodes, which is crucial for task performance. Thus, distilling the topological knowledge from the pre-trained teacher model into the pre-trained student model is essential to boost the performance of the pre-trained student model. To address this, we propose the Graph Topology Distillation loss. To quantify the similarity between node features \mathcal{X}_i of node v_i and \mathcal{X}_j of node v_j in the graph, we employ a similarity kernel function (Joshi et al., 2022). This function computes the similarity \mathcal{Z}_{ij} for both \mathcal{G}^h and \mathcal{G}^l . Specifically, we adopt the *Linear Kernel* as the node similarity function $\mathcal{F}(\cdot)$, defined as follows:

$$\begin{aligned}\mathcal{Z}_{ij}^h &= \mathcal{F}(\mathcal{X}_i^h, \mathcal{X}_j^h) = \mathcal{X}_i^h \cdot \mathcal{X}_j^h \\ \mathcal{Z}_{ij}^l &= \mathcal{F}(\mathcal{X}_i^l, \mathcal{X}_j^l) = \mathcal{X}_i^l \cdot \mathcal{X}_j^l\end{aligned}\quad (8)$$

Note that $\{v_i, v_j\} \in (\mathcal{V}^h \cap \mathcal{V}^l)$ and $\{v_i, v_j\} \notin \mathcal{V}^d$. Guided by positive and negative pairs in \mathcal{A}^h and the influence of the GTD loss, we aim to pull similar positive node pairs \mathcal{P}_{ij}^+ closer and push dissimilar negative node pairs \mathcal{P}_{ij}^- farther apart in \mathcal{A}^l . This process first requires defining and selecting \mathcal{P}_{ij}^+ and \mathcal{P}_{ij}^- for both HD and LD graphs. As described earlier, three node sets are involved: the complete/HD set, the deleted/removed set, and the remaining/LD set. Since LD graphs are formed by removing certain electrodes/channels/nodes \mathcal{V}^d from \mathcal{V}^h , the removed electrodes \mathcal{V}^d significantly influence the topological structure of LD graphs.

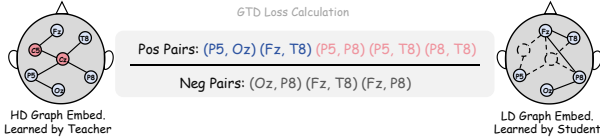


Figure 2: The intuitive diagram for the selection of positive and negative sample pairs in GTD Loss.

Positive and Negative Nodes Selection: As described in the equations above, in LD graphs, if two nodes v_i and v_j are either directly connected (1-hop) or indirectly connected (2-hop) through a removed node v_k in \mathcal{V}^d , acting as a mediator in the graph embedding of HD graphs learned by the teacher model, these node pairs v_{ij} are treated as positive contrastive pairs. Conversely, if node pairs v_{ij} are connected in the embedding LD graphs learned by the student model but are neither directly nor indirectly connected in the embedding learned by the pre-trained teacher model, they are treated as negative contrastive pairs.

As shown in Fig. 2, we provide an intuitive illustration for the selection of positive and negative pairs. In the HD graph, the blue nodes represent the nodes that are retained in the LD graph, while the pink nodes represent those that exist only in the HD graph but are absent in the LD graph (i.e., the information lost in the LD graph, which is what we aim to

distill). In the LD graph, the pink nodes and their associated edges are removed. As shown, (P5, Oz) and (Fz, T8) are directly connected as first-order neighbors in the HD graph, making them positive pairs. Meanwhile, (P5, P8), (P5, T8), and (P8, T8) are second-order neighbors via the pink nodes, so they are also considered positive pairs according to our rules. However, (Oz, P8), (Fz, T8), and (Fz, P8) are neither first-order neighbors nor second-order neighbors via the pink nodes in the HD graph. Yet, they become first-order neighbors in the LD graph, thus qualifying as negative pairs that need correction.

Once the positive and negative pairs are identified, we apply KL divergence as the distillation function. In the numerator, it is used to align the kernel feature distributions learned by the pre-trained teacher and student models for positive pairs, encouraging the pre-trained student model to replicate the topological distribution of the pre-trained teacher model and increase the similarity of positive pairs in the embeddings learned by the student model. In the denominator, KL divergence is also employed to adjust the erroneous topological distribution learned from negative pairs by the student model.

$$\begin{aligned}\mathcal{L}_{Pos} &= \sum_{(i,j) \in \mathcal{P}^+} KL(\text{softmax}(\mathcal{Z}_{ij}^l) \parallel \text{softmax}(\mathcal{Z}_{ij}^h)) \\ \mathcal{L}_{Neg} &= \sum_{(i,j) \in \mathcal{P}^-} KL(\text{softmax}(\mathcal{Z}_{ij}^l) \parallel \text{softmax}(\mathcal{Z}_{ij}^h))\end{aligned}\quad (9)$$

The final GTD loss function \mathcal{L}_{Dis}^{GTD} in the contrastive format is as follows:

$$\mathcal{L}_{Dis}^{GTD} = \frac{\mathcal{L}_{Pos}/\mathcal{C}_{Pos}}{\mathcal{L}_{Neg}/\mathcal{C}_{Neg} + \epsilon} \quad (10)$$

where \mathcal{C}_{Pos} and \mathcal{C}_{Neg} are the counts of \mathcal{P}_{ij}^+ and \mathcal{P}_{ij}^- . ϵ is a constant to avoid division by zero errors.

Finally, we integrate all the loss functions to form the total *Fine-tune* loss $\mathcal{L}_{Finetune}$:

$$\mathcal{L}_{Finetune} = \mathcal{L}_{CE} + \mathcal{L}_{Dis}^{Logits} + \mathcal{L}_{Dis}^{GTD} \quad (11)$$

3.4. Special Case for the Proposed GTD Loss

The GTD loss is primarily designed to distill topological knowledge from \mathcal{G}^h to \mathcal{G}^l . However, there is a special case known as H2H distillation, where \mathcal{G}^l and \mathcal{G}^h have the same number of nodes, meaning $\mathcal{V}^l = \mathcal{V}^h$ and $\mathcal{V}^d = \emptyset$. In this scenario, no nodes are removed, and only the connections in \mathcal{A}^l and \mathcal{A}^h may differ. With slight modifications, our loss function can also be applied to this special case. The modified GTD loss for the H2H distillation scenario is given as follows:

$$\mathcal{P}_{ij}^+ = \mathbb{I}(\mathcal{A}_{ij}^h > 0) \quad \mathcal{P}_{ij}^- = \mathbb{I}(\mathcal{A}_{ij}^l > 0 \text{ and } \mathcal{A}_{ij}^h = 0) \quad (12)$$

In this special case, GTD loss does not consider \mathcal{V}^d . The learning objective becomes utilizing the learned \mathcal{A}^h learned from the teacher model to correct incorrectly edges in \mathcal{A}^l learned from the student model, thereby making \mathcal{A}^l as close to \mathcal{A}^h as possible.

4. Algorithm Pipeline of GTD Loss

To clarify the GTD loss calculation, we present the pipeline as shown in Algorithm 1.

Algorithm 1 GTD Loss Calculation

Input: $\mathcal{X}^l, \mathcal{V}^l, \mathcal{A}^l, \mathcal{X}^h, \mathcal{V}^h, \mathcal{A}^h, \mathcal{V}^d$

Parameter: $\mathcal{F}(\cdot), \theta, \epsilon$

Output: \mathcal{L}_{Dis}^{GTD}

```

1: Normalize  $\mathcal{A}^l, \mathcal{A}^h$ 
2: Apply threshold:  $\mathcal{A}^l \leftarrow (\mathcal{A}^l > \theta), \mathcal{A}^h \leftarrow (\mathcal{A}^h > \theta)$ 
3: Compute kernel matrices:  $\mathcal{Z}^l = \mathcal{F}(\mathcal{X}^l), \mathcal{Z}^h = \mathcal{F}(\mathcal{X}^h)$ 
4: Assert:  $|\mathcal{V}^l| \leq |\mathcal{V}^h|$ 
5: if  $|\mathcal{V}^l| \neq |\mathcal{V}^h|$  then
6:   Extract sub-matrices  $\mathcal{A}_{sub}^h = \mathcal{A}^h[\mathcal{V}^l, \mathcal{V}^l]$ 
7:   Direct connections:  $\mathcal{A}_{1-hop}^h = (\mathcal{A}_{sub}^h > 0)$ 
8:   Indirect connections:  $\mathcal{A}_{2-hop}^h = \mathcal{A}^h[\mathcal{V}^d, : \mathcal{V}^l]$ 
9:    $\mathcal{L}_{Pos} = KL(\mathcal{Z}^l || \mathcal{Z}^h) | (\mathcal{A}_{1-hop}^h \vee \mathcal{A}_{2-hop}^h)$ 
10:   $\mathcal{L}_{Neg} = KL(\mathcal{Z}^l || \mathcal{Z}^h) | (\mathcal{A}^l \wedge \neg(\mathcal{A}_{1-hop}^h \vee \mathcal{A}_{2-hop}^h))$ 
11: else
12:   $\mathcal{L}_{Pos} = KL(\mathcal{Z}^l || \mathcal{Z}^h) | (\mathcal{A}^h > 0)$ 
13:   $\mathcal{L}_{Neg} = KL(\mathcal{Z}^l || \mathcal{Z}^h) | (\mathcal{A}^l > 0 \wedge \mathcal{A}^h = 0)$ 
14: end if
15:  $\mathcal{L}_{PosAvg} = \frac{\mathcal{L}_{Pos}}{C_{pos}}$ 
16:  $\mathcal{L}_{NegAvg} = \frac{\mathcal{L}_{Neg}}{C_{neg}}$ 
17: return  $\mathcal{L}_{Dis}^{GTD} = \frac{\mathcal{L}_{PosAvg}}{\mathcal{L}_{NegAvg} + \epsilon}$ 
    
```

5. Experiments

5.1. Implementation details

During pre-training, we used a batch size of 128. For downstream fine-tuning, we used a batch size of 32. Both pre-training and fine-tuning were optimized using the Adam optimizer.

5.2. EEG Datasets and Downstream Tasks

We evaluated our EEG-DisGCMAE framework on two clinical datasets with rs-EEG time series: the Establishing Moderators and Biosignatures of Antidepressant Response in Clinical Care (EMBARC) (Trivedi et al., 2016) and the Healthy Brain Network (HBN) (Alexander et al., 2017). The EMBARC dataset comprises EEG data from 308 eye-open and 308 eye-closed samples, while the HBN dataset includes 1,594 eye-open and 1,745 eye-closed samples. De-

tailed dataset preprocessing information is provided in the appendices. For EMBARC, we performed binary classification tasks: sex classification in Major Depressive Disorder (MDD) patients (Male vs Female) and depression severity classification based on the Hamilton Depression Rating Scale ($HAMD_{17}$) (Williams, 1988) (Mild vs Severe Depression) (Boessen et al., 2013). For HBN, we conducted binary classifications for MDD (Healthy vs MDD) and Autism Spectrum Disorder (ASD) (Healthy vs ASD). Additional details can be found in the appendices. We tested three EEG electrode density levels: high-density (HD), medium-density (MD), and low-density (LD). In EMBARC, these densities correspond to the 10-20 EEG system electrode distributions of 64 (HD), 32 (MD), and 16 (LD) electrodes, respectively. For HBN, the densities correspond to 128 (HD), 64 (MD), and 32 (LD) electrodes.

5.3. Comparative Experiment Analysis

We compared the proposed EEG-DisGCMAE against five classes of methods: Traditional machine learning methods (MLP, LSTM), GNN-based models (GCN, GFormer, HyperGCN (Feng et al., 2019)), EEG-specific models (EEGNet (Lawhern et al., 2018), DGCNN, EEG-Conformer (Song et al., 2022), RGNN), Graph contrastive pre-training Models (GCC, GraphCL, GRACE), and Graph generative pre-training Models (GraphMAE, GPT-GNN, GraphMAE2). In addition, we compare our method with LaBraM (Jiang et al.), a time-series EEG model that performs pre-training directly on the raw time-series data. As demonstrated in Table 1, our model outperforms all other state-of-the-art methods. Notably, pre-training-based models, including those based on GCL-PT (GCC, GraphCL, GRACE) and GMAE-PT (GraphMAE, GPT-GNN, GraphMAE2), utilize large Graph Transformers as their backbone in this study. In contrast, our method can be suitable to both spatial-based Graph Transofrmer and spectral-based vanilla GCNs (DGCNN) as the backbone. We evaluated both tiny and large model sizes. As illustrated in Fig. 3(a), our tiny model, with only 1.3M parameters, performs comparably to pre-training-based methods with larger models (5.7M parameters). Moreover, our large-tiny model, despite having a similar parameter size to others, significantly outperforms them by about 5% in both AUROC and accuracy. This demonstrates that our approach achieves a superior balance between performance and efficiency, delivering high performance with a more compact parameter set. As illustrated in Fig. 3(b), we investigated the relationship between model parameters and performance across three factors: model size, model type, and varying input EEG densities. It is evident that when the model type and input EEG density are fixed, the large-size model outperforms the tiny-size model. For a given model, reducing the input density (i.e., using LD data) leads to a decline in performance compared to using

Table 1: Performance comparison of different methods on two clinical EEG datasets for different classification tasks. Our teacher and student model can adopt both spectral-based GCNs (DGCNN) or spatial-based Graph Transformer as the backbone, whereas other graph pre-training models utilize large Graph Transformers. The experiments encompass both high-density and low-density EEG scenarios. Metrics are reported as AUROC(%) / ACC(%).

Method	HBN MDD		HBN ASD		EMBARC Sex		EMBARC Severity	
	HD	LD	HD	LD	HD	LD	HD	LD
MLP	73.2/75.7	71.6/72.5	58.3/61.1	56.3/59.4	68.0/71.3	65.7/67.4	61.5/63.7	59.4/62.6
LSTM	76.7/79.2	73.7/76.8	60.3/64.6	58.4/61.8	69.0/71.8	67.3/69.3	62.8/66.0	61.2/63.8
GCN	75.8/77.6	72.3/76.4	60.5/63.7	59.2/61.8	69.1/72.8	66.7/69.6	63.5/66.2	60.7/63.3
GFormer	80.4/83.6	76.3/80.4	62.7/64.2	61.5/62.8	71.8/74.4	68.1/71.6	66.2/69.8	64.7/66.8
Hyper-GCN	77.6/80.8	75.4/77.7	60.1/64.5	59.7/63.1	70.5/73.8	67.6/72.3	64.7/68.3	63.0/66.2
EEGNet	80.6/82.9	76.3/80.1	62.0/64.6	59.3/62.8	71.1/74.0	66.6/70.3	65.4/70.1	62.4/65.2
EEG Conformer	79.3/83.1	77.5/79.8	61.6/64.3	60.3/62.4	72.2/74.8	67.9/70.7	64.7/66.3	64.2/64.8
RGNN	79.4/82.5	76.8/79.2	60.3/62.6	58.4/63.2	71.8/73.5	68.7/71.5	64.7/66.2	62.5/65.1
DGCNN	77.1/81.7	74.2/78.7	61.3/63.8	59.3/62.7	70.6/73.2	66.6/72.3	65.4/67.8	62.7/64.2
GCC	82.2/85.1	80.4/82.8	64.3/66.1	63.2/62.1	72.9/75.7	69.5/73.7	68.5/71.1	66.1/68.7
GRACE	83.7/84.6	79.9/81.8	63.7/66.8	61.6/63.9	73.2/74.9	70.7/72.8	67.3/71.8	66.7/67.6
GraphCL	81.7/83.9	78.6/80.6	64.6/65.4	63.4/64.1	72.8/75.4	68.5/73.5	69.4/72.6	67.3/69.2
GraphMAE	82.8/85.3	79.5/83.3	65.1/64.7	62.5/62.9	72.6/76.3	70.2/73.8	69.4/69.8	65.8/68.5
GPT-GNN	83.3/85.2	80.7/82.2	65.6/66.9	64.3/65.0	71.2/74.7	68.5/72.9	68.3/70.4	65.4/67.1
GraphMAE2	83.5/85.7	81.3/83.0	65.3/65.9	62.5/65.9	72.2/75.6	69.5/73.2	70.5/70.0	66.2/68.3
Ours-Tiny (DGCNN)	84.8/85.4	81.6/82.4	66.1/66.4	63.4/64.1	73.4/76.7	71.8/75.6	68.6/71.9	66.8/69.3
Ours-Large (DGCNN)	86.6/87.4	84.4/85.3	67.3/68.8	66.7/65.9	75.4/77.8	74.5/76.3	71.5/74.6	69.2/72.8
Ours-Tiny (Gformer)	85.3/86.8	82.6/84.3	66.6/67.8	64.7/65.7	75.2/77.1	73.3/75.3	68.7/73.5	67.3/72.1
Ours-Large (Gformer)	87.4/87.8	84.8/86.9	68.6/69.4	66.8/67.4	76.6/77.9	75.7/76.8	72.3/77.2	70.6/74.0

HD data. However, after pre-training and distillation, the performance of the initially less effective tiny-size model improves significantly, reaching a level comparable to that of the large-size teacher model using HD data without pre-training. This demonstrates that our GCMAE-PT and GTD loss can enhance model performance while maintaining a lightweight parameter set without compromising efficiency.

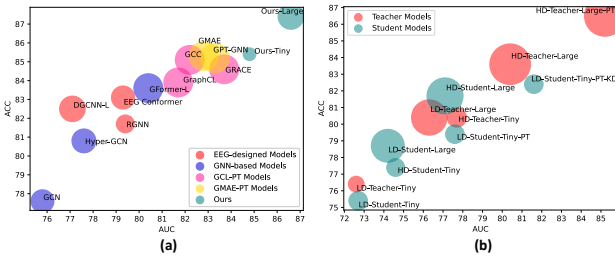


Figure 3: (a) compares model sizes and performance where 'L' denotes large-size models, and (b) examines the same factors across different model types. Both analyses were conducted on the HBN dataset for MDD classification. Models of the same color belong to the same category, and circle size represents the number of parameters, with larger circles indicating higher parameter counts. In (b), the model backbone is a vanilla GCN (DGCNN).

5.3.1. ANALYSIS OF EEG PATTERNS FOR MASKING AND RECONSTRUCTION

To illustrate the effectiveness of our proposed pre-training method, we visualized EEG data patterns across various

densities, masking ratios, and reconstruction methods, as depicted in Fig. 4. Fig. 4(a) shows clear and well-connected activated regions with no masking. As we increased the masking ratio in Figs. 4(b), 4(c), and 4(d), the activated regions diminish and connectivity deteriorates, reflecting increased information loss. Fig. 4(e) demonstrates the effectiveness of our reconstruction method with 50% masking, revealing a pattern that closely resembles the unmasked data in Fig. 4(a), with improved activation and high reconstruction accuracy.

5.4. Ablation Study Analysis

5.4.1. ELECTRODE DENSITY AND MODEL SIZE

Table 2 presents ablation experiments examining EEG graphs with varying densities (HD/MD/LD) and model types (teacher/student) with different sizes (tiny/large). The results reveal that as electrode density decreases, performance on EEG recognition tasks deteriorates. The decline is more pronounced when reducing density from MD to LD than from HD to MD. This is because, while the reduction from HD to MD removes redundant electrodes, MD still retains essential information, preserving performance. However, reducing from MD to LD results in the loss of critical electrodes, leading to a significant performance drop. Additionally, ablation experiments comparing different model sizes, including tiny and large versions of the spatial-based graph transformer and spectral-based DGCNN, indicate that the teacher model consistently outperforms the student model of the same size. The tiny teacher model performs

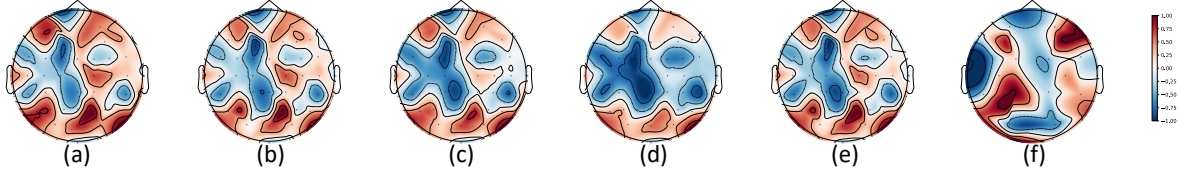


Figure 4: Ablation studies of EEG patterns on the EMBARC datasets for MDD severity classification task. (a) is the pattern of original HD EEG w/o masking. (b), (c) and (d) are patterns of HD EEG w/ 25%, 50% and 75% masking ratios and reconstructed by vanilla GMAE-PT, respectively. (e) is the pattern of HD EEG w/ 50% masking ratio and reconstructed by our proposed GCMAE-PT. (f) is the pattern of original MD EEG w/o masking. The MSE losses value for (b), (c), (d), and (e) are 0.25, 0.31, 0.44, and 0.17, respectively.

similarly to the large student model, and within the same model type, the large model substantially exceeds the tiny model in performance.

Table 2: The ablation explores the impact of varying EEG densities (HD/MD/LD), model types (teacher/student), and sizes (tiny/large) on performance. ‘T’ denotes teacher models and ‘S’ denotes student models. The experiments were conducted on the HBN dataset for the MDD classification task, with all models evaluated without pre-training. Metrics are reported as AUROC(%) / ACC(%).

Density	GFormer-Large (T)		Gformer-Tiny (S)	
	Tiny	Large	Tiny	Large
LD	72.6/76.4	76.3/80.4	72.7/75.4	74.2/78.7
MD	75.6/78.1	78.7/82.5	74.3/77.2	76.0/80.5
HD	77.7/80.4	80.4/83.6	75.6/78.4	77.1/81.7

5.4.2. ABLATION STUDY OF DIFFERENT DISTILLATION LOSS FUNCTIONS

We compared the proposed GTD loss with several commonly used graph distillation loss functions. As shown in the Table 3, our GTD loss outperforms other graph distillation losses. Furthermore, we observed that combining our GTD loss with the traditional logits distillation loss yields the best distillation performance, as it enables the model to distill both semantic information from the logits and structural information from the topology learned via GTD loss.

5.4.3. ANALYSIS OF PRE-TRAINING METHODS

As detailed in Table 4, we compared our GCMAE-PT with three other pre-training approaches: graph contrastive pre-training (GCL-PT) (You et al., 2020), graph masked autoencoder pre-training (GMAE-PT) (Hou et al., 2022), and a sequential combination of GCL-PT and GMAE-PT (Seq. Comb.). Following pre-training, we evaluated the models on downstream classification tasks. The results indicate that our framework surpasses GCL-PT, GMAE-PT, and their sequential combination. This underscores that sequentially

Table 3: Ablation studies on logits distill loss and our GTD loss. T and S denote teacher and student. The experiments are conducted on the HBN dataset for MDD classification.

GKD Methods	w/o Pre-Training		w/ Pre-Training	
	HD-T-Large	LD-S-Tiny	HD-T-Large	LD-S-Tiny
Baseline	80.4/83.6	72.7/75.4	85.2/86.5	77.6/79.4
+ Logits	-	73.6/77.3	-	79.6/80.7
+ Proposed	-	73.8/78.5	-	80.2/81.5
+ Union	-	75.0/79.2	-	81.6/82.4
+ LSP	-	73.1/76.7	-	78.8/80.7
+ G-CRD	-	73.4/77.8	-	79.3/80.3

combining contrastive and generative pre-training methods does not achieve optimal performance. Our approach, which seamlessly integrates these techniques into a cohesive framework with explicit and implicit mutual supervision, delivers superior results.

Table 4: Ablation studies were conducted on our GCMAE-PT, comparing it with GCL-PT, GMAE-PT, and their sequential combination (Seq. Comb.). ‘T’ and ‘S’ represent the teacher and student models, respectively. The experiments were performed on the HBN and EMBARC datasets. Baseline results can be found in Table 2. The teacher model uses HD inputs with a large-size configuration, while the student model uses LD inputs with a tiny-size configuration.

GPT Methods	HBN MDD		EMBARC Sex	
	HD-T-Large	LD-S-Tiny	HD-T-Large	LD-S-Tiny
Baseline	80.4/83.6	72.7/75.4	71.8/74.4	64.2/69.3
GCL-PT	82.2/85.1	74.3/75.9	72.9/75.7	66.4/71.5
GMAE-PT	82.8/85.3	75.1/76.6	72.6/76.3	67.1/71.7
Seq. Comb.	83.3/85.9	75.6/78.1	73.3/77.1	67.6/72.5
GCMAE-PT	85.2/86.5	77.6/79.4	74.7/78.4	68.8/73.3

6. Experiments for Model Robust Analysis

To evaluate the robustness of our model, we introduce perturbations to the EEG data by adding Gaussian noise to the raw signals and randomly dropping electrode channels. As

Table 5: Performance on HBN MDD before and after EEG perturbations.

Model	Before Pertur.	After Perturbation	
		Add Noise	Drop Electrodes
GCN	76.4%	72.8%	71.4%
GFormer	80.4%	74.6%	75.7%
GraphMAE	83.3%	78.5%	77.8%
Ours	86.9%	83.7%	84.0%

shown in the Table 5, all models exhibit performance degradation under perturbations. However, our model shows the smallest performance drop, indicating superior robustness and stronger resistance to adversarial perturbations. This advantage can be attributed to our carefully designed pre-training strategy, which incorporates robustness-aware components, such as the reconstruction of masked electrodes in GMAE-PT, thereby enhancing the model’s ability to handle corrupted inputs.

7. Analysis of (Pre-)Training and Distillation

As shown in Fig. 5, we visualized the optimization process of the loss curves, including contrastive loss, reconstruction loss, and GTD loss, during both the pre-training and fine-tuning stages. Fig. 5(a) shows that during pre-training, we jointly optimized the contrastive loss and reconstruction loss for both the teacher and student models. All four losses converge effectively during optimization. Notably, the contrastive loss for both the teacher and student models exhibit similar optimization trends, as do the reconstruction losses. Fig. 5(b) illustrates the impact of the proposed GCMAT-PT and GTD loss on downstream classification tasks. We present the optimization curves for the general Cross-Entropy (CE) loss, as well as the optimization curves after applying GCMAT-PT, GTD loss, and both combined. It is clear that the CE loss is better optimized with the application of GCMAT-PT and GTD loss. This confirms that both pre-training with GCMAT-PT and GTD loss enhance the performance of downstream classification tasks.

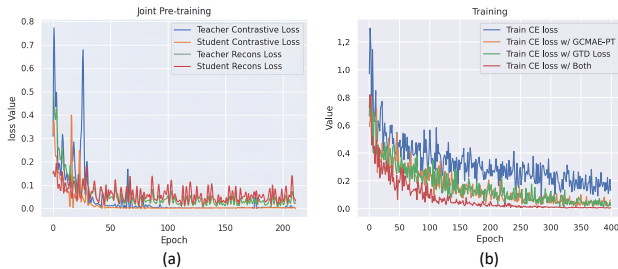


Figure 5: Illustrations of loss curves in both the pre-training stage (a) and fine-tuning stage (b). We applied early stopping to prevent overfitting. This also indicates that GTD loss effectively accelerates convergence and avoids overfitting.

8. Conclusion

In this paper, we present an innovative framework for EEG pre-training and distillation, which effectively integrates contrastive-based and generative-based graph pre-training paradigms. Furthermore, our framework incorporates a specifically designed EEG graph topology distillation loss function, tailored for the distillation process from high-density to low-density EEG data.

Impact Statement

Our proposed framework not only holds significant algorithmic innovation but also demonstrates strong practical application value. At the algorithmic level, our framework is the first to integrate the two mainstream pre-training paradigms within the realm of graph networks, while also combining graph pre-training with graph distillation. From an application perspective, we address a novel yet highly practical problem: how to distill high-density EEG into low-density EEG. The practical value of this lies in enabling our algorithm to achieve performance comparable to HD EEG data using more affordable and accessible LD EEG data, thereby reducing the challenges and costs associated with obtaining advanced equipment.

Acknowledgement

This work was supported in part by NIH grants (R01MH129694, R21AG080425, R21MH130956), NSF grants (2319451, 2215789), Alzheimer’s Association Grant (AARG-22-972541), and Lehigh University FIG (FIGAWD35) and CORE grants. Portions of this research were conducted on Lehigh University’s Research Computing infrastructure partially supported by NSF Award 2019035.

References

- Alexander, L. M., Escalera, J., Ai, L., Andreotti, C., Febre, K., Mangone, A., Vega-Potler, N., Langer, N., Alexander, A., Kovacs, M., et al. An open resource for transdiagnostic research in pediatric mental health and learning disorders. *Scientific data*, 4(1):1–26, 2017.
- Boessen, R., Groenwold, R. H., Knol, M. J., Grobbee, D. E., and Roes, K. C. Comparing hamd17 and hamd subscales on their ability to differentiate active treatment from placebo in randomized controlled trials. *Journal of affective disorders*, 145(3):363–369, 2013.
- Cataldo, A., Criscuolo, S., De Benedetto, E., Masciullo, A., Pesola, M., Schiavoni, R., and Invitto, S. A method for optimizing the artifact subspace reconstruction performance in low-density eeg. *IEEE Sensors Journal*, 22(21): 21257–21265, 2022.

- Davidson, R. J., Pizzagalli, D., Nitschke, J. B., and Putnam, K. Depression: perspectives from affective neuroscience. *Annual review of psychology*, 53(1):545–574, 2002.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Duan, R.-N., Zhu, J.-Y., and Lu, B.-L. Differential entropy feature for eeg-based emotion classification. In *2013 6th international IEEE/EMBS conference on neural engineering (NER)*, pp. 81–84. IEEE, 2013.
- Feng, Y., You, H., Zhang, Z., Ji, R., and Gao, Y. Hypergraph neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 3558–3565, 2019.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.
- Hinton, G., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Hou, Z., Liu, X., Cen, Y., Dong, Y., Yang, H., Wang, C., and Tang, J. Graphmae: Self-supervised masked graph autoencoders. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 594–604, 2022.
- Hou, Z., He, Y., Cen, Y., Liu, X., Dong, Y., Kharlamov, E., and Tang, J. Graphmae2: A decoding-enhanced masked self-supervised graph learner. In *Proceedings of the ACM web conference 2023*, pp. 737–746, 2023.
- Hu, Y., You, H., Wang, Z., Wang, Z., Zhou, E., and Gao, Y. Graph-mlp: Node classification without message passing in graph. *arXiv preprint arXiv:2106.04051*, 2021.
- Hu, Z., Dong, Y., Wang, K., Chang, K.-W., and Sun, Y. Gpt-gnn: Generative pre-training of graph neural networks. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 1857–1867, 2020.
- Huang, Z., Jin, X., Lu, C., Hou, Q., Cheng, M.-M., Fu, D., Shen, X., and Feng, J. Contrastive masked autoencoders are stronger vision learners. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- Jiang, W., Zhao, L., and Lu, B.-l. Large brain model for learning generic representations with tremendous eeg data in bci. In *The Twelfth International Conference on Learning Representations*.
- Joshi, C. K., Liu, F., Xun, X., Lin, J., and Foo, C. S. On representation knowledge distillation for graph neural networks. *IEEE transactions on neural networks and learning systems*, 2022.
- Justesen, A. B., Foged, M. T., Fabricius, M., Skaarup, C., Hamrouni, N., Martens, T., Paulson, O. B., Pinborg, L. H., and Beniczky, S. Diagnostic yield of high-density versus low-density eeg: The effect of spatial sampling, timing and duration of recording. *Clinical Neurophysiology*, 130(11):2060–2064, 2019.
- Kingma, D. P. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Knyazev, G. G. Motivation, emotion, and their inhibitory control mirrored in brain oscillations. *Neuroscience & Biobehavioral Reviews*, 31(3):377–395, 2007.
- Kuang, F., Shu, L., Hua, H., Wu, S., Zhang, L., Xu, X., Liu, Y., and Jiang, M. Cross-subject and cross-device wearable eeg emotion recognition using frontal eeg under virtual reality scenes. In *2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp. 3630–3637. IEEE, 2021.
- Lashgari, E., Liang, D., and Maoz, U. Data augmentation for deep-learning-based electroencephalography. *Journal of Neuroscience Methods*, 346:108885, 2020.
- Lawhern, V. J., Solon, A. J., Waytowich, N. R., Gordon, S. M., Hung, C. P., and Lance, B. J. Eegnet: a compact convolutional neural network for eeg-based brain-computer interfaces. *Journal of neural engineering*, 15(5):056013, 2018.
- Liu, F., Yang, P., Shu, Y., Liu, N., Sheng, J., Luo, J., Wang, X., and Liu, Y.-J. Emotion recognition from few-channel eeg signals by integrating deep feature aggregation and transfer learning. *IEEE Transactions on Affective Computing*, 2023.
- Pizzagalli, D. A. Frontocingulate dysfunction in depression: toward biomarkers of treatment response. *Neuropsychopharmacology*, 36(1):183–206, 2011.
- Price, J. L. and Drevets, W. C. Neural circuits underlying the pathophysiology of mood disorders. *Trends in cognitive sciences*, 16(1):61–71, 2012.
- Qiu, J., Chen, Q., Dong, Y., Zhang, J., Yang, H., Ding, M., Wang, K., and Tang, J. Gcc: Graph contrastive coding for graph neural network pre-training. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 1150–1160, 2020.

- Sanei, S. and Chambers, J. A. *EEG signal processing*. John Wiley & Sons, 2013.
- Siuly, S., Li, Y., and Zhang, Y. Eeg signal analysis and classification. *IEEE Trans Neural Syst Rehabil Eng*, 11: 141–144, 2016.
- Song, T., Zheng, W., Song, P., and Cui, Z. Eeg emotion recognition using dynamical graph convolutional neural networks. *IEEE Transactions on Affective Computing*, 11 (3):532–541, 2018.
- Song, Y., Zheng, Q., Liu, B., and Gao, X. Eeg conformer: Convolutional transformer for eeg decoding and visualization. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31:710–719, 2022.
- Stoyell, S. M., Wilmskoetter, J., Dobrota, M.-A., Chinapen, D. M., Bonilha, L., Mintz, M., Brinkmann, B. H., Herman, S. T., Peters, J. M., Vulliemoz, S., et al. High-density eeg in current clinical practice and opportunities for the future. *Journal of clinical neurophysiology*, 38(2): 112–123, 2021.
- Tang, S., Dunnmon, J. A., Saab, K., Zhang, X., Huang, Q., Dubost, F., Rubin, D. L., and Lee-Messer, C. Self-supervised graph neural networks for improved electroencephalographic seizure analysis. *arXiv preprint arXiv:2104.08336*, 2021.
- Thibodeau, R., Jorgensen, R. S., and Kim, S. Depression, anxiety, and resting frontal eeg asymmetry: a meta-analytic review. *Journal of abnormal psychology*, 115(4): 715, 2006.
- Trivedi, M. H., McGrath, P. J., Fava, M., Parsey, R. V., Kurian, B. T., Phillips, M. L., Oquendo, M. A., Bruder, G., Pizzagalli, D., Toups, M., et al. Establishing moderators and biosignatures of antidepressant response in clinical care (embarc): Rationale and design. *Journal of psychiatric research*, 78:11–23, 2016.
- Williams, J. B. A structured interview guide for the hamilton depression rating scale. *Archives of general psychiatry*, 45(8):742–747, 1988.
- Xie, Y., Xu, Z., Zhang, J., Wang, Z., and Ji, S. Self-supervised learning of graph neural networks: A unified review. *IEEE transactions on pattern analysis and machine intelligence*, 45(2):2412–2429, 2022.
- Yang, L., Tian, Y., Xu, M., Liu, Z., Hong, S., Qu, W., Zhang, W., Bin, C., Zhang, M., and Leskovec, J. Vqgraph: Rethinking graph representation space for bridging gnns and mlps. In *The Twelfth International Conference on Learning Representations*, 2024.
- Yang, Y., Qiu, J., Song, M., Tao, D., and Wang, X. Distilling knowledge from graph convolutional networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 7074–7083, 2020.
- You, Y., Chen, T., Sui, Y., Chen, T., Wang, Z., and Shen, Y. Graph contrastive learning with augmentations. *Advances in neural information processing systems*, 33:5812–5823, 2020.
- Yun, S., Jeong, M., Kim, R., Kang, J., and Kim, H. J. Graph transformer networks. *Advances in neural information processing systems*, 32, 2019.
- Zhang, C., Liu, J., Dang, K., and Zhang, W. Multi-scale distillation from multiple graph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 4337–4344, 2022a.
- Zhang, C., Zhang, C., Song, J., Yi, J. S. K., Zhang, K., and Kweon, I. S. A survey on masked autoencoder for self-supervised learning in vision and beyond. *arXiv preprint arXiv:2208.00173*, 2022b.
- Zhong, P., Wang, D., and Miao, C. Eeg-based emotion recognition using regularized graph neural networks. *IEEE Transactions on Affective Computing*, 13(3):1290–1301, 2020.
- Zhu, Y., Xu, Y., Yu, F., Liu, Q., Wu, S., and Wang, L. Deep graph contrastive representation learning. *arXiv preprint arXiv:2006.04131*, 2020.

A. Preliminaries of Dynamic GNNs

In traditional GNNs, the adjacency matrix \mathcal{A} is static. However, in this paper, we adopt dynamic GNNs, where the adjacency matrix can be dynamically adjusted during training to suit the specific task better. This approach allows the model to adapt the graph structure based on the input data and learning objectives. In such models, the edge weights α_{ij} between nodes (i, j) are learned during training. The edge weights can be computed as:

$$\alpha_{ij} = \sigma(f(\mathcal{X}_i, \mathcal{X}_j)) \quad (13)$$

where $f(\cdot)$ is a function for calculating edge weights, and σ is an activation function (e.g., Sigmoid). The dynamic adjacency matrix \mathcal{A} is then updated based on these weights, typically using a thresholding mechanism:

$$\mathcal{A}_{ij} = \begin{cases} 1, & \text{if } \alpha_{ij} > \theta \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

where θ is a threshold. During message passing, the dynamic adjacency matrix influences how messages are aggregated from neighboring nodes:

$$\mathbf{m}_i = \sum_{j \in \mathcal{N}(i)} \alpha_{ij} \mathbf{W} x_j \quad (15)$$

Here, α_{ij} represents the dynamically computed edge weight used to weight the messages from neighbors. Node features are updated as follows:

$$\mathcal{X}_i^{(l+1)} = \sigma(\mathbf{W}^{(l)} \mathcal{X}_i^{(l)} + \mathbf{b}^{(l)} + \mathbf{m}_i) \quad (16)$$

By dynamically adjusting the adjacency matrix, dynamic GNNs can capture more complex and evolving relationships within the graph, thereby enhancing flexibility and overall performance.

B. Details of Motivation and Problem

B.1. GTL for Unlabeled/Labeled EEG

Many existing methods primarily focus on training models with limited labeled EEG data, overlooking the potential of abundant unlabeled data. These methods emphasize novel GNN architectures but fail to fully leverage the available data. Additionally, they do not exploit high-density (HD) EEG data to improve models for low-density (LD) scenarios. This underscores the need for strategies that integrate both labeled and unlabeled data, and use HD data to enhance performance in LD contexts.

Moreover, most pre-training methods are directly applied to EEG time series, with very few addressing the issue from the perspective of large-scale graph pre-training. In contrast, our approach proposes pre-training EEG graph models using a graph-based pre-training perspective. This not only aims to transfer knowledge from unlabeled EEG data to tasks on labeled EEG data but also benefits HD-to-LD distillation. This is based on the following observation:

Observation: An LD EEG graph can be viewed as an HD EEG graph with specific nodes removed. In graph contrastive self-supervised pre-training, contrastive views are obtained by graph augmentation, such as removing nodes and edges. Another graph pre-training method, graph masked autoencoders pre-training, operates by masking node features and then reconstructing them. The relationships between these methods are formulated as follows:

$$\underbrace{\text{Density Decrease}}_{\text{Electrodes Loss}} \iff \underbrace{\text{Node Dropping}}_{\text{GCL Augmentation}} \iff \underbrace{\text{Node Masking}}_{\text{GMAE Masking}} \quad (17)$$

Based on this observation, we propose a novel unified graph self-supervised pre-training paradigm called GCMAE-PT. This approach intricately combines Graph Contrastive Pre-training (Qiu et al., 2020; You et al., 2020) with Graph Masked Autoencoders Pre-training (Hou et al., 2022), allowing us to model and capture the relationships among the three entities described in Eq. 18.

B.2. GKD for High/Low-Density EEG

As previously mentioned, an LD EEG graph can be viewed as an HD EEG graph with specific nodes removed. Consequently, HD EEG contains many features that LD EEG lacks. We naturally formulate this as a graph knowledge distillation (GKD)

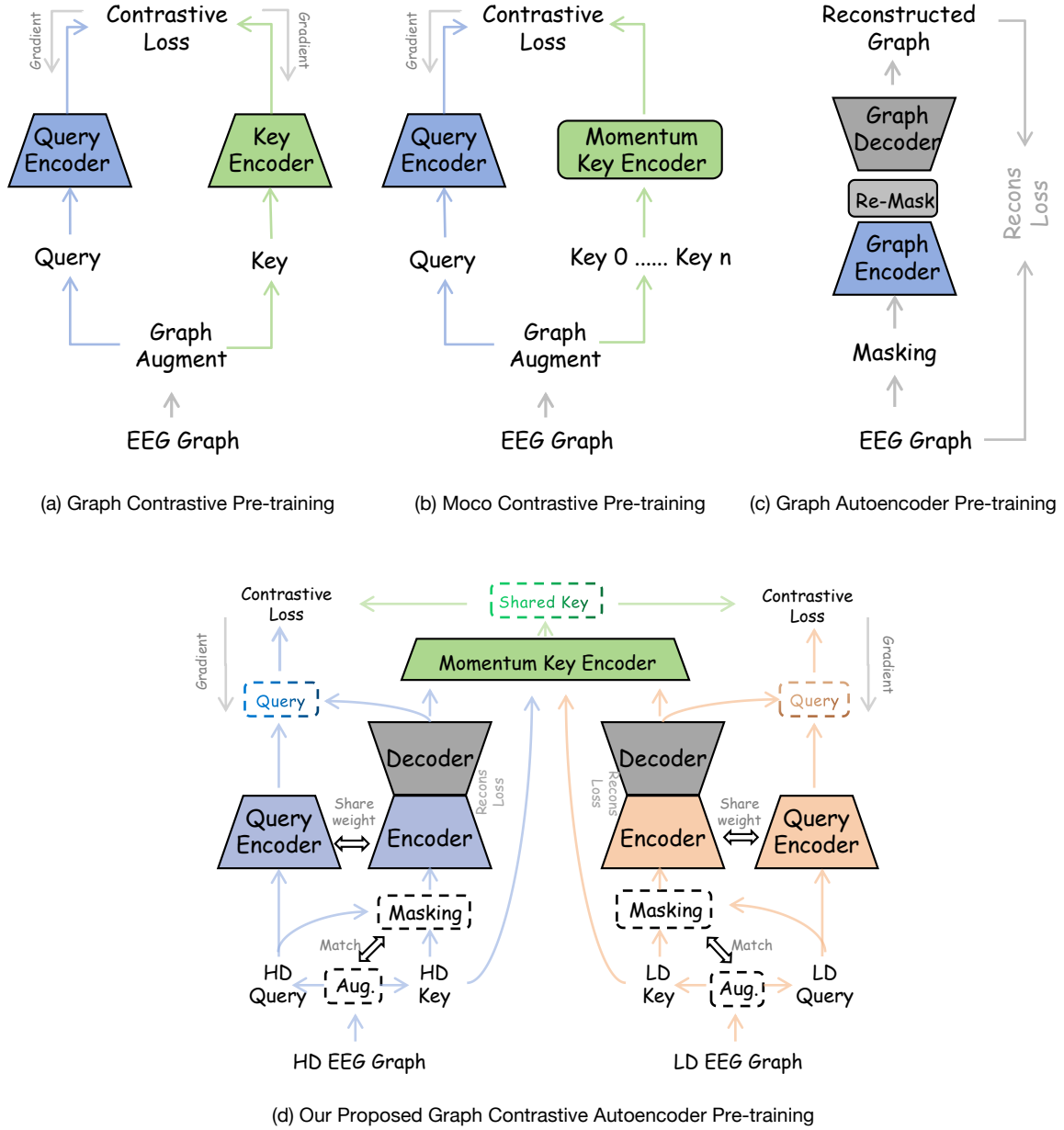


Figure 6: Illustration of our proposed pipeline and other previous contrastive pre-training and masked autoencoder pre-training. (a) Classical graph contrastive pre-training, which contrasts query and key samples. (b) Moco-like contrastive pre-training via a key momentum encoder. (c) Graph masked autoencoder pre-training through masking and then reconstructing graph samples. (d) Our proposed graph contrastive masked autoencoder pre-training framework. We simultaneously pre-train the teacher and student models by contrasting their reconstruction samples and reconstructing their contrastive samples (query and key).

task, focusing on how to transfer information from HD EEG data to LD EEG applications, which is a data-level distillation process. Additionally, if a more complex teacher model with a larger number of parameters is used to extract features from HD EEG data, and a simpler student model with fewer parameters is used for LD EEG data, this involves model-level distillation. The aim is to deploy the lightweight student model while ensuring that its performance approaches, or even surpasses, that of the more cumbersome teacher model.

Therefore, the GKD process can be represented by the following formula:

$$\underbrace{\text{Teacher Model}}_{\text{HD EEG Data}} \xrightarrow[\text{Distill (Data-level)}]{\text{Compress (Model-level)}} \underbrace{\text{Student Model}}_{\text{LD EEG Data}} \quad (18)$$

C. Illustrations of the Proposed Two Assumptions

As illustrated in Fig. 7, we have visually presented the overall pipeline of our model through an illustrative diagram. As shown, the entire pipeline is constructed based on the two assumptions we proposed in Section 3.2.

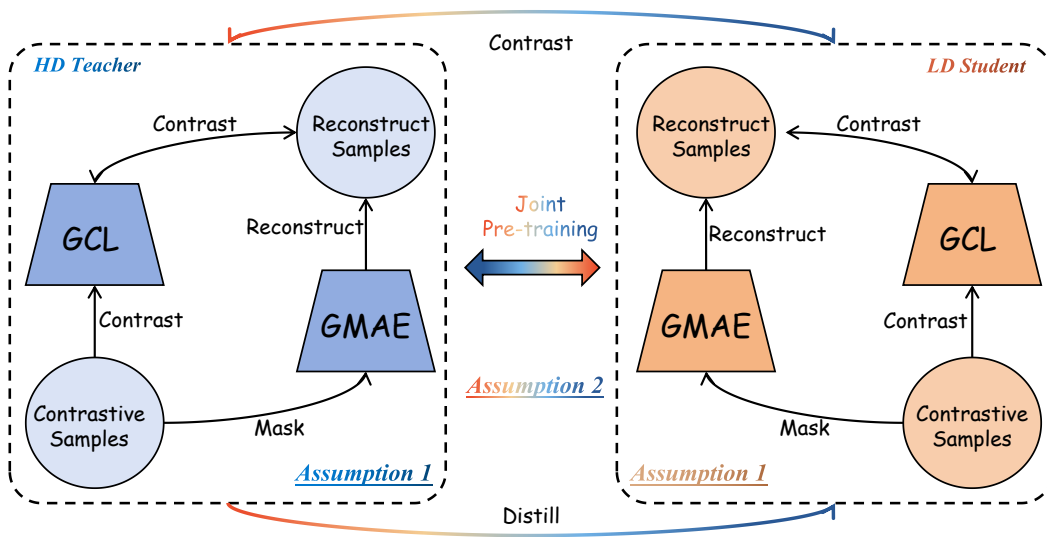


Figure 7: The schematic diagram of two assumption in Section 3.2. For assumption 1, we unify the GCL and GMAE pre-training paradigms by contrasting the reconstruction samples of GMAE in GCL and reconstructing the contrastive samples of GCL in GMAE. For assumption 2, we simultaneously pre-train the teacher and student models by contrasting their contrastive and reconstruction samples of GCL and GMAE, enable the encoders to be distillers via the help of GTD Loss.

D. Data Collection and Pre-processing

D.1. EEG Data Quantity Statistics

The EMBARC dataset consists of EEG signals collected from 308 subjects in both eye-open and eye-closed states. The EEG time series were sampled at 250Hz, with each trial lasting approximately 200 seconds. Similarly, in the HBN dataset, EEG signals were collected in both eye-open and eye-closed states, with 1,594 subjects for the eye-open condition and 1,764 subjects for the eye-closed condition. The duration of the recordings is also around 200 seconds, with the same sampling frequency of 250Hz. Both EMBARC and HBN datasets use the 10-20 EEG standard system, with EMBARC employing a 64-electrode cap and HBN using a 128-electrode cap.

D.2. Explanation of Unlabeled Data

Collecting EEG recordings, each patient diagnosed with a particular mental disorder can be classified as a labeled subject. Patients with EEG diagnosed as other disorders or healthy controls, are categorized as unlabeled data. In clinical, the amount of labeled data diagnosed as certain disorders was limited. Therefore, models trained exclusively on such sparse

labeled data are prone to underfitting, undermining their predictive performance. However, by broadening the scope to include aggregated data from a range of disorders to form a comprehensive unlabeled or mixed-labeled dataset, pre-training models on this enriched dataset can mitigate the constraints imposed by data scarcity. This approach enhances the model’s generalizability and improves performance, even in the face of limited labeled examples.

D.3. Construction and Augmentation of Pre-Training Graph Datasets

To construct the pre-training dataset, we combined the data from both the eye-open and eye-closed states from these two datasets. For EEG data augmentation, we applied a sliding window sampling method for each subject in the EMBARC and HBN datasets. EEG time series segments were extracted every 50 seconds, with a 20-second overlap between consecutive segments. The formula for calculating the number of segments for each subject is as follows:

$$\text{Segments per Subject} = \left\lfloor \frac{\text{Total length} - \text{Window length}}{\text{Window length} - \text{Overlap length}} \right\rfloor + 1 \quad (19)$$

Additionally, we combined the entire time series for each subject with the extracted segments. For each time series segment, we computed the Power Spectral Density features and then constructed the EEG graph samples. The formula for calculating the total number of samples used in the construction of the pre-training datasets is as follows:

$$\text{Total Samples} = (\text{Segments per Subject} + 1) \times \text{Subjects} \quad (20)$$

Note that the term *Subjects* here refers to the combination of EEG segments from both EMBARC and HBN datasets, including both eye-open and eye-closed EEG samples. Ultimately, we obtain approximately 4,000 samples ($308 + 308 + 1,594 + 1,764 \approx 4,000$), resulting in about 24,000 EEG graph samples for the graph self-supervised pre-training corpus.

D.4. (Pre-)Training and Evaluation Settings

For pre-training on the EMBARC dataset, we addressed the issue of dataset size disparity between EMBARC and the HBN dataset, which both originate from the same EEG system. Specifically, we downsampled the 128 channels of the HBN data to 64, 32, and 16 channels, maintaining the same arrangement. These downsampled data were then combined with the corresponding density datasets from EMBARC to create a unified pre-training dataset. Note that, as the EMBARC dataset does not include 128 channels, the 128-channel HD pre-training dataset does not incorporate data from EMBARC (HBN only).

For downstream task fine-tuning, due to the limited amount of labeled data, we employed 10-fold cross-validation with 10 runs for all model training. The Adam optimizer (Kingma, 2014) was used to optimize the training process. Pre-training was performed over 200 epochs, while downstream fine-tuning was carried out for 400 epochs.

D.5. Construction of Downstream Datasets

Table 6 provides the quantity of labeled data for four downstream classification tasks across the EMBARC and HBN datasets.

In the EMBARC dataset, the number of subjects is consistent across eye-open and eye-closed conditions. For the MDD sex classification task, there are 296 subjects with varying levels of depression (all diagnosed with depression) and 12 normal subjects. Among the depressed individuals, there are 194 males and 102 females. For the depression severity classification task, 166 subjects are diagnosed with severe depression (HAMD_{17} score > 17) (Boessen et al., 2013), and 130 subjects are diagnosed with mild depression (HAMD_{17} score ≤ 17).

The HBN dataset, which includes a range of diseases, has significantly fewer labeled samples compared to the total data volume due to the high number of samples without explicit MDD and ASD diagnostic labels. Additionally, the number of labeled subjects differs between eye-open and eye-closed conditions. In the eye-open data, there are 178 healthy controls, 109 MDD patients, and 234 ASD patients. In the eye-closed data, there are 187 healthy controls, 120 MDD patients, and 245 ASD patients.

To ensure a large-scale pre-training dataset, we utilized slicing operations to expand the dataset size. However, for constructing labeled datasets for downstream tasks, slicing was not employed. Instead, we calculated the PSD features for the entire 200-second EEG time series.

D.6. Comparison between the Pre-training Dataset and Downstream Datasets

For the pre-training dataset, which includes both labeled and unlabeled data, we applied slicing operations to significantly increase the dataset size. In contrast, for the downstream dataset, particularly for the HBN data, the labeled data constitutes only a small fraction of the total dataset, and no slicing operations were performed. In this context, it is crucial to leverage the pre-training dataset effectively to enhance model performance on the limited labeled data available.

Table 6: Labeled data distribution of the EMBARC and HBN datasets. 'HC' means healthy control.

Datasets	EMBARC		HBN	
	Sex	Severity	MDD	ASD
Eye-open	Female: 194 Male: 102	Severe: 166 Mild: 130	Patient: 109 HC: 178	Patient: 234 HC: 178
Eye-closed	Female: 194 Male: 102	Severe: 166 Mild: 130	Patient: 120 HC: 187	Patient: 245 HC: 187

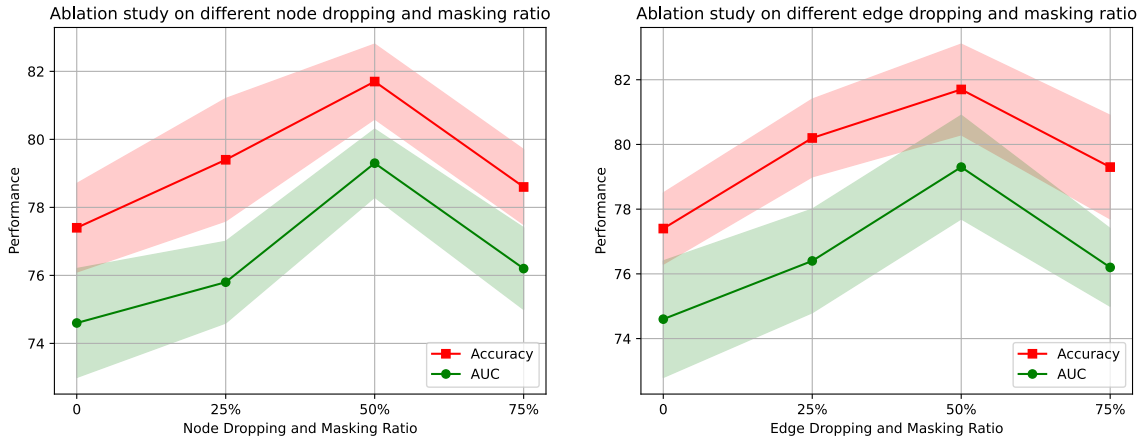


Figure 8: Ablation studies on different node and edge dropping (for GCL-PT) and masking (for GMAE-PT) ratios. A 50% masking ratio for both nodes and edges achieves the best performance.

E. Different Configurations of Teacher and Student Models

Table 7: Comparison of Model Configurations. Note that DGCNN is a **spectral-based** vanilla GCNs model (DGCNN), while GFormer means the **spatial-based** Graph Transformer model. We considered both types of graph models to demonstrate the versatility of our pipeline.

Model	Encoder	Sizes (S)	Layers (L)	Dimensions (D)	Heads (H)	Position Embedding (P)	Params (PM)
Teacher	DGCNN	Large	8	128	-	✗	5.7M
	GFormer	Large	8	128	8	✓	6.9M
Student	DGCNN	Tiny	4	64	-	✗	1.3M
	GFormer	Tiny	4	64	4	✓	1.4M

DGCNN and Graph Transformer are two representative types of graph neural networks (GNNs). DGCNN is a message-passing-based GNN, while the Graph Transformer is a spatial-attention-based GNN. In terms of performance, GNNs represented by DGCNN are relatively lightweight but tend to achieve lower accuracy. In contrast, GNNs represented by Graph Transformers generally yield better performance, albeit with higher model complexity and computational cost.

Table 8: Subject-dependent and subject-independent results on EMBARC (PSD+Pearson) and SEED (DE+Coherence) using LD EEG data.

Model	Sex Classification (EMBARC)		Emotion Recognition (SEED)	
	Subject-Dependent	Subject-Independent	Subject-Dependent	Subject-Independent
Graph Transformer	71.6%	68.2%	86.4%	75.4%
GraphMAE	73.8%	70.6%	88.6%	78.1%
Ours	76.8%	74.1%	93.6%	84.3%

F. Subject-Dependent/Independent Experiments

Table 8 reports the performance of different models on two representative datasets, EMBARC and SEED (Duan et al., 2013), under both subject-dependent and subject-independent settings. For the EMBARC dataset, we constructed EEG graphs using Power Spectral Density (PSD) as node features and Pearson correlation as the connectivity metric. In contrast, the SEED dataset employed Differential Entropy (DE) (Duan et al., 2013) features and coherence-based functional connectivity to construct the EEG graph. Notably, our proposed model achieves the best performance across all settings, significantly outperforming baselines.

More importantly, despite being pre-trained or optimized primarily on medical data, our model generalizes remarkably well to the SEED dataset for the task of emotion recognition. This demonstrates the strong transferability of our model across different EEG datasets, feature types, and functional connectivity metrics, highlighting its robustness and versatility in real-world cross-dataset scenarios.

G. Ablation Study on Different Connectivity (Spatial/Functional) for EEG Graph Construction

We construct EEG graphs using two primary approaches: functional connectivity-based methods and spatial distance-based methods. The functional connectivity-based methods include Pearson correlation, coherence, and mutual information. As shown in the Table 9, the coherence-based graph construction yields the best performance, while the spatial distance-based method performs the worst. This is likely because the strength of functional connectivity between EEG electrodes does not necessarily correlate with their physical distance, electrodes that are far apart spatially may still exhibit strong connections due to similar functional activity.

Table 9: Comparison of EEG graph construction methods using different connectivity measures.

Model	Functional Connectivity			Spatial Connectivity
	Pearson Correlation	Coherence	Mutual Information	Spatial Distance
Graph Transformer	71.6%	73.1%	71.9%	67.5%
GraphMAE	73.8%	74.7%	73.8%	68.5%
Ours	76.8%	78.1%	77.5%	72.8%

H. Ablation study on different EEG bands

As shown in Table 10, we performed ablation studies on various EEG frequency bands across four tasks on two datasets and observed that the alpha band consistently yielded the best performance across all tasks. Consequently, we selected the alpha band as our primary configuration.

Table 10: Ablation experiments on performance of different EEG bands. The model employed is a **tiny-sized student** model obtained through pre-training and distillation with **LD EEG input**.

Datasets	Downstream Tasks	Alpha	Beta	Gamma	Theta	All Bands
HBN	MDD Diagnosis	84.8/85.4	82.6/84.1	81.3/82.5	79.6/80.6	86.3/87.6
	ASD Diagnosis	66.1/66.4	61.4/64.1	63.3/63.5	60.7/62.6	68.7/70.2
EMBARC	MDD Sex	73.4/76.7	71.6/72.4	70.3/70.9	68.5/72.9	74.6/79.0
	MDD Severity	68.6/71.9	66.3/67.2	64.5/66.8	63.7/66.2	70.1/74.5

I. Similarity Kernel Selection for GTD loss

We follow (Joshi et al., 2022) and try different similarity kernels to measure the distance between nodes. All four kernels are shown in the following formula:

$$\begin{aligned}
 \text{Linear Kernel: } \mathcal{Z}_{ij} &= \mathcal{F}(\mathcal{X}_i, \mathcal{X}_j) = \mathcal{X}_i \cdot \mathcal{X}_j \\
 \text{Euclidean Kernel: } \mathcal{Z}_{ij} &= \mathcal{F}(\mathcal{X}_i, \mathcal{X}_j) = \|\mathcal{X}_i - \mathcal{X}_j\|_2 \\
 \text{Polynomial Kernel: } \mathcal{Z}_{ij} &= \mathcal{F}(\mathcal{X}_i, \mathcal{X}_j) = (\mathcal{X}_i \cdot \mathcal{X}_j + c)^d \\
 \text{RBF Kernel: } \mathcal{Z}_{ij} &= \mathcal{F}(\mathcal{X}_i, \mathcal{X}_j) = \exp(-\gamma \|\mathcal{X}_i - \mathcal{X}_j\|_2^2)
 \end{aligned} \tag{21}$$

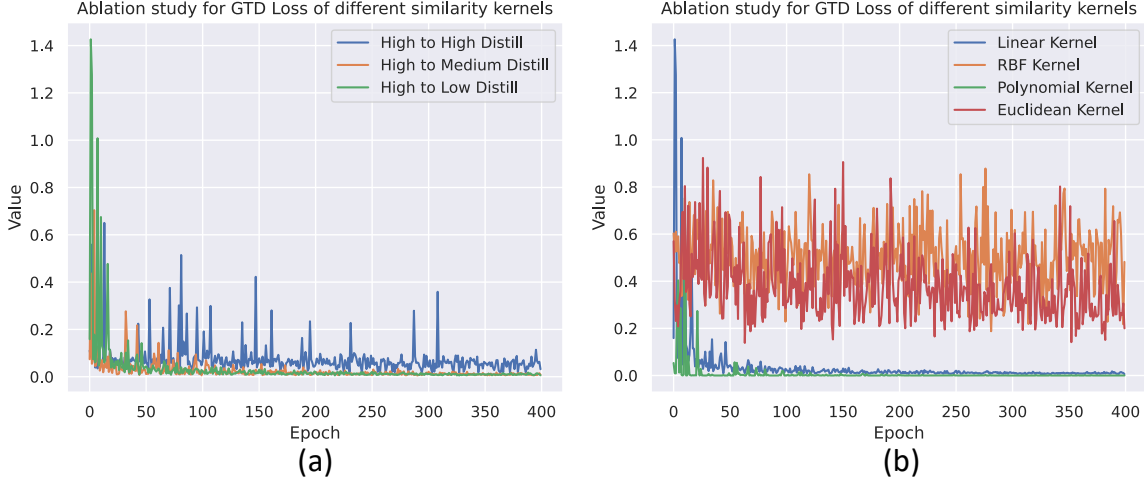


Figure 9: The ablation studies of distillation across different density settings (a) and kernels (b).

As shown in Fig. 9, we conducted ablation experiments on the GTD loss. Figure 9(a) illustrates the results of distillation in three scenarios: high-to-low (H2L), high-to-medium (H2M), and high-to-high (H2H). Note that H2H is a special case. Although the GTD loss is designed primarily for high-to-low density distillation, it can also be applied to high-to-high density distillation as an exception.

The optimization curves for H2L and H2M show good convergence. However, in the special case of H2H, while the optimization curve also converges, the gradient descent is less pronounced. This suggests that although GTD loss can still be applied in the H2H scenario, it is less effective. This is because GTD loss mainly focuses on nodes that are removed, and since no nodes are removed in H2H, the distillation’s primary goal is to correct the student model’s misinterpretation of connectivity. Consequently, there is less knowledge to distill compared to H2L and H2M scenarios, resulting in a less noticeable decrease in the optimization curve. In contrast, the optimization curve for H2L shows the most significant decrease, followed by H2M.

Figure 9(b) presents the results of ablation experiments with different similarity kernels. The experimental results indicate that only the GTD loss using polynomial and linear kernels achieved good convergence during optimization. Among these, the linear kernel provided the best distillation effect, which is why we selected it as the primary kernel for our GTD loss.

J. Clinical Interpretation of EEG Patterns

By grounding our visual findings in these established studies, we provide a clearer link between the reconstructed EEG patterns and their clinical implications, emphasizing the robustness and diagnostic utility of our approach.

To address the clinical relevance of our EEG pattern reconstructions, we link the visual patterns presented in Figure 3 to established clinical findings in MDD research. The unmasked EEG pattern in Figure 4(a) reveals clear activations in the frontal and central regions, which are crucial areas involved in cognitive processing and emotional regulation. These regions are frequently highlighted in MDD studies due to their role in mood and executive function. Specifically, the prefrontal

cortex, anterior cingulate cortex, and related regions are implicated in emotional processing and regulation, with MDD patients often showing disrupted activity in these areas (Davidson et al., 2002; Price & Drevets, 2012). Reduced activation in these areas can reflect difficulties in cognitive control and emotional regulation, key features of depressive symptoms.

As the masking ratio increases (Figures 4b-d), the patterns show a noticeable decline in activation and connectivity, particularly in the frontal and central regions. This aligns with findings in MDD literature, where disrupted functional connectivity, especially in the frontocingulate networks, is a well-documented feature of the disorder. For example, alterations in prefrontal connectivity are often associated with the severity of depressive symptoms and the inability to regulate negative emotions (Pizzagalli, 2011). The degradation observed in Figures 4b-d is consistent with the hypothesis that higher masking ratios simulate information loss, highlighting the importance of intact frontal connectivity for accurate MDD classification. Critically, Figure 4(e), which displays the reconstructed pattern using our proposed GCMAE-PT with 50% masking, closely resembles the unmasked pattern seen in Figure 4(a). The reconstructed data retain key activations in the frontal and parietal regions, indicating that our method effectively preserves clinically relevant EEG features even under challenging conditions. This preservation is crucial because altered activity in these regions, particularly in the alpha and theta bands, is often linked to cognitive and emotional dysregulation in MDD patients (Thibodeau et al., 2006; Knyazev, 2007). For instance, lower alpha activity in the frontal regions has been associated with greater emotional dysregulation, while changes in theta activity are linked to altered cognitive processes, both of which are core characteristics of MDD.

The preserved patterns in Figure 4(e) suggest that GCMAE-PT can maintain these clinically significant EEG characteristics, which are essential for accurate classification of MDD severity. This finding not only demonstrates the robustness of our reconstruction method but also aligns with known clinical markers of MDD, supporting the practical relevance of our approach. Furthermore, the ability to accurately reconstruct these key patterns contributes directly to classification tasks, as regions showing consistent and clinically significant alterations are critical for distinguishing between different severity levels of MDD. By maintaining the integrity of essential EEG features under masking, our method ensures that the reconstructed data remain informative and diagnostically valuable, potentially leading to better predictive performance in clinical settings.

K. Experiments on Very-Low Density Situation

To further test the generalization ability of our model in extreme scenarios, we evaluated it under a very low-density (VLD) condition. Specifically, we tested the extreme case with EEG data using only 8 electrodes. As shown in Table 11, our proposed pre-training framework and the corresponding GTD loss are able to tackle the extreme case with very few electrodes.

Table 11: Experiments on the very low-density (VLD) situation. HD -> LD/VLD means high-density to (very)low-density distillation.

PT Methods	FT Loss	HD -> LD		HD -> VLD	
		Sex	Severity	Sex	Severity
GCL-PT	w/o GTD	1.8%↑	1.7%↑	1.5%↑	2.0%↑
GMAE-PT	w/o GTD	1.6%↑	1.5%↑	2.1%↑	2.4%↑
GCMAE-PT (Ours)	w/o GTD	2.9%↑	3.8%↑	3.5%↑	4.4%↑

L. Experiments of Different Fine-tuning Paradigms

Table 12: Experiments on the effectiveness and efficiency of different fine-tuning (FT) methods. The experiments are conducted on HBN dataset for MDD classification. The unit of fine-tuning speed is seconds (s), and the unit of memory cost is gigabytes (G). The input data consists of 128-channel HD EEG graphs, and the model uses a large-size graph transformer.

Fine-tuning Methods	Effectiveness		Efficiency	
	AUROC	ACC	FT Speed	Memory Cost
Vanilla FT	80.4%	83.6%	183s	1.0G
Parameter-Efficient FT	78.7%	83.4%	86s	0.3G

As shown in Table 12, we evaluated two distinct fine-tuning paradigms. The first, termed Vanilla FT, involves fine-tuning all parameters of the pre-trained encoder. The second, referred to as parameter-efficient FT, entails freezing the lower layers of the pre-trained encoder and fine-tuning only the parameters of the upper layers, such as the fully connected layers. It is evident that parameter-efficient FT, which requires fewer parameters to be optimized, results in a fine-tuning speed three times faster and memory usage one-third that of Vanilla FT. However, this approach incurs a slight performance trade-off compared to Vanilla FT.

M. Analysis of Shared Key Pool Queue

The reason for implementing the proposed teacher-student shared key pool queue is that we have two types of original input data: HD and LD EEG graphs. Through the key pool queue, we allow high-density and low-density EEG key samples to share the same gradient update process within a batch. This approach also enables both the teacher and student models to simultaneously capture shared patterns between these two types of data.

N. Experiments on Held-Out Validation

In the pre-training dataset of the previous experiment, as shown in Table 13, we integrated heterogeneous EEG data from different diseases to pretrain our model. To further validate the reliability of our model, we conducted a held-out validation experiment.

Table 13: Ablation studies on the proposed pre-training framework and the GTD loss. The held-out validation means we pre-train the model only on the HBN dataset and fine-tune the model to the EMBARC dataset. HD -> LD means high-density to low-density distillation. 'All' means HBN + EMBARC. ↑ means performance improvement in terms of accuracy. The downstream task is conducted on the EMBARC dataset for MDD severity classification task. The backbone model is the spatial-based Graph Transformer.

PT Methods	FT Loss	Datasets		Distill Performance	
		Pre-Train	Fine-Tune	HD -> MD	HD -> LD
GCL-PT	w/o GTD	HBN	EMBARC	1.5%↑	1.7%↑
GMAE-PT	w/o GTD	HBN	EMBARC	1.7%↑	1.5%↑
Seq. Comb.	w/o GTD	HBN	EMBARC	2.0%↑	2.1%↑
GCMAT-PT (Held-Out)	w/o GTD	HBN	EMBARC	3.1%↑	3.0%↑
GCMAT-PT (Ours)	w/o GTD	All	EMBARC	3.7%↑	3.8%↑
GCL-PT	w/ GTD	HBN	EMBARC	2.3%↑	2.2%↑
GMAE-PT	w/ GTD	HBN	EMBARC	2.2%↑	2.5%↑
Seq. Comb.	w/ GTD	HBN	EMBARC	2.7%↑	3.1%↑
GCMAT-PT (Held-Out)	w/ GTD	HBN	EMBARC	4.4%↑	5.0%↑
GCMAT-PT (Ours)	w/ GTD	All	EMBARC	4.7%↑	5.6%↑

O. Challenges, Limitations, and Future Works

EEG graph self-supervised pre-training offers a promising avenue for leveraging extensive EEG data, paving the way for large-scale graph-based EEG models. Our proposed GCMAT-PT method is well-suited as a pre-training approach for large-scale EEG foundation model. However, a key challenge is unifying data with varying electrode configurations across different EEG systems to address data heterogeneity. In our study, while constructing a unified EEG pre-training dataset from multiple sources, we faced the constraint of all datasets being from the same EEG system (10-20 system). To standardize the data, we reduced the number of electrodes in datasets with more electrodes to match those with fewer electrodes, creating a unified pre-training dataset. This approach, however, leads to a loss of information from removed electrodes and restricts the use of datasets with fewer electrodes for pre-training on datasets with more electrodes. Addressing the challenge of integrating EEG data with differing electrode counts from various systems, while preserving electrode precision, is crucial for developing a comprehensive pre-training dataset. Successfully overcoming this issue could enable large-scale graph pre-training and establish a robust EEG graph foundation model, representing a significant advancement in the field.