# GMSS: Graph-Based Multi-Task Self-Supervised Learning for EEG Emotion Recognition

Yang Li, *Member, IEEE,* Ji Chen, Fu Li*, *Member, IEEE,* Boxun Fu, Hao Wu, Youshuo Ji, Yijin Zhou, Yi Niu, Guangming Shi, *Fellow, IEEE,* Wenming Zheng, *Senior Member, IEEE*

**Abstract**—Previous electroencephalogram (EEG) emotion recognition relies on single-task learning, which may lead to overfitting and learned emotion features lacking generalization. In this paper, a graph-based multi-task self-supervised learning model (GMSS) for EEG emotion recognition is proposed. GMSS has the ability to learn more general representations by integrating multiple self-supervised tasks, including spatial and frequency jigsaw puzzle tasks, and contrastive learning tasks. By learning from multiple tasks simultaneously, GMSS can find a representation that captures all of the tasks thereby decreasing the chance of overfitting on the original task, i.e., emotion recognition task. In particular, the spatial jigsaw puzzle task aims to capture the intrinsic spatial relationships of different brain regions. Considering the importance of frequency information in EEG emotional signals, the goal of the frequency jigsaw puzzle task is to explore the crucial frequency bands for EEG emotion recognition. To further regularize the learned features and encourage the network to learn inherent representations, contrastive learning task is adopted in this work by mapping the transformed data into a common feature space. The performance of the proposed GMSS is compared with several popular unsupervised and supervised methods. Experiments on SEED, SEED-IV, and MPED datasets show that the proposed model has remarkable advantages in learning more discriminative and general features for EEG emotional signals.

**Index Terms**—EEG emotion recognition, multi-task learning, self-supervised learning, graph neural network.

---

## 1 INTRODUCTION

EMOTION is close to everyone and plays an important role in our daily lives [1]. It is a complex and comprehensive psychological and physiological state that can be characterized by behavioral and physiological signals [2]. Neuroscience research indicates that physiological signals are closer to the source of emotion than behavioral signals [3]. As a physiological signal, EEG has the advantage of being difficult to disguise and hide compared with behavioral signals, such as facial expressions and voice [4]. Moreover, EEG signals significantly benefited from the technological developments in non-invasive EEG recording methods, and are widely used in the research on emotion recognition [5] [6]. In recent years, emotion recognition has become a research hotspot in human-computer interaction and affective computing [7].

A wide variety of methods has been proposed to effectively analyze EEG emotional signals over the past decades. Traditional machine learning methods typically adopt a two-stage model to implement emotional recognition. For example, Lin et al. [8] extracted power spectrum density, differential asymmetry power, and rational asymmetry power

as features of EEG signals, and then classified them using a support vector machine to study the relationship between emotion and EEG signals. Jenke et al. [9] studied and compared the effects of EEG emotion features extracted from the time domain, the frequency domain, and the time-frequency domain on EEG emotion signal recognition. However, traditional machine learning methods rely on handcrafted features and expert experience [10]. With the spectacular success of deep learning methods in the field of computer vision and language recognition, many researchers have considered deep learning models for EEG emotion signals for their ability to automatically extract complex features [11] [12]. For instance, some researchers utilized convolutional neural networks (CNNs) and recurrent neural networks (RNNs) to handle emotion recognition [13] [14]. Recently, the topological structure of EEG signals has been increasingly studied in EEG emotion recognition owing to the superior performance of graph neural networks (GNNs) in irregular data structures [15]. Wang et al. [16] proposed a multichannel EEG emotion recognition method based on phase locking value (PLV) graph convolutional neural networks (P-GCNN) to extract the spatio-temporal characteristics and the inherent information in functional connections.

Based on the literature, most EEG-based emotion recognition methods basically face three challenges: (1) how to generalize the emotion recognition model, and correctly classify new data; (2) how to make full use of EEG characteristics to capture more discriminative data representation for emotion recognition; and (3) how to solve the problem of emotional noise labels. Regarding the first challenge, EEG displays a highly heterogeneous and nonstationary

- *Yang Li, Ji Chen, Fu Li, Boxun Fu, Hao Wu, Youshuo Ji, Yijin Zhou, Yi Niu, Guangming Shi are with the Key Laboratory of Intelligent Perception and Image Understanding of Ministry of Education, the School of Artificial Intelligence, Xidian University, Xi'an, 710071, China. (E-mail: fuli@mail.xidian.edu.cn).*
- *Wenming Zheng is with the Key Laboratory of Child Development and Learning Science (Ministry of Education), School of Biological Sciences and Medical Engineering, Southeast University, Nanjing, Jiangsu, 210096, China.*
*\*Corresponding author*

behavior because emotional signals usually consist of many neural process [17]. The enormous data distribution shift leads to a lack of generalization for data from different subject or new situation of the current subject [18]. Thus, some researchers have adopted the domain adaptation (DA) method to improve generalization. For example, Li et al. [18] proposed a bi-hemisphere domain adversarial neural network (BiDANN) that contains three domain discriminators to assist with the learning of discriminative emotional features, narrowing the distribution gap between training and testing data, and improving the generality of the recognition model. However, most DA-based methods achieve generality by training the model on labeled training data and unlabeled testing data, which is not suitable for real applications. Thus, it is meaningful and applicable to explore other methods that learn general data representations without test data. For the second challenge, handcrafted features such as power spectrum density, statistical measure, and discrete wavelet transform are frequently used for generic EEG signal classification tasks. However, these features are not specially designed for EEG emotion signal [18]. This issue has also been discussed in recent deep-learning literature on EEG emotion recognition. For example, Zheng et al. [19] employed a deep belief network (DBN) to directly model the EEG emotion signal. Even though these handcrafted and deep features have been able to extract certain emotion discriminative information, they do not sufficiently exploit specific emotion-related information in EEG emotion recognition tasks. Thus, it is necessary to utilize the characteristics of EEG signals to extract high-level features. Regarding the third challenge, the emotion labels in the collected EEG data may be noisy and inconsistent as participants may not always produce the expected emotions when watching emotions stimulate stimuli [20]. Consequently, it is challenging and meaningful to explore how to solve the problem of emotional noise labels that are often ignored in EEG emotion recognition research.

To address the above three major issues in EEG emotion recognition tasks, in this paper we propose GMSS, which can learn general EEG emotion representation and improve EEG emotion recognition ability by solving three self-supervised pretext tasks. To improve generality, GMSS adopts multiple EEG emotion-related tasks that share learned knowledge to generate more general features and avoid overfitting [21]. GMSS consists of two graph-based jigsaw puzzle tasks and a contrastive learning task, making it capable to study the impact of emotional expression on spatial and frequency information. The spatial jigsaw puzzle task enables the predefined distant electrodes to become neighbor electrodes and more emotion-related spatial information is learned in return. Meanwhile, the frequency jigsaw puzzle task explores crucial frequency bands for EEG emotion recognition. Utilizing the augmented samples of the above jigsaw puzzle tasks, the contrastive learning task further standardizes the feature space and enhances the generalization ability of the model. These self-supervised pretext tasks, which are based on the intrinsic attributes of EEG emotion data, allow GMSS to deal with EEG noise labels without semantic labeling. In this study, both unsupervised and supervised approaches of GMSS were evaluated. The experimental results show that GMSS achieves state-of-the-

art (SOTA) performance on three public datasets.

In summary, the contributions of this work can be outlined as follows:

- To the best of our knowledge, this is the first work that adopts multi-task learning to improve model generalization capability and avoid overfitting in EEG emotion recognition.
- Through the pretext tasks of jigsaw puzzles and contrastive learning, GMSS learns more discriminative features and alleviates the problem of emotional noise labels, which further improves EEG emotion recognition.
- The experimental results, based on both unsupervised and supervised learning approaches, demonstrate that GMSS can achieve SOTA performance on three benchmark datasets.

The rest of this paper is organized as follows: Section II provides an overview of previous studies on EEG emotion recognition, graph neural networks, multi-task learning, and self-supervised learning. Section III specifies the GMSS method and its application to EEG emotion recognition. In section IV the proposed method is evaluated for EEG emotion recognition through extensive experiments. Finally, section V concludes the paper.

## 2 RELATED WORKS

In this section, related works on EEG-based emotion recognition, graph neural networks, multi-task learning, and self-supervised learning are introduced.

### 2.1 EEG-based Emotion Recognition

The general process of EEG emotion recognition includes feature extraction and classification. Traditional machine learning-based methods typically adopt the statistical measure, discrete wavelet transform, or power spectrum density [8] as features and then classify the extracted features using SVM, LDA, or LR [22]. However, deep learning based methods generally extract features by designing feature extraction neural networks followed by linear layers to achieve classification. Many deep learning methods such as CNN, RNN and GNN have been introduced to effectively distinguish different emotional states in EEG emotional signals. Li et al. [23] proposed a hierarchical spatial-temporal neural network (R2G-STNN) based on a bidirectional long short-term memory (BiLSTM) network to capture the intrinsic spatial relationships of EEG electrodes within the brain region and between brain regions for EEG emotion recognition. Song et al. [24] proposed a multichannel EEG emotion recognition method based on a novel dynamic graph convolutional neural network (DGCNN) to dynamically learn the intrinsic relationship between different EEG channels to assist with features classification. Zhong et al. [20] proposed a regularized graph neural network (RGNN) with two regularizers to deal with cross-subject EEG variations and the noise label problem, and achieved promising results. Li et al. [25] proposed a bi-hemispheric discrepancy model (BiHDM) to learn discrepancy information between two hemispheres to improve EEG emotion recognition ability.

## 2.2 Graph Neural Network

The traditional convolutional neural network is excellent for dealing with Euclidean data. However, GNN is suitable for handling non-Euclidean data and has shown great promise in the field of social networks, recommendation systems, and knowledge maps [26] [27] [28]. The GNN fall into two categories, spectral-based and spatial-based. The spectral-based method specifies graphic convolution by introducing a filter from the perspective of graphic signal processing, where the graphic convolution operation is viewed as noise removal from the graphic signal. The spatial-based method is based on the recurrent neural network theory and defines the graph convolution through information propagation [29]. Defferrard et al. [30] argued that the original spectrum convolution suffers from the disadvantages of a large number of parameters and high complexity, and proposed a fast localized convolution algorithm using a recursive formulation of the K-order Chebyshev polynomials to approximate the filters. Kipf et al. [15] proposed a graph convolutional network (GCN) with a faster localized graph convolutional operation, which is the first-order approximation of Chebyshev polynomials, that is, $K = 1$. Veli et al. [31] proposed a graph attention network (GAT), which stacking layers in nodes that are able to attend over their neighborhoods' features, specifying different weights to different nodes in a neighborhood, without requiring costly matrix operation or depending on knowing the graph structure upfront.

Bianchi et al. [32] proposed a graph convolutional layer that provides a flexible frequency response, which is more robust to noise, and better captures the global graph structure. Bouritsas et al. [33] proposed a graph substructure network that is more expressive than Weisfeiler-Leman graph isomorphism test, which allows the model retains multiple attractive properties of standard GNNs, while being able to eliminate even hard instances of graph isomorphism. Ciano et al. [34] proposed a mixed inductive–transductive GNN model, study its properties and introduce an experimental strategy that help to understand and distinguish the role of inductive and transductive learning. Tiezzi et al. [35] proposed an approach to learning in GNNs based on constrained optimization in the Lagrangian framework. Learning both the transition function and the node states is the outcome of a joint process, in which the state convergence procedure is implicitly expressed by a constraint satisfaction mechanism, avoiding iterative epoch-wise procedures and the network unfolding.

However, in EEG emotion recognition, some GNN-based methods [24] only consider second-order or third-order neighbors to avoid over-smoothing, which may result in the loss of valuable information between distant nodes. Thus, the spatial jigsaw puzzle was applied to challenge this problem.

## 2.3 Multi-Task Learning

Multi-task learning is an effective machine learning method and has shown its advantages in many fields, including computer vision [36] [37], natural language processing [38], and speech recognition [39]. Ruder et al. [21] introduced two commonly used multi-task learning methods in deep learning, clarifying the working principle of multi-task learning as well as pointing out that properly designed pretext tasks can encourage the model to learn a more general representation while decreasing the risk of overfitting. Compared with a single task, multi-task learning combines multiple related tasks and utilizes all the data from each task so that the knowledge on each task is shared. Additional information on the associated tasks is also obtained in multi-task learning models, resulting in significant improvements in the learning ability, generalization capability, and robustness of the model [40]. However, considering the different significance of each task, the weight of each task should be dynamic. Sener et al. [41] regarded multi-task learning as a multi-objective optimization problem and proved that optimizing the upper bound of the multi-objective loss can obtain the Pareto optimal solution. Kendall et al. [37] proposed a principled approach to multi-task deep learning that weighs multiple loss functions by considering the homoscedastic uncertainty of each task to avoid the cost of manual tuning. Benefiting from these advantages, in this work, multi-task learning framework is adopted to learn more generalization features and reduce the risk of overfitting.

## 2.4 Self-Supervised Learning

Self-supervised learning is a popular method for learning intrinsic information using unlabeled data [42]. Generally, self-supervised learning applies the attributes of data to generate pseudo labels as opposed to human-annotated labels to train the network. Based on the different data attributes used in the design, there are four categories of pretext tasks: generation-based, context-based, free semantic label-based, and cross-modal-based [42]. In the visual feature learning field, context-based pretext tasks mainly employs spatial structure, temporal structure, and context similarity for the design. Many studies learn the general features of images by predicting the relative position of the patches to solve jigsaw puzzle tasks, thereby solving the problem of image classification [43] [44] [45]. Gidaris et al. [46] applied a 2D rotation to the image to construct the pretext task and then predicted the rotation angle to enable the model to learn the position, type, and posture of objects in the image. Carno et al. [47] used a clustering method to generate pseudo labels for images and combined learning neural network parameters and result features to obtain more abundant semantic information. Mathilde et al. [48] proposed a method for unsupervised learning of visual features by contrasting cluster assignments (SwAV), which takes advantage of contrastive methods without requiring to compute pairwise comparisons. He et al. [49] suggested that momentum contrast (MoCo) would significantly narrow the gap between unsupervised representation learning and supervised representation learning. The performance of the contrastive SimCLR framework proposed by Chen et al. [50] on ImageNet surpasses that of supervised learning based models. Xinlei et al. [51] proposed a simple Siamese (SimSiam) network that achieved the best results without negative samples, large batches, and momentum encoders. In addition, the contrastive learning method was applied in the field of video processing, and achieved excellent performance at the time it was proposed [52] [53]. In the
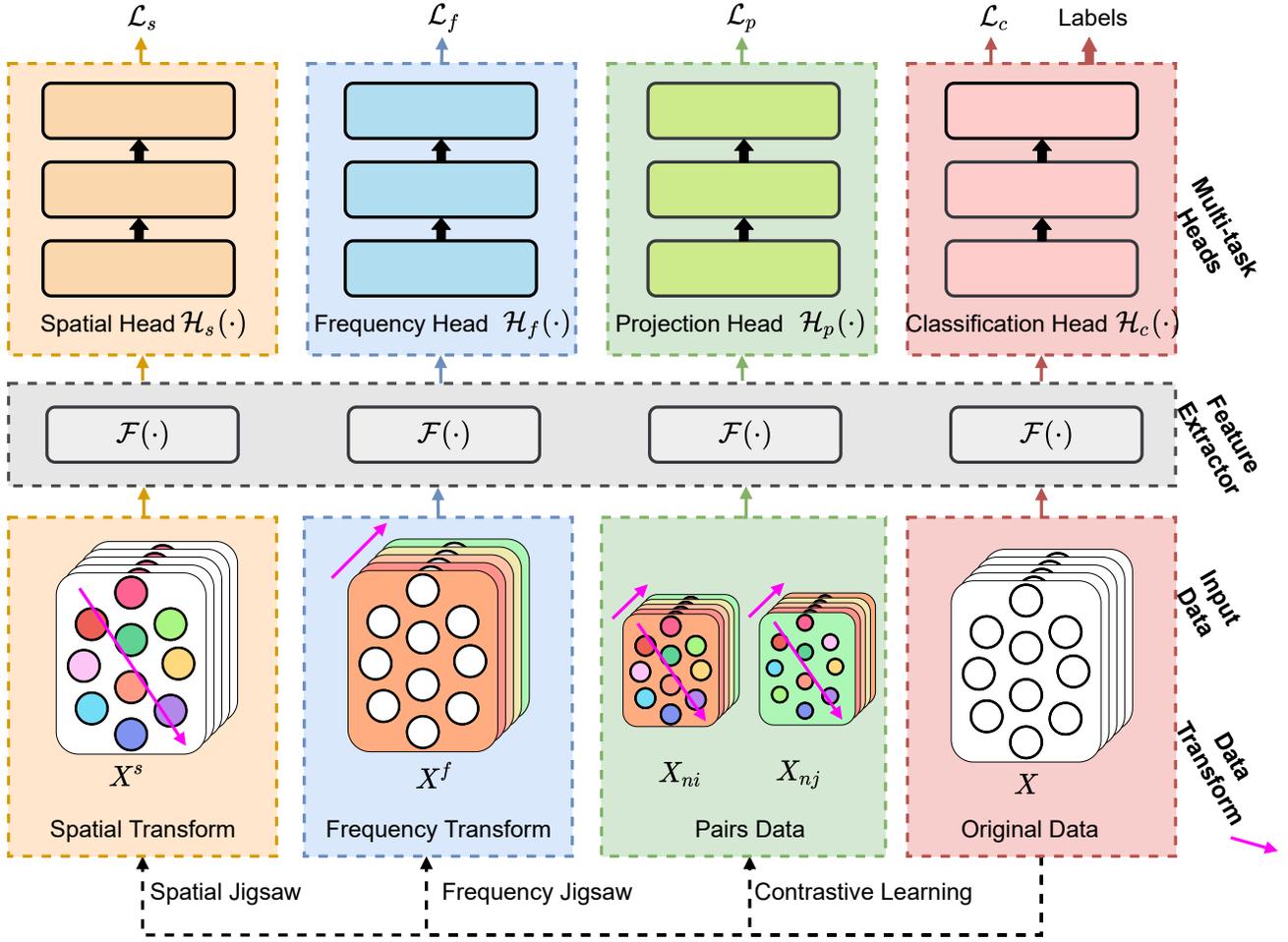
Fig. 1. Framework of GMSS. In the unsupervised training mode, for the upstream task, the original graph data are not used to train the network. For the downstream task, the feature extractor is frozen and only the pink part with a substituted one linear layer is executed. In the supervised training mode, all the parts are executed simultaneously.

field of EEG emotion recognition, Xie et al. [54] proposed an innovative solution which contains six different transformations to learn high-level EEG representation (SSL-EEG). Mohsenvand et al. [55] present a framework for learning representations from EEG signals via contrastive learning which recombines channels from multi-channel recordings and trains a channel-wised feature extractor to learn EEG emotion representation (SeqCLR). Inspired by self-supervised learning, in this work, two jigsaw puzzle tasks and a contrastive learning task were designed to assist with the learning of general EEG emotional features while circumventing the problem of EEG emotion noise labels. Further, DeepCluster is a method based on clustering, and SwAV use a swapped prediction mechanism to predict the cluster assignment of a view from the representation of another view, SSL-EEG learn the EEG representations from complex signal transformation, while MoCo, SimCLR, SimSiam and SeqCLR are methods based on maximizing the similarity between positive pairs. Compared with these methods above, our GMSS is a multi-task framework that incorporates multiple emotion-related tasks that utilizes all the data from each task so that the knowledge on each task is shared. This will helpful to obtain the additional information on the associated tasks that results in improving the learning ability, generalization capability, and robustness of the model. Another difference is that our self-supervised model concentrates on the characteristics of EEG emotion signal. For example, the jigsaw puzzle learning will force our model focus on the important brain regions and frequency bands of EEG signal, which are very important for emotion expression.

TABLE 1
EEG electrodes associated with each brain region in the experiment.

| Brain region | Electrode name |
|---|---|
| Pre-Frontal | AF3, FP1, FPZ, FP2, AF4 |
| Frontal | F1, FZ, F2, FC1, FCZ, FC2 |
| Left Frontal | F7, F5, F3, FT7, FC5, FC3 |
| Right Frontal | F4, F6, F8, FC4, FC6, FT8 |
| Left Temporal | T7, C5, C3, TP7, CP5, CP3 |
| Right Temporal | C4, C6, T8, CP4, CP6, TP8 |
| Central | C1, CZ, C2, CP1, CPZ, CP2, P1, PZ, P2 |
| Left Parietal | P7, P5, P3, PO7, PO5, CB1 |
| Right parietal | P4, P6, P8, PO6, PO8, CB2 |
| Occipital | PO3, POZ, PO4, O1, OZ, O2 |

# 3 GRAPH-BASED MULTI-TASK SELF-SUPERVISED LEARNING FOR EEG EMOTION RECOGNITION

The goal of GMSS is to capture general and discriminative EEG emotion features using multi-task self-supervised learning, as illustrate in Fig. 1. Three self-supervised tasks are designed to achieve this goal under unsupervised and supervised modes. These tasks share a common feature extractor. There are four task heads, i.e., Spatial Head $\mathcal{H}_s(\cdot)$, Frequency Head $\mathcal{H}_f(\cdot)$, Projection Head $\mathcal{H}_p(\cdot)$, Classification Head $\mathcal{H}_c(\cdot)$. $\mathcal{H}_s$ and $\mathcal{H}_f$ are employed for spatial puzzle and frequency puzzle respectively. $\mathcal{H}_p$ is adopted to project the learned representation into feature space. $\mathcal{H}_c$ is used for emotion recognition. Each head consists of three fully connected layers.
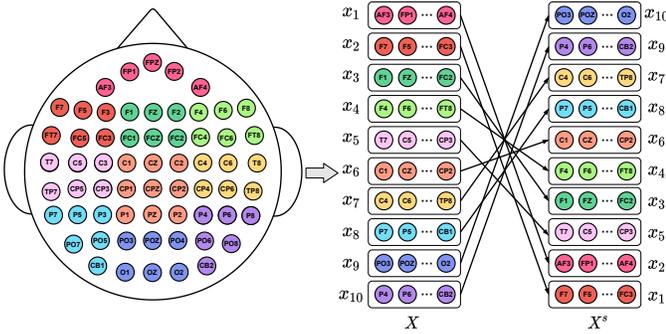


Fig. 2. Spatial jigsaw puzzle. The 62 electrodes are divided into 10 blocks according to the location of brain regions. The placement of these channels is relocated while keeping the original connection based on the topology of the scalp. The spatial jigsaw puzzle task is to identify which of the 128 classes the channels reorganized by blocks belong.

## 3.1 Multiple Self-Supervised Tasks

To learn more generalized and discriminative features and alleviate the noise problem of EEG emotion labels, multiple self-supervised learning tasks are considered, including the spatial jigsaw puzzle task, the frequency jigsaw puzzle task, and the contrastive learning task. Each of these three pretext tasks is described in depth.

### 3.1.1 Spatial Jigsaw Puzzle

The spatial jigsaw puzzle aims to capture the spatial patterns of EEG electrodes in different brain regions. Due to the different effects of brain regions on emotion expression, the spatial jigsaw puzzle task is defined as a series of brain region permutations [56] [57] [58] [59]. As shown in Table 1, the original EEG data $X \in \mathbb{R}^{n \times d}$ are partitioned into 10 blocks according to the location of the brain regions, denoted as $X = (\widetilde{X}_1, \widetilde{X}_2, \cdots, \widetilde{X}_{10})^{\mathsf{T}}$ where $\widetilde{X}_i \in \mathbb{R}^{n_i \times d}$, $\sum_{i=1}^{10} n_i = n, n_i > 0$. Then, all brain region permutations can be obtained:

$$\begin{cases} \hat{X}_1 & = (\widetilde{X}_1, \widetilde{X}_2, \cdots, \widetilde{X}_{10} | y_1), \\ \hat{X}_2 & = (\widetilde{X}_1, \widetilde{X}_2, \cdots, \widetilde{X}_9 | y_2), \\ & \vdots \\ \hat{X}_{10!} & = (\widetilde{X}_{10}, \widetilde{X}_9, \cdots, \widetilde{X}_1 | y_{10!}), \end{cases} \quad (1)$$

where $\hat{X}_i$ and $y_i$ represent the i-th permutation and its serial number, respectively. There are $10! = 3628800$ permutations in total. The goal is to distinguish which permutation the spatial transformed data corresponds to. However, it is quite challenging to distinguish these massive permutations for self-supervised pretext tasks. Therefore, we develop a $R_k(\cdot)$ operator. $R_k(\cdot)$ selects the $k$ permutations with maximum Hamming distance from the full permutation of Eq. (1) and randomly transformed the input data to one of the $k$ permutations. We define a unique pseudo label for each of these $k$ permutations, generating $k$ different kinds of pseudo labels in total, with a range from 1 to $k$. Each input data is randomly transformed into one of the $k$ permutations and the corresponding unique pseudo labels are obtained. $k$ is set to 128. The overall permutation is displayed in Fig. 2, and is formulated as follows:

$$(X^s, y^s) = R_{128}(X), \quad (2)$$

where $X^s$ is the generated EEG data with pseudo label $y^s \in \mathbb{Z}_+^{128}$.

To recognize these spatial jigsaw puzzles, a classification head $\mathcal{H}_s(\cdot)$ is applied, and cross entropy is adopted as the loss function. Formally, the loss of spatial jigsaw puzzle tasks can be expressed as $\mathcal{L}_s$:

$$\mathcal{L}_s = -\sum_{i=1}^{N} \bar{y}_i^s log(\mathcal{H}_s(\mathcal{F}(X_i^s))), \quad (3)$$

where $\mathcal{F}(\cdot)$ is the shared feature extractor, $\bar{y}_i^s$ is the one-hot encoding of the corresponding pseudo label $y_i^s$, and N is the number of training samples.



Fig. 3. Frequency jigsaw puzzle. The frequency jigsaw puzzle transforms the frequency bands of each channel of an EEG emotion data in the same way. The goal of the frequency jigsaw puzzle is to figure out which of the 120 classes the scrambled EEG emotion data belong.

### 3.1.2 Frequency Jigsaw Puzzle

The frequency jigsaw puzzle task is designed to learn the inner relationship between frequency bands, explore the crucial frequency bands for EEG emotion recognition and improve the discrimination ability of the model. In general, as illustrated in Fig. 3, the energy features of the EEG data are extracted from five emotion expression-related frequency bands, including $\delta$ (1-3 Hz), $\theta$ (4-7 Hz), $\alpha$ (8-13 Hz), $\beta$ (14-30 Hz), $\gamma$ (31-50 Hz). Similar to the spatial jigsaw puzzle, the original EEG data $X$ are divided into five blocks according to different frequency bands, denoted as $(x_1, x_2, \cdots, x_5)$, where $x_j \in \mathbb{R}^{n \times 1}$. The goal is to identify the corresponding permutation of the frequency

transformed data. All frequency bands permutations can be obtained:

$$
\begin{cases}
X_1' &= (x_1, x_2, \cdots, x_5 | y_1'), \\
X_2' &= (x_1, x_2, \cdots, x_4 | y_2'), \\
&\vdots \\
X_{5!}' &= (x_5, x_4, \cdots, x_1 | y_{5!}'),
\end{cases}
\tag{4}
$$

where $X_j'$ and $y_j'$ represent the j-th permutation and its serial number, respectively. In the frequency jigsaw puzzle, the operator $R_k(\cdot)$ is applied to generate transformed data with pseudo label, and $k = 120$:

$$
(X^f, y^f) = R_{120}(X),
\tag{5}
$$

where $X^f$ is the generated EEG data with pseudo label $y^f \in \mathbb{Z}_+^{120}$.

To recognize these frequency jigsaw puzzles, a classification head $\mathcal{H}_f(\cdot)$ is applied and cross entropy is adopted as the loss function. Formally, the loss in the frequency jigsaw puzzle task can be expressed as follows:

$$
\mathcal{L}_f = -\sum_{j=1}^{N} \bar{y}_j^f log(\mathcal{H}_f(\mathcal{F}(X_j^f))),
\tag{6}
$$

where $\mathcal{F}(\cdot)$ is the shared feature extractor and $\bar{y}_j^f$ is the one-hot encoding of the corresponding pseudo label $y_j^f$.
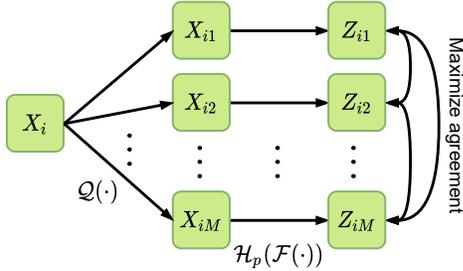


Fig. 4. Contrastive learning. The original data applied with spatial transformation and frequency transformation to generate the pairs data.

### 3.1.3 Contrastive Learning

To further regularize feature learning and encourage the network to learn inherent representations, contrastive learning is adopted to map the transformed data into a common feature space. The purpose is to maximize the agreement between the different augmented data of the same EEG emotion data, as shown in Fig. 4. To ensure that positive pairs move closer and negative pairs move far away in feature space, a data augmentation operation $\mathcal{Q}(\cdot)$ is defined to consider the spatial and frequency transformations of the same original EEG emotion data. For each original EEG emotion data $X_i, i \in \{1, 2, \cdots, N\}$, M augmented data $\{X_{i1}, X_{i2}, \cdots, X_{iM}\} = \mathcal{Q}(X_i)$ are obtained. As a result, each augmented data has $(M-1)$ positive pairs and $(N-1) \times M$ negative pairs. In total, $N \times M$ augmented data are obtained by:

$$
\begin{aligned}
\{X_{nm}; n \in \{1, 2, \cdots, N\}, m \in \{1, 2, \cdots, M\}\} = \\
\mathcal{Q}(X_1) \cup \mathcal{Q}(X_2) \cup \cdots \cup \mathcal{Q}(X_N),
\end{aligned}
\tag{7}
$$

where $X_{nm} \in \mathbb{R}^{n \times d}$ is the m-th transformation of the n-th EEG sample.

Similar to SimCLR [50], a projection head $\mathcal{H}_p(\cdot)$ is applied to map the EEG emotion data onto the feature space, that is, $Z_{nm} = \mathcal{H}_p(\mathcal{F}(X_{nm}))$. The similarity of two data points is quantitatively described by the dot product, which normalizes u and v through the $\ell_2$-norm. i.e., $sim(u, v) = u^\mathsf{T} v / \|u\| \|v\|$. Then, the loss of all positive pairs $\ell_n$ of sample $X_n$ is calculated as follows:

$$
\ell_n = -log \frac{g_+}{g_+ + g_-},
\tag{8}
$$

$$
g_+ = \sum_{i=1}^{M-1} \sum_{j=i+1}^{M} exp(sim(Z_{ni}, Z_{nj})/\tau),
\tag{9}
$$

$$
g_- = \sum_{o=1}^{M} \sum_{t=1}^{N} \sum_{w=1}^{M} exp(sim(Z_{no}, Z_{tw})/\tau), t \neq n,
\tag{10}
$$

where $(Z_{ni}, Z_{nj})$ are positive pairs, and $(Z_{no}, Z_{tw})$ are negative pairs. $\tau$ is the temperature parameter and is set to 0.5. Furthermore, the arithmetic average of the loss of all positive pairs' $\ell_n$ of all samples is calculated for backpropagation as follows:

$$
\mathcal{L}_p = \frac{1}{N} \sum_{n=1}^{N} \ell_n,
\tag{11}
$$

## 3.2 Training Mode for EEG Emotion Recognition

Two modes of training are provided: unsupervised and supervised. The feature extractor $\mathcal{F}(\cdot)$ in both modes are the same. When training the feature extractor, the distinction is in the presence or absence of ground-truth emotion labels. In the unsupervised mode, instead of using ground-truth emotion labels, the feature extractor $\mathcal{F}(\cdot)$ is trained only on the self-supervised tasks mentioned above. Then, the frozen feature extractor $\mathcal{F}(\cdot)$ is transferred to the downstream task and the performance is verified using a linear classifier. In the supervised mode, a joint training strategy is adopted. The network is simultaneously trained on self-supervised tasks and supervised tasks. To avoid manually tuning the weights of the different loss functions, the total loss function is defined by considering the homoscedastic uncertainty of each task [37]. In particular, the training loss $\mathcal{L}$ is calculated as follows:

$$
\begin{aligned}
\mathcal{L} = \frac{1}{\sigma_{\mathcal{L}_s}^2}\mathcal{L}_s + \frac{1}{\sigma_{\mathcal{L}_f}^2}\mathcal{L}_f + \frac{1}{2\sigma_{\mathcal{L}_p}^2}\mathcal{L}_p + log(\sigma_{\mathcal{L}_s}\sigma_{\mathcal{L}_f}\sigma_{\mathcal{L}_p}) \\
+ \psi \cdot (\frac{1}{\sigma_{\mathcal{L}_c}^2}\mathcal{L}_c + log(\sigma_{\mathcal{L}_c})),
\end{aligned}
\tag{12}
$$

$$
\psi = \begin{cases}
0, & unsupervised \quad mode, \\
1, & supervised \quad mode,
\end{cases}
\tag{13}
$$

where $\mathcal{L}_c$ is the cross entropy loss of supervised EEG emotion classification task; $\sigma_{\mathcal{L}_s}, \sigma_{\mathcal{L}_f}, \sigma_{\mathcal{L}_p}$ and $\sigma_{\mathcal{L}_c}$ are the observation noise scalars of the corresponding tasks [37]. $\psi$ is the mode-selection operator. The observation noise scalar $\sigma$ is a principled approach to multi-task deep learning which weighs multiple loss functions by the homoscedastic uncertainty of each task. This allows us to simultaneously learn

various quantities with different units or scales in both classification and regression settings, which can balance these weightings optimally, resulting in superior performance. These scalars can be calculated as learnable parameters which change constantly during the model training process and the initial values are 1.
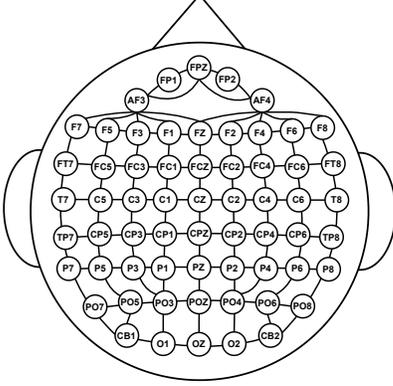


Fig. 5. EEG graph structure and adjacency matrix $A$ construction.

### 3.3 Feature Extractor of GMSS

As shown in Fig. 5, an undirected graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ is employed to model the EEG data. Meanwhile, the adjacency matrix $A$ of the EEG data was obtained. In $\mathcal{G}(\mathcal{V}, \mathcal{E})$, $\mathcal{V}$ denotes the set of nodes where $|\mathcal{V}| = n$; each node has $d$ dimensions. As a result, nodes can be represented by the feature matrix $X \in \mathbb{R}^{n \times d}$. $\mathcal{E}$ denotes a set of edges between nodes. $(v_i, v_j)$ is the edge between nodes $v_i$ and node $v_j$, that is, $(v_i, v_j) \in \mathcal{E}$. The adjacency matrix $A \in \mathbb{R}^{n \times n}$ contains the topological information of the undirected graph, that is, the EEG data. $D$ is the degree matrix of the vertices, and $L = D - A$ is the combinatorial Laplacian matrix. In this study, $n$ denotes the channel number of EEG data; $d$ denotes the number of frequency bands and $d = 5$. The energy feature is extracted from five bands, namely, $\delta$ (1-3 Hz), $\theta$ (4-7 Hz), $\alpha$ (8-13 Hz), $\beta$ (14-30 Hz), $\gamma$ (31-50 Hz).

In the GMSS model, Chebyshev polynomials are employed instead of the convolution kernel of SCNN [60] in the spectral domain, so that there are only k parameters in the convolution kernel, and feature decomposition is not required, reducing the computational load. Thus, the feature extractor $\mathcal{F}(\cdot)$ of the GMSS can be formulated as:

$$\mathcal{F}(X) = \sigma\left(\sum_{k=0}^{K-1} \beta_k T_k(\widetilde{L}) X\right), \quad (14)$$

where $\sigma(\cdot)$ is the activation function; $X$ is the input EEG emotion data; $\beta_k$ refers to the learning parameters in network training; and $T_k(\cdot)$ is the Chebyshev polynomial of order $K$. Additionally, $\widetilde{L} = 2L/\lambda_{max} - I$, where $\lambda_{max}$ is the maximum eigenvalue of Laplace matrix $L$. In this study, we set $K = 2$ to avoid over-smoothing.

## 4 EXPERIMENTS

In this section, experiments were conduct on the following three datasets to evaluate the performance of our model: SEED [19], SEED-IV [64], and MPED [63]. All three datasets

were collected while subjects watched emotional video clips in a quiet, comfortable, and non-interfering environment. All three datasets were generated by recording EEG signals through the ESI NeuroScan system using 62 electrode channels positioned according to the 10-20 system [56]. These three datasets are introduced next along with the experimental results.

### 4.1 Experimental Dataset

**SEED**. In the SEED dataset, there are a total of 15 subjects. There are three sessions associated with each subject. In each session, there are a total of 15 film clips to induce happy, neutral, and sad emotions, and there are 5 film clips for each emotion. That is, there are 15 trials per session, and each trial has 185-238 samples, resulting in approximately 3400 samples per session.

**SEED-IV**. In the SEED-IV dataset, similar to SEED, there are 15 subjects, and three sessions for each subject. The difference is that each session includes four kinds of emotions: happy, neutral, sad, and fear. Each emotion has 6 different film clips. As a result, there are 24 trials, and each trial has 12-64 samples for each session. Consequently, each session has approximately 830 samples.

**MPED**. In the MPED dataset there are 30 subjects, and each subject has only one session. In a session, there are seven types of emotions: joy, funny, neutral, sad, fear, disgust, and anger. Each type of emotion has 4 related film clips. Therefore, there are 28 trials per session. Each trial consists of 120 samples and there are a total of 3360 samples in one session.

### 4.2 Experimental Protocol

To fully evaluate our model, two types of experiments are implemented: subject-dependent and subject-independent experiment. For the subject-dependent experiment, the training data and testing data are obtained from different EEG trials of the same subject. For the subject-independent experiment, the training data and testing data are obtained from different subjects.

For the subject-dependent experiment, the same experimental protocol is applied as in [9] [19] [25] [63]. That is, for the SEED dataset, the EEG data of the first nine trials are used in each session as training data and the remaining six trials in the session as testing data for each subject. For the SEED-IV dataset, the first sixteen trials of the session are used for each subject as training data and the remaining eight trials as testing data. For the MPED dataset, the EEG data of the first twenty-one trials in the session are adopted for the training data and the remaining seven trials in this session are the testing data for each subject.

For the subject-independent experiment, the leave-one-subject-out (LOSO) cross-validation strategy is used in [25] [65] for each subject. Namely, one subject's EEG emotion data constituted the testing data, and the remaining subjects' EEG emotion data constituted the training data. The process continued until all subjects' EEG emotion data are tested once.

TABLE 2
Subject-dependent and subject-independent classification accuracy (mean/std) for unsupervised mode on SEED, SEED-IV, and MPED datasets

| Model | SEED | | SEED-IV | | MPED | |
| --- | --- | --- | --- | --- | --- | --- |
| | dependent | independent | dependent | independent | dependent | independent |
| DeepCluster [61] | 74.60/12.17* | 59.01/17.65* | 49.60/10.28* | 44.54/09.88* | 26.38/05.59* | 23.25/04.86* |
| MoCo [49] | 76.58/10.72* | 58.26/15.05* | 49.40/10.99* | 46.19/10.04* | 27.47/05.27* | 23.86/04.66* |
| SwAV [48] | 77.81/10.15* | 58.65/16.66* | 52.03/14.71* | 49.28/10.44* | 27.91/05.05* | 23.50/04.81* |
| SimCLR [50] | 81.79/11.15* | 63.45/15.96* | 52.47/11.57* | 50.07/11.17* | 29.53/05.36* | 24.21/05.10* |
| SimSiam [51] | 80.18/10.53* | 63.95/11.95* | 53.71/11.98* | 51.24/12.47* | 28.19/05.88* | 24.31/04.61* |
| SSL-EEG [54] | 83.32/09.20* | 67.52/12.73* | 63.59/19.82* | 53.62/08.47* | 25.22/04.25* | 21.87/02.53* |
| SeqCLR [55] | 82.91/08.97* | 64.56/11.89* | 63.13/15.41* | 50.75/07.71* | 30.47/06.07* | 23.33/03.89* |
| GMSS | **89.18/09.74** | **76.04/11.91** | **65.61/17.33** | **62.13/08.33** | **34.81/06.88** | **26.97/05.01** |

* indicates the experiment results obtained by our own implementation.
Note: For the subject-dependent experiment, we calculate the average accuracy based on the results of all the sessions. While for the subject-independent experiment, we calculate the average accuracy based on the results of all the subjects.

TABLE 3
Subject-dependent and subject-independent classification accuracy (mean/std) for supervised mode on SEED, SEED-IV, and MPED datasets

| Model | SEED | | SEED-IV | | MPED | |
| --- | --- | --- | --- | --- | --- | --- |
| | dependent | independent | dependent | independent | dependent | independent |
| SVM [62] | 83.99/09.72 | 56.73/16.29 | 56.61/20.05 | 37.99/12.52 | 32.39/09.53 | 19.66/03.96 |
| DGCNN [24] | 90.40/08.49 | 79.95/09.02 | 69.88/16.29 | 52.82/09.23 | 32.37/06.08 | 25.12/04.20 |
| DANN [37] | 91.36/08.30 | 75.08/11.18 | 63.07/12.66 | 47.59/10.01 | 35.04/06.52 | 22.36/04.37 |
| BiDANN [18] | 92.38/07.04 | 83.28/09.60 | 70.29/12.63 | 65.59/10.39 | 37.71/06.04 | 25.86/04.92 |
| A-LSTM [63] | 88.61/10.16 | 72.18/10.85 | 69.50/15.65 | 55.03/09.28 | 38.99/07.53 | 24.06/04.58 |
| BiHDM [25] | 93.12/06.06 | 85.40/07.53 | 74.35/14.09 | 69.03/08.66 | 40.34/07.53 | 28.27/04.99 |
| RGNN [20] | 94.24/05.95 | 85.30/06.72 | 79.37/10.54 | 73.84/08.02 | — | — |
| BiHDM w/o DA | 91.07/08.21 | 81.55/09.74 | 72.22/14.69 | 67.47/08.22 | 38.55/07.22 | 27.43/04.96 |
| RGNN w/o DA | — | 81.92/09.35 | — | 71.65/09.34 | — | — |
| GMSS | **96.48/04.63** | **86.52/06.22** | **86.37/11.45** | **73.48/07.41** | **40.16/06.08** | **28.49/04.42** |

— indicates the experiment results are not reported on that dataset.
Note: For the subject-dependent experiment, we calculate the average accuracy based on the results of all the sessions. While for the subject-independent experiment, we calculate the average accuracy based on the results of all the subjects.
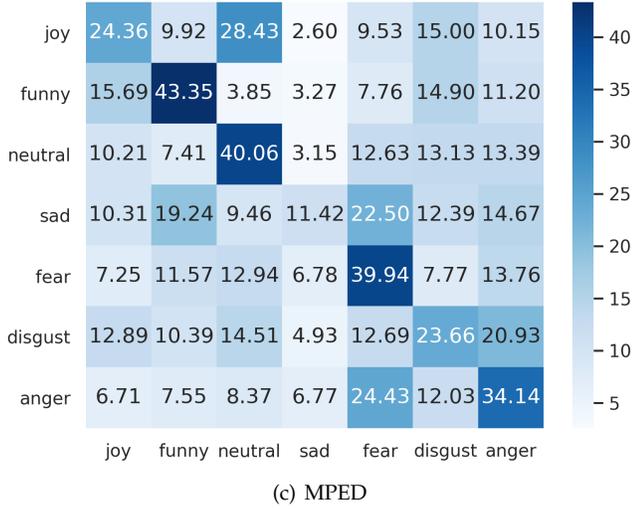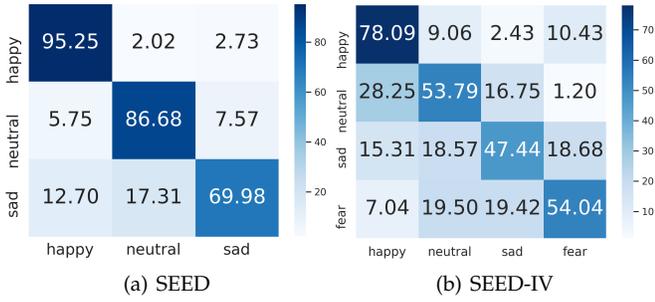
## 4.3 Experimental Details

In the experiments, the released differential entropy (DE) in SEED and SEED-IV, and the short-time Fourier transform (STFT) in MPED are feed into the model as input. The size of the input X is $62 \times 5$; the output dimensions of each electrode is 32; and $K = 2$, that is, the graph convolution aggregated the information of the second-order neighbors. In particular, GMSS is implemented by pytorch on a Nvidia 3080 GPU. The model is trained using the Adam optimizer with a batch size of 100. The learning rate is 0.001, and the weight decay rate is 8e-5. The mean accuracy (ACC) and standard deviation (STD) are employed as evaluation criteria in all datasets. The code of GMSS can be found at https://github.com/CHEN-XDU/GMSS.

## 4.4 Experimental Results

### 4.4.1 Unsupervised Mode

In the upstream task, the model is trained by self-supervised pretext tasks, consisting of two jigsaw puzzle tasks and one contrastive learning task. In the downstream task, the frozen feature extractor is applied and a linear classifier is used to evaluate the performance of GMSS. We compared GMSS with two self-supervised EEG emotion recognition methods SSL-EEG [54] and SeqCLR [55]. In addition, since there are few methods based on self-supervised EEG emotion recognition and the code is not released, we also compared with some popular self-supervised methods in other fields such as DeepCluster [61], MoCo [49], SwAV [48], SimCLR [50] and SimSiam [51]. These methods are reproduced and maintain the experimental protocol consistent with GMSS.

For a fair comparison, all of these methods of other fields adopted the same feature extraction operation as GMSS. To fit the EEG emotion recognition task, MoCo,

**(a) SEED**

|  | happy | neutral | sad |
|---|---|---|---|
| happy | 95.25 | 2.02 | 2.73 |
| neutral | 5.75 | 86.68 | 7.57 |
| sad | 12.70 | 17.31 | 69.98 |

**(b) SEED-IV**

|  | happy | neutral | sad | fear |
|---|---|---|---|---|
| happy | 78.09 | 9.06 | 2.43 | 10.43 |
| neutral | 28.25 | 53.79 | 16.75 | 1.20 |
| sad | 15.31 | 18.57 | 47.44 | 18.68 |
| fear | 7.04 | 19.50 | 19.42 | 54.04 |

**(c) MPED**

|  | joy | funny | neutral | sad | fear | disgust | anger |
|---|---|---|---|---|---|---|---|
| joy | 24.36 | 9.92 | 28.43 | 2.60 | 9.53 | 15.00 | 10.15 |
| funny | 15.69 | 43.35 | 3.85 | 3.27 | 7.76 | 14.90 | 11.20 |
| neutral | 10.21 | 7.41 | 40.06 | 3.15 | 12.63 | 13.13 | 13.39 |
| sad | 10.31 | 19.24 | 9.46 | 11.42 | 22.50 | 12.39 | 14.67 |
| fear | 7.25 | 11.57 | 12.94 | 6.78 | 39.94 | 7.77 | 13.76 |
| disgust | 12.89 | 10.39 | 14.51 | 4.93 | 12.69 | 23.66 | 20.93 |
| anger | 6.71 | 7.55 | 8.37 | 6.77 | 24.43 | 12.03 | 34.14 |

**(1) Subject-dependent experimental results**

**(d) SEED**

|  | happy | neutral | sad |
|---|---|---|---|
| happy | 84.16 | 10.24 | 5.60 |
| neutral | 12.07 | 70.27 | 17.66 |
| sad | 15.95 | 26.72 | 57.33 |

**(e) SEED-IV**

|  | happy | neutral | sad | fear |
|---|---|---|---|---|
| happy | 59.29 | 22.61 | 12.68 | 5.42 |
| neutral | 11.84 | 67.53 | 15.34 | 5.29 |
| sad | 18.91 | 28.59 | 45.38 | 7.12 |
| fear | 26.63 | 17.57 | 18.48 | 37.31 |

**(f) MPED**

|  | joy | funny | neutral | sad | fear | disgust | anger |
|---|---|---|---|---|---|---|---|
| joy | 11.75 | 24.41 | 27.37 | 4.10 | 13.78 | 10.08 | 8.50 |
| funny | 5.89 | 50.06 | 12.59 | 4.18 | 12.47 | 7.17 | 7.64 |
| neutral | 7.22 | 17.77 | 42.62 | 3.73 | 14.26 | 7.53 | 6.87 |
| sad | 6.30 | 25.04 | 17.61 | 6.93 | 20.97 | 11.17 | 11.99 |
| fear | 5.37 | 15.23 | 17.33 | 3.68 | 39.23 | 6.88 | 12.30 |
| disgust | 8.13 | 20.52 | 22.68 | 5.14 | 15.20 | 15.95 | 12.38 |
| anger | 6.39 | 18.38 | 16.74 | 4.98 | 23.43 | 10.83 | 19.25 |

**(2) Subject-independent experimental results**

Fig. 6. Confusion matrices in unsupervised mode. (a)-(c) and (d)-(f) are the subject-dependent and subject-independent results on SEED, SEED-IV and MPED datasets, respectively.
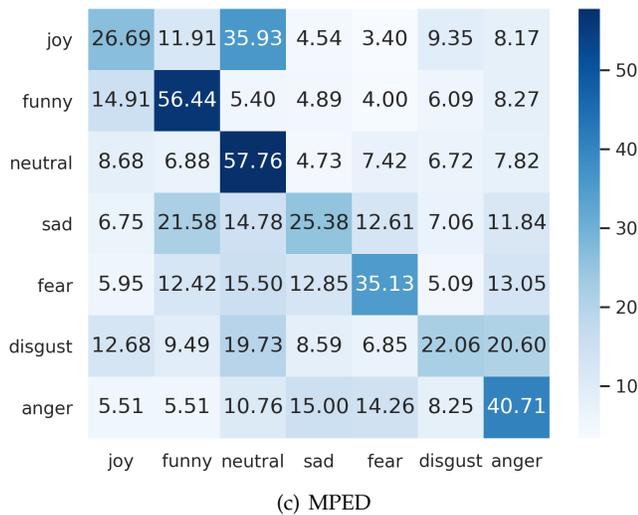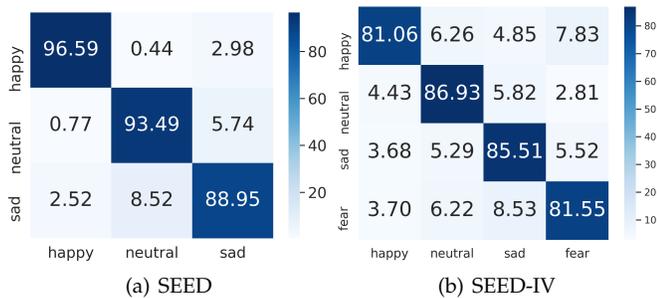
SwAV, SimCLR and SimSiam adopted the same data augmented as GMSS. The experimental results are shown in Table 2. Concretely, GMSS improves the accuracy by 5.86%, 8.52%, 2.02%, 8.51%, 4.34%, and 2.66% compared with the existing SOTA methods in the subject-dependent and subject-independent experiments on SEED, SEED-IV, and MPED datasets, respectively. Especially compared with MoCo, SwAV, SimCLR, SimSiam and SeqCLR which are also contrastive learning-based methods, GMSS achieves better results. This is attributed to GMSS having more positive and negative pairs (We set M = 8), and two more pretext tasks, that is, spatial and frequency jigsaw puzzle tasks, which are helpful in learning more discriminant and general EEG emotion representation. In summary, from the results of Table 2, in the unsupervised mode, it is observed that GMSS achieves an acceptable results without labels, making it more relevant to practical applications.

To better understand the confusion matrix of GMSS in recognizing different emotions, the unsupervised confusion matrices of all the experiments are displayed in Fig. 6. There are two observations:
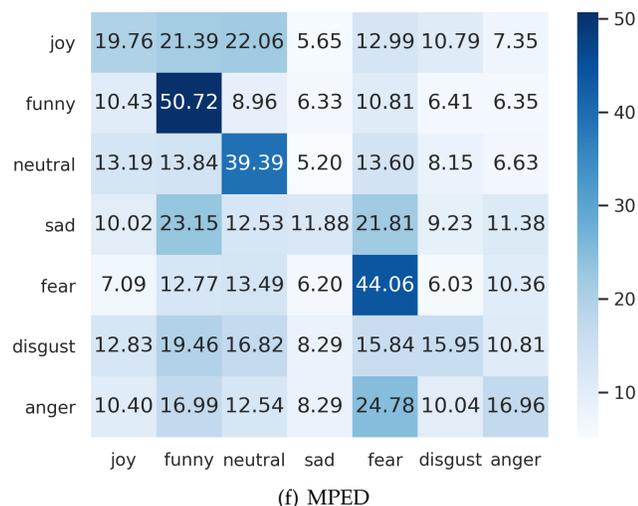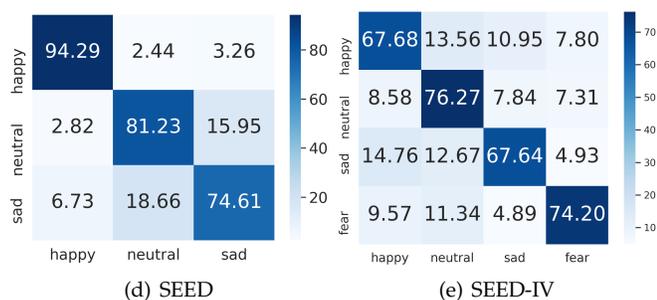
(1) For the subject-dependent experiment shown in Fig. 6(1), it is observed that happy is the easiest emotion recognized by SEED dataset. This is also observed in the results of the SEED-IV dataset. For MPED, which contains seven emotions, GMSS shows its superiority when identifying funny, neutral, fear, and anger. In addition, we can find joy is most easily confused with neutral. This may be because joy is more difficult to induce than other emotions.

(2) From the results of the subject-independent task shown in Fig. 6(2), for SEED, it is obvious that the accuracy of the happy emotion is much higher than neutral and sad, which is similar to the observation in Fig. 6(1). With SEED-IV, we can notice that neutral emotion achieves the highest accuracy since other emotions such as neutral lead to confusion. For MPED, funny, neutral, and sad emotions are much easier to recognize. It should be noted that, in the cross-subject task, the focus is on the sad emotion, which is difficult to identify from our observation.

### 4.4.2 Supervised Mode

In this section, a joint-training strategy is adopted. Based on the self-supervised training approaches, ground-truth emotion labels are used to train the feature extractor simultaneously. To evaluate the advantages of GMSS, the experiments conducted were the same as those of other methods, including linear support vector machine (SVM) [62], dynamical graph convolutional neural network (DGCNN) [24], regularized graph neural network (RGNN) [20], domain adversarial neural networks (DANN) [37], bi-hemisphere domain adversarial neural network (BiDANN) [18], attention-long short-term memory (A-LSTM) [63], and bi-hemispheric discrepancy model for EEG emotion recognition (BiHDM) [25]. All these methods are representative of previous studies on emotion recognition. Their results are directly quoted or reproduced from the literature to ensure a convincing comparison with the proposed method, and are summarized in Table 3.

**(a) SEED**

| | happy | neutral | sad |
|---|---|---|---|
| happy | 96.59 | 0.44 | 2.98 |
| neutral | 0.77 | 93.49 | 5.74 |
| sad | 2.52 | 8.52 | 88.95 |

**(b) SEED-IV**

| | happy | neutral | sad | fear |
|---|---|---|---|---|
| happy | 81.06 | 6.26 | 4.85 | 7.83 |
| neutral | 4.43 | 86.93 | 5.82 | 2.81 |
| sad | 3.68 | 5.29 | 85.51 | 5.52 |
| fear | 3.70 | 6.22 | 8.53 | 81.55 |

**(c) MPED**

| | joy | funny | neutral | sad | fear | disgust | anger |
|---|---|---|---|---|---|---|---|
| joy | 26.69 | 11.91 | 35.93 | 4.54 | 3.40 | 9.35 | 8.17 |
| funny | 14.91 | 56.44 | 5.40 | 4.89 | 4.00 | 6.09 | 8.27 |
| neutral | 8.68 | 6.88 | 57.76 | 4.73 | 7.42 | 6.72 | 7.82 |
| sad | 6.75 | 21.58 | 14.78 | 25.38 | 12.61 | 7.06 | 11.84 |
| fear | 5.95 | 12.42 | 15.50 | 12.85 | 35.13 | 5.09 | 13.05 |
| disgust | 12.68 | 9.49 | 19.73 | 8.59 | 6.85 | 22.06 | 20.60 |
| anger | 5.51 | 5.51 | 10.76 | 15.00 | 14.26 | 8.25 | 40.71 |

**(1) Subject-dependent experimental results**

**(d) SEED**

| | happy | neutral | sad |
|---|---|---|---|
| happy | 94.29 | 2.44 | 3.26 |
| neutral | 2.82 | 81.23 | 15.95 |
| sad | 6.73 | 18.66 | 74.61 |

**(e) SEED-IV**

| | happy | neutral | sad | fear |
|---|---|---|---|---|
| happy | 67.68 | 13.56 | 10.95 | 7.80 |
| neutral | 8.58 | 76.27 | 7.84 | 7.31 |
| sad | 14.76 | 12.67 | 67.64 | 4.93 |
| fear | 9.57 | 11.34 | 4.89 | 74.20 |

**(f) MPED**

| | joy | funny | neutral | sad | fear | disgust | anger |
|---|---|---|---|---|---|---|---|
| joy | 19.76 | 21.39 | 22.06 | 5.65 | 12.99 | 10.79 | 7.35 |
| funny | 10.43 | 50.72 | 8.96 | 6.33 | 10.81 | 6.41 | 6.35 |
| neutral | 13.19 | 13.84 | 39.39 | 5.20 | 13.60 | 8.15 | 6.63 |
| sad | 10.02 | 23.15 | 12.53 | 11.88 | 21.81 | 9.23 | 11.38 |
| fear | 7.09 | 12.77 | 13.49 | 6.20 | 44.06 | 6.03 | 10.36 |
| disgust | 12.83 | 19.46 | 16.82 | 8.29 | 15.84 | 15.95 | 10.81 |
| anger | 10.40 | 16.99 | 12.54 | 8.29 | 24.78 | 10.04 | 16.96 |

**(2) Subject-independent experimental results**

Fig. 7. Confusion matrices in supervised mode. (a)-(c) and (d)-(f) are the subject-dependent and subject-independent results on SEED, SEED-IV and MPED datasets, respectively.

For the subject-dependent experiments, in Table 3, it is observed that GMSS attains the best performance on three public EEG emotional datasets compared with all aforementioned methods above. In particular, the results on SEED and SEED-IV, with GMSS are 2.24% and 7% higher than those of the most advanced method RGNN. Meanwhile, it is also observed that GMSS achieves a performance very close to BiHDM on MPED, that is, 40.16% vs. 40.34%. This is because BiHDM is trained not only on the labeled training data but also on the unlabeled testing data. However, the GMSS is trained only on the training data. For a fair comparison, the domain discriminator of BiHDM is ablated and the experiments are conducted on the same input data as GMSS, which is denoted as BiHDM w/o DA. The experimental results show that GMSS improves the classification accuracy by 1.61% compared with BiHDM w/o DA. Furthermore, GMSS outperforms the BiHDM by 3.36% and 12.02% on SEED and SEED-IV datasets, respectively. These results verify that GMSS has a better discrimination capability under subject-dependent experiments. Additionally, our GMSS has a considerable running speed. On the SEED dataset of subject-dependent experiments, the average training time and average testing time for one epoch are 3762.7ms and 331.39ms respectively. Subject-independent experiment are also performed. It is observed that GMSS achieves the SOTA performance on SEED and MPED, which is 1.12% and 0.22% higher than the previous best method BiHDM, respectively. Moreover, GMSS achieves a performance close to that of RGNN on SEED-IV, that is, 73.48% vs. 73.84%, respectively. However, while RGNN removes its node-wise domain adversarial training component (NodeDAT), that is, training with the labeled training data as well as without the unlabeled testing data, denoted as RGNN w/o DA, the accuracy of RGNN w/o DA is 1.83% lower than that of GMSS. Furthermore, GMSS outperforms RGNN by 1.22% on the SEED dataset. In addition, compared with these advanced methods training without the unlabeled testing data, that is, BiHDM w/o DA and RGNN w/o DA, GMSS is 4.6%, 1.83%, and 1.06% higher, respectively. This indicates that our model can extract more general data representations for different subjects. Besides, compared with all baselines on all datasets and both experimental protocols, GMSS achieves the lowest standard deviation in accuracy, indicating the excellent discrimination and generalization capability of our model. We argue that the main reason can be attributed to the multi-task framework and self-supervised learning tasks.

Similar to the unsupervised mode, the confusion matrices of all experiments are also applied in the supervised mode to better understand the confusion of GMSS in recognizing different emotions as shown in Fig. 7. There are two observations:

(1) For the results of subject-dependent EEG emotion recognition experiment in Fig. 7(1), the classification accuracy for the three emotions is approximately 90% for the SEED dataset. In particular, for happy, the accuracy is above 95%. The happy and neutral emotions are easier to recognize than the sad emotion. For SEED-IV, which contains four emotions, we can notice that the accuracy of all emotions is above 80%. For MPED, which is a complex dataset that consists

of seven types of emotions, it is observed that funny and neutral emotions are much easier to recognize than other emotions. Moreover, for negative emotions, fear and anger are easier to recognize than sad and disgust.

(2) From the results of the subject-independent EEG emotion recognition experiment, for SEED, the happy emotion is much easier to be recognize than neutral and sad emotions. For SEED-IV, neutral and fear emotions are much easier to recognize. For MPED, which is a hard seven classification task, only funny, neutral and fear achieve acceptable results, which suggests that researchers should pay attention to joy, sad, disgust and anger in cross-subject emotion recognition.

## 4.5 Discussion

In this section, the representations of the visualization and ablation studies are presented.

### 4.5.1 Representation Visualization

To verify the discriminating ability of GMSS, the features obtained by GMSS on the MPED dataset are visualized. Fig. 8 shows the representation visualization of the subject-dependent experiment in supervised mode using t-distributed stochastic neighbor embedding (t-SNE) [66] on the MPED dataset. As shown in Fig. 8(1), it is difficult to separate the different classes from the original EEG data. However, for the learned EEG representation in Fig. 8(2), for the same emotion clusters, there are clear borders between different emotions, which verify that GMSS can discriminate features for EEG emotion recognition. Moreover, it is observed that the funny is more distinguishable than other emotions. This may be because funny induced more easily. In addition, comparing with Fig. 8(1) and Fig. 8(2), it is observed that GMSS has the potential to clarify the borders of various emotions and brings the same emotions closer together in feature space.

### 4.5.2 Ablation study

To assess the contribution of each essential pretext task in our model, experiments are conducted with the ablated GMSS models in both unsupervised and supervised modes. The ablation research verifies the influence of each pretext task and the combination of multiple tasks on the performance of EEG emotion recognition. In Table 4, the results are presented for the subject-dependent experiments in both unsupervised and supervised modes. In the unsupervised mode, GMSS-S, -F, and -C denote that the only spatial jigsaw puzzle task, frequency jigsaw puzzle task, and contrastive learning task are taken into consideration in the ablation model. Furthermore, GMSS-SF, -SC, -FC denote the spatial and frequency jigsaw puzzle tasks, spatial jigsaw puzzle and contrastive learning tasks, frequency jigsaw puzzle and contrastive learning tasks respectively taken into consideration by the ablation model, simultaneously. Similarly, in the supervised mode, GMSS-F, -C, -SF, -SC, and -FC represent the same ablation methods but are trained on the ground-truth emotion labels instead.

In the case of one self-supervised pretext task, GMSS-S achieves the best performance on four out of six results. This indicates that the spatial jigsaw puzzle task is extremely helpful in improving the discrimination of EEG emotional signals. Moreover, GMSS-F achieves the best performance on two out of six results, which implies that the frequency jigsaw puzzle task is helpful as well. The above results demonstrate that only one jigsaw puzzle task could improve the ability to distinguish EEG emotion signals. In the case of two tasks, GMSS-SF achieves the best performance on all datasets except MPED in the supervised mode, which is slightly lower than that of GMSS-SC. This further proves the effectiveness of the jigsaw puzzle task. In addition, compared with the corresponding results of only one task, the combination of the two tasks improve the accuracy of emotion recognition. This indicates that the three self-supervised tasks that were proposed are relevant and can promote model learning and more discriminative emotional representation. Furthermore, we can see that GMSS adopts all pretext tasks, achieving the best performance. This proves the effectiveness of our graph-based multi-task self-supervised learning framework.

### 4.5.3 Parameter analysis - Chebyshev filter size

As a hyper-parameter, Chebyshv filter size $K$, namely, $K$-order neighbor, will impact the performance of EEG emotion recognition. Thus, in this section, we conduct additional experiment to analyze the results of different Chebyshv filter size $K$ on SEED dataset. Here we set $K = 1, 2, ..., 10$ separately. And the results are shown in Fig. 9. It is obvious that GMSS achieves the best performance when $K = 2$. When $K$ is greater than 2, the performance of the model has a relatively noticeable downward trend. When $K$ is greater than 4, it tends to be stable gradually. We attribute the decline to the influence of over-smoothing.

## 5 CONCLUSION

In this paper, a graph-based multi-task self-supervised learning model is proposed for EEG emotion recognition. Our model is inspired by the multi-task learning theory and self-supervised learning theory, which combines different self-supervised tasks to improve model generalization and the ability to recognize EEG emotional signals. Several self-supervised tasks assist in improving the resilience of the model to emotion noise labels. The spatial pattern of EEG emotion signals is studied through the spatial jigsaw puzzle task. To reveal the intrinsic frequency bands for EEG emotion recognition, the frequency jigsaw puzzle task is employed, and the feature space is further standardized by the contrastive learning tasks. The experimental results validate the effectiveness of the proposed model. In future work, multi-task self-supervised learning will be further investigated to explore how to further improve EEG emotion recognition.

## REFERENCES

[1] R. J. Dolan, "Emotion, cognition, and behavior," *Science*, vol. 298, no. 5596, pp. 1191–1194, 2002.

(a) Subject-01　(b) Subject-02　(c) Subject-03　(d) Subject-04　(e) Subject-05

(1) Original EEG emotion data

(f) Subject-01　(g) Subject-02　(h) Subject-03　(i) Subject-04　(j) Subject-05

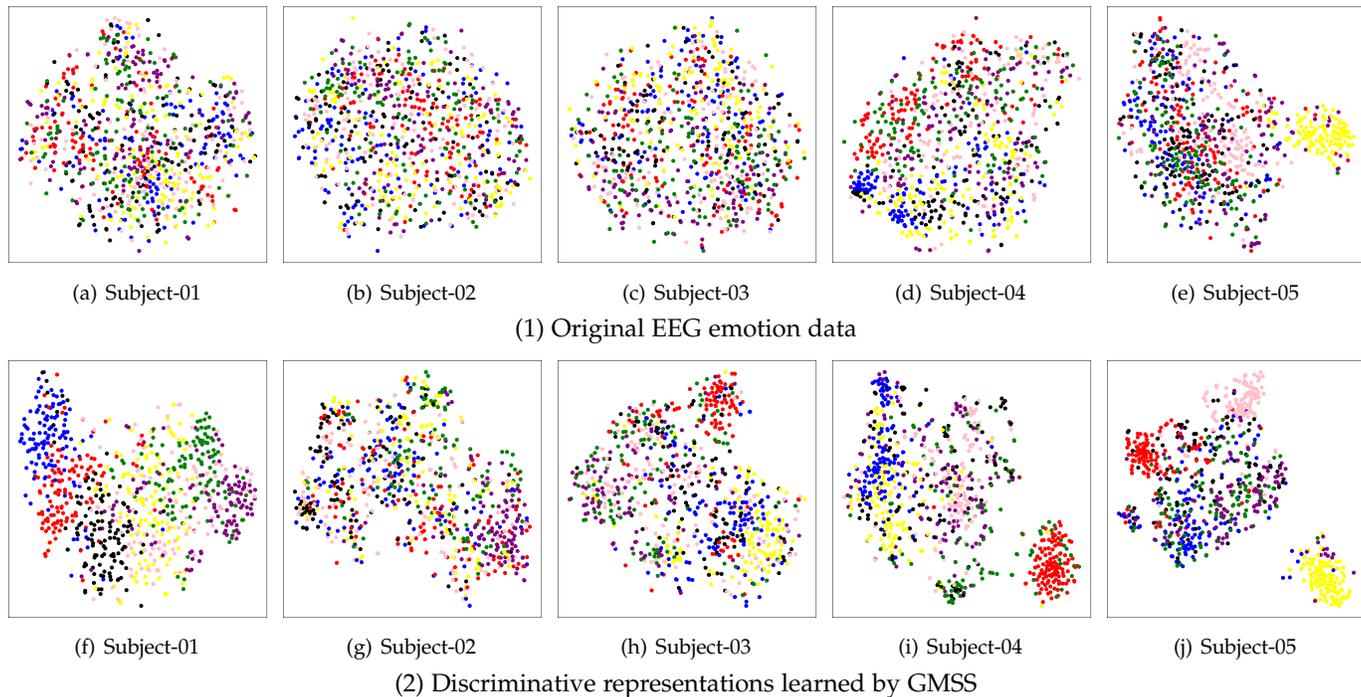(2) Discriminative representations learned by GMSS

Fig. 8. t-SNE visualization based on original EEG emotion data and discriminative representations learned by GMSS. (a)-(e) are the distributions of original EEG emotion data before being fed into network; (f)-(j) are the learned discriminative representations by GMSS. Blue, red, yellow, green, purple, black and pink dots denote joy, funny, neutral, sad, fear, disgust and anger emotions, respectively.

TABLE 4
Ablation study of subject-dependent classification accuracy (mean/std) for unsupervised mode and supervised mode on SEED, SEED-IV, and MPED datasets

| Ablation Models | SEED | | SEED-IV | | MPED | |
|---|---|---|---|---|---|---|
| | unsupervised | supervised | unsupervised | supervised | unsupervised | supervised |
| GMSS-S | **86.43/09.36** | **88.82/08.81** | **63.29/16.50** | 79.01/16.13 | 31.91/06.09 | **35.82/06.17** |
| GMSS-F | 84.84/10.68 | 86.75/08.64 | 62.31/16.24 | **79.56/14.31** | **33.32/06.45** | 34.98/06.13 |
| GMSS-C | 84.14/10.65 | 85.92/09.78 | 59.77/17.64 | 77.68/15.02 | 32.28/06.07 | 35.59/05.99 |
| GMSS-SF | **88.24/09.77** | **94.98/09.34** | **64.21/14.92** | **84.54/14.30** | 33.81/05.67 | 37.78/05.95 |
| GMSS-SC | 86.81/10.37 | 93.94/09.57 | 62.66/17.47 | 83.42/11.83 | 32.42/06.42 | **38.06/05.65** |
| GMSS-FC | 86.35/10.15 | 92.93/08.29 | 62.83/17.29 | 83.83/12.49 | 33.65/06.66 | 37.11/05.97 |
| GMSS | **89.18/09.74** | **96.48/04.63** | **65.61/17.33** | **86.37/11.45** | **34.81/06.88** | **40.16/06.08** |

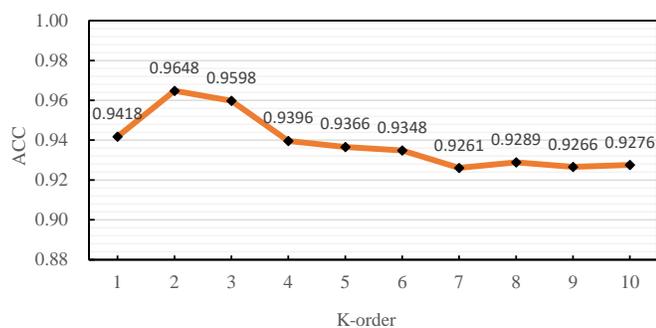

Fig. 9. Experiment results based on different Chebyshv filter sizes.

[2] S. Koelstra, C. Muhl, M. Soleymani, J.-S. Lee, A. Yazdani, T. Ebrahimi, T. Pun, A. Nijholt, and I. Patras, "Deap: A database for emotion analysis; using physiological signals," *IEEE Transactions on Affective Computing*, vol. 3, no. 1, pp. 18–31, 2011.

[3] J. C. Britton, K. L. Phan, S. F. Taylor, R. C. Welsh, K. C. Berridge, and I. Liberzon, "Neural correlates of social and nonsocial emotions: An fmri study," *Neuroimage*, vol. 31, no. 1, pp. 397–409, 2006.

[4] Y. Liu, O. Sourina, and M. K. Nguyen, "Real-time eeg-based emotion recognition and its applications," in *Transactions on Computational Science XII*. Springer, 2011, pp. 256–277.

[5] S. M. Alarcao and M. J. Fonseca, "Emotions recognition using eeg signals: A survey," *IEEE Transactions on Affective Computing*, vol. 10, no. 3, pp. 374–393, 2017.

[6] U. R. Acharya, V. K. Sudarshan, H. Adeli, J. Santhosh, J. E. Koh, and A. Adeli, "Computer-aided diagnosis of depression using eeg signals," *European Neurology*, vol. 73, no. 5-6, pp. 329–336, 2015.

[7] R. W. Picard, *Affective computing*. MIT press, 2000.

[8] Y.-P. Lin, C.-H. Wang, T.-P. Jung, T.-L. Wu, S.-K. Jeng, J.-R. Duann, and J.-H. Chen, "Eeg-based emotion recognition in music listening," *IEEE Transactions on Biomedical Engineering*, vol. 57, no. 7, pp. 1798–1806, 2010.

[9] R. Jenke, A. Peer, and M. Buss, "Feature extraction and selection for emotion recognition from eeg," *IEEE Transactions on Affective Computing*, vol. 5, no. 3, pp. 327–339, 2014.

[10] Y. Roy, H. Banville, I. Albuquerque, A. Gramfort, T. H. Falk, and J. Faubert, "Deep learning-based electroencephalography analysis: a systematic review," *Journal of Neural Engineering*, vol. 16, no. 5, p. 051001, 2019.

[11] R. T. Schirrmeister, J. T. Springenberg, L. D. J. Fiederer, M. Glasstetter, K. Eggensperger, M. Tangermann, F. Hutter, W. Burgard, and T. Ball, "Deep learning with convolutional neural networks for eeg decoding and visualization," *Human Brain Mapping*, vol. 38, no. 11, pp. 5391–5420, 2017.

[12] Z. Gao, X. Wang, Y. Yang, C. Mu, Q. Cai, W. Dang, and S. Zuo, "Eeg-based spatio–temporal convolutional neural network for driver fatigue evaluation," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 9, pp. 2755–2763, 2019.

[13] D. Zhang, L. Yao, X. Zhang, S. Wang, W. Chen, R. Boots, and B. Benatallah, "Cascade and parallel convolutional recurrent neural networks on eeg-based intention recognition for brain computer interface," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.

[14] T. Zhang, W. Zheng, Z. Cui, Y. Zong, and Y. Li, "Spatial–temporal recurrent neural network for emotion recognition," *IEEE Transactions on Cybernetics*, vol. 49, no. 3, pp. 839–847, 2018.

[15] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.

[16] Z. Wang, Y. Tong, and X. Heng, "Phase-locking value based graph convolutional neural networks for emotion recognition," *IEEE Access*, vol. 7, pp. 93 711–93 722, 2019.

[17] B. García-Martínez, A. Martinez-Rodrigo, R. Alcaraz, and A. Fernández-Caballero, "A review on nonlinear methods using electroencephalographic recordings for emotion recognition," *IEEE Transactions on Affective Computing*, 2019.

[18] Y. Li, W. Zheng, Z. Cui, T. Zhang, and Y. Zong, "A novel neural network model based on cerebral hemispheric asymmetry for eeg emotion recognition," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 2018, pp. 1561–1567.

[19] W.-L. Zheng and B.-L. Lu, "Investigating critical frequency bands and channels for eeg-based emotion recognition with deep neural networks," *IEEE Transactions on Autonomous Mental Development*, vol. 7, no. 3, pp. 162–175, 2015.

[20] P. Zhong, D. Wang, and C. Miao, "Eeg-based emotion recognition using regularized graph neural networks," *IEEE Transactions on Affective Computing*, 2020.

[21] S. Ruder, "An overview of multi-task learning in deep neural networks," *arXiv preprint arXiv:1706.05098*, 2017.

[22] F. Lotte, M. Congedo, A. Lécuyer, F. Lamarche, and B. Arnaldi, "A review of classification algorithms for eeg-based brain–computer interfaces," *Journal of Neural Engineering*, vol. 4, no. 2, p. R1, 2007.

[23] Y. Li, W. Zheng, L. Wang, Y. Zong, and Z. Cui, "From regional to global brain: a novel hierarchical spatial-temporal neural network model for eeg emotion recognition," *IEEE Transactions on Affective Computing*, 2019.

[24] T. Song, W. Zheng, P. Song, and Z. Cui, "Eeg emotion recognition using dynamical graph convolutional neural networks," *IEEE Transactions on Affective Computing*, vol. 11, no. 3, pp. 532–541, 2018.

[25] Y. Li, L. Wang, W. Zheng, Y. Zong, L. Qi, Z. Cui, T. Zhang, and T. Song, "A novel bi-hemispheric discrepancy model for eeg emotion recognition," *IEEE Transactions on Cognitive and Developmental Systems*, 2020.

[26] Z. Li, Z. Cui, S. Wu, X. Zhang, and L. Wang, "Fi-gnn: Modeling feature interactions via graph neural networks for ctr prediction," in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 2019, pp. 539–548.

[27] D. Nathani, J. Chauhan, C. Sharma, and M. Kaul, "Learning attention-based embeddings for relation prediction in knowledge graphs," *arXiv preprint arXiv:1906.01195*, 2019.

[28] S. Fan, J. Zhu, X. Han, C. Shi, L. Hu, B. Ma, and Y. Li, "Metapath-guided heterogeneous graph neural network for intent recommendation," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 2478–2486.

[29] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, "A comprehensive survey on graph neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.

[30] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," *arXiv preprint arXiv:1606.09375*, 2016.

[31] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.

[32] F. M. Bianchi, D. Grattarola, L. Livi, and C. Alippi, "Graph neural networks with convolutional arma filters," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

[33] G. Bouritsas, F. Frasca, S. P. Zafeiriou, and M. Bronstein, "Improving graph neural network expressivity via subgraph isomorphism counting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

[34] G. Ciano, A. Rossi, M. Bianchini, and F. Scarselli, "On inductive–transductive learning with graph neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 2, pp. 758–769, 2021.

[35] M. Tiezzi, G. Marra, S. Melacci, and M. Maggini, "Deep constraint-based propagation in graph neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

[36] I. Kokkinos, "Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6129–6138.

[37] A. Kendall, Y. Gal, and R. Cipolla, "Multi-task learning using uncertainty to weigh losses for scene geometry and semantics," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7482–7491.

[38] R. Collobert and J. Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *Proceedings of the 25th International Conference on Machine Learning*, 2008, pp. 160–167.

[39] J.-T. Huang, J. Li, D. Yu, L. Deng, and Y. Gong, "Cross-language knowledge transfer using multilingual deep neural network with shared hidden layers," in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2013, pp. 7304–7308.

[40] Y. Zhang and Q. Yang, "A survey on multi-task learning," *IEEE Transactions on Knowledge and Data Engineering*, 2021.

[41] O. Sener and V. Koltun, "Multi-task learning as multi-objective optimization," *arXiv preprint arXiv:1810.04650*, 2018.

[42] L. Jing and Y. Tian, "Self-supervised visual feature learning with deep neural networks: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[43] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *European Conference on Computer Vision*. Springer, 2016, pp. 69–84.

[44] Z. Yang, H. Yu, H. He, Z.-H. Mao, and A. Mian, "Self-supervised learning with fully convolutional networks," *arXiv preprint arXiv:2012.10017*, 2020.

[45] F. M. Carlucci, A. D'Innocente, S. Bucci, B. Caputo, and T. Tommasi, "Domain generalization by solving jigsaw puzzles," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2229–2238.

[46] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," *arXiv preprint arXiv:1803.07728*, 2018.

[47] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 132–149.

[48] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9912–9924, 2020.

[49] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9729–9738.

[50] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International Conference on Machine Learning*. PMLR, 2020, pp. 1597–1607.

[51] X. Chen and K. He, "Exploring simple siamese representation learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 750–15 758.

[52] Y. Liu, K. Wang, H. Lan, and L. Lin, "Temporal contrastive graph for self-supervised video representation learning," *arXiv preprint arXiv:2101.00820*, 2021.

[53] L. Lin, S. Song, W. Yang, and J. Liu, "Ms2l: Multi-task self-supervised learning for skeleton based action recognition," in

*Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 2490–2498.

[54] Z. Xie, M. Zhou, and H. Sun, "A novel solution for eeg-based emotion recognition," in *2021 IEEE 21st International Conference on Communication Technology*. IEEE, 2021, pp. 1134–1138.

[55] M. N. Mohsenvand, M. R. Izadi, and P. Maes, "Contrastive representation learning for electroencephalogram classification," in *Machine Learning for Health*. PMLR, 2020, pp. 238–253.

[56] R. Oostenveld and P. Praamstra, "The five percent electrode system for high-resolution eeg and erp measurements," *Clinical Neurophysiology*, vol. 112, no. 4, pp. 713–719, 2001.

[57] K. A. Lindquist, T. D. Wager, H. Kober, E. Bliss-Moreau, and L. F. Barrett, "The brain basis of emotion: a meta-analytic review," *The Behavioral and Brain Sciences*, vol. 35, no. 3, p. 121, 2012.

[58] W. Heller and J. B. Nitscke, "Regional brain activity in emotion: A framework for understanding cognition in depresion," *Cognition & Emotion*, vol. 11, no. 5-6, pp. 637–661, 1997.

[59] R. J. Davidson, "Affective style, psychopathology, and resilience: brain mechanisms and plasticity." *American Psychologist*, vol. 55, no. 11, p. 1196, 2000.

[60] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.

[61] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 132–149.

[62] J. A. Suykens and J. Vandewalle, "Least squares support vector machine classifiers," *Neural Processing Letters*, vol. 9, no. 3, pp. 293–300, 1999.

[63] T. Song, W. Zheng, C. Lu, Y. Zong, X. Zhang, and Z. Cui, "Mped: A multi-modal physiological emotion database for discrete emotion recognition," *IEEE Access*, vol. 7, pp. 12 177–12 191, 2019.

[64] W.-L. Zheng, W. Liu, Y. Lu, B.-L. Lu, and A. Cichocki, "Emotionmeter: A multimodal framework for recognizing human emotions," *IEEE Transactions on Cybernetics*, vol. 49, no. 3, pp. 1110–1122, 2018.

[65] W.-L. Zheng and B.-L. Lu, "Personalizing eeg-based affective models with transfer learning," in *Proceedings of the Twenty-fifth International Joint Conference on Artificial Intelligence*, 2016, pp. 2732–2738.

[66] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *Journal of Machine Learning Research*, vol. 9, no. 11, 2008.