

# GRAPH-BASED LEARNING OF SPECTRO-TOPOGRAPHICAL EEG REPRESENTATIONS WITH GRADIENT ALIGNMENT FOR BRAIN-COMPUTER INTERFACES

Prithila Angkan<sup>1,2</sup>, Amin Jalali<sup>1,2</sup>, Paul Hungler<sup>1,3</sup>, Ali Etemad<sup>1,2</sup>

<sup>1</sup>Ingenuity Labs Research Institute

<sup>2</sup>Department of Electrical and Computer Engineering

<sup>3</sup>Department of Chemical Engineering

Queen's University, Kingston, Canada

{prithila.angkan, amin.jalali, paul.hungler, ali.etemad}@queensu.ca

## ABSTRACT

We present a novel graph-based learning of EEG representations with gradient alignment (GEEGA) that leverages multi-domain information to learn EEG representations for brain-computer interfaces. Our model leverages graph convolutional networks to fuse embeddings from frequency-based topographical maps and time-frequency spectrograms, capturing inter-domain relationships. GEEGA addresses the challenge of achieving high inter-class separability, which arises from the temporally dynamic and subject-sensitive nature of EEG signals by incorporating the center loss and pairwise difference loss. Additionally, GEEGA incorporates a gradient alignment strategy to resolve conflicts between gradients from different domains and the fused embeddings, ensuring that discrepancies, where gradients point in conflicting directions, are aligned toward a unified optimization direction. We validate the efficacy of our method through extensive experiments on three publicly available EEG datasets: BCI-2a, CL-Drive and CLARE. Comprehensive ablation studies further highlight the impact of various components of our model.

**Index Terms**— EEG, BCI, Graph, Gradient alignment

## 1. INTRODUCTION

Electroencephalography (EEG) is a non-invasive technique that captures the electrical activity of the brain. Its cost-effectiveness and high temporal resolution make it widely used for brain-computer interfaces (BCI) in various research areas [1–3]. However, EEG presents challenges due to its low signal-to-noise ratio, subject-dependency, and low spatial resolution [4]. Prior EEG studies leverage information from various domains such as time, frequency, and topographical mapping to enhance representations [5, 6]. However, learning effective multi-domain representations from EEG poses two nuanced challenges. First, obtaining distinct class-specific clusters with large inter-class separation has proven challenging, especially in multi-domain setups [7]. Second, to

learn multi-domain information, gradient conflicts can arise, resulting in suboptimal training [8].

To address these challenges, we propose a novel approach using Graph-based learning of spectro-topographical EEG representations with Gradient Alignment (GEEGA). GEEGA encodes EEG from frequency-based topography maps and time-frequency spectrograms, maps embeddings onto a shared feature space using graph convolutional networks, and aligns gradients to reduce domain conflicts. Our method calculates class centers and pulls positive pairs toward them while pushing negatives apart for maximum inter-class separation. We evaluate our method on three publicly available EEG datasets, CLARE [9], CL-Drive [10], and BCI-2a [11]. Our approach achieves state-of-the-art performance across all three benchmarks.

The contributions in this work are summarized as follows. (1) We propose a new model, GEEGA, for EEG representation learning. Our model successfully learns multi-domain spectro-topographical information from EEG through graph-based fusion. (2) Our model effectively resolves gradient conflicts by aligning the gradients of the fused embeddings, ensuring that discrepancies, where gradients from each domain point in different directions, are addressed and guided toward a unified direction. This ensures balanced optimization across all domains causing the fused embeddings effectively capture complementary information from different domains, leading to enhanced performance. To the best of our knowledge, this is the first attempt to resolve gradient conflicts in the context of BCI as well as the first effort toward addressing such conflicts in a *multi-domain* setting in any context. (3) Moreover, our model incorporates class centers, enhancing inter-class separability by pulling positive pairs toward their respective class centers while pushing negative pairs apart. (4) GEEGA shows strong performances across several datasets and outperforms prior works. Detailed ablation studies demonstrate the positive impact of different components of our method.

## 2. RELATED WORK

Transformers have recently become popular in EEG representation learning. In [2], EEG-Deformer was proposed combining CNNs with transformers to capture coarse and fine-grained temporal dynamics. In [6] parallel transformers were used for spatial-temporal feature extraction with CNN integration, while [12] employed CNNs for channel-wise feature extraction followed by transformer processing. EEG channel-attention with Swin Transformer for motor imagery was integrated in [13] and [14] and utilized multi-dimensional global attention for spectral-spatial-temporal features. In [15] self-supervised masked autoencoders for cognitive load classification were applied, while [16] implemented Bayesian transformers for sleep staging.

Graph-based architectures have gained traction for EEG classification. GCN was used in [17] for sleep stage classification to learn intrinsic channel connections. In [18], graph and 1D convolutions were combined for intra- and inter-channel interactions, while [19] integrated GCNs with LSTMs for emotion classification. GCN and attention mechanisms were fused in [20] for structural relationships and long-range dependencies. Another graph-based network was used in [21], leveraging the spatial and temporal dependencies of EEG for emotion recognition. Finally [22] dynamically adjusted graph connections per instance using multi-level graph convolutions and coarsening.

## 3. METHOD

### 3.1. Problem Statement

Given a set of EEG signals,  $X = [X_1, X_2, \dots, X_c] \in \mathbb{R}^c$  with  $c$  channels, we aim to extract complementary representations: frequency domain  $E_{\text{freq}} \in \mathbb{R}^{M_1}$  and time-frequency domain  $E_{\text{time-freq}} \in \mathbb{R}^{M_2}$ , where  $M_1$  and  $M_2$  are the size of the embeddings. Training a unified multi-domain model faces the challenge of misaligned gradients. Specifically, the gradients measured by the loss function over a mini-batch  $B$  for the frequency domain ( $\nabla_B^{\text{freq}}$ ), for the time-frequency domain ( $\nabla_B^{\text{time-freq}}$ ), and the fused domain ( $\nabla_B^{\text{joint}}$ ), often point to conflicting directions, hindering effective training. Our goal is to align these gradients for unified optimization while achieving high inter-class separability.

### 3.2. Our Approach

**Multi-domain encoding.** We encode the pre-processed EEG signals  $X$  into multi-spectral topography maps  $X_{\text{topo}} \in \mathbb{R}^{B \times k \times h \times w}$  (frequency domain) and spectrograms  $X_{\text{spectro}} \in \mathbb{R}^{B \times c \times h \times w}$  (time-frequency domain), where  $B, k, c, h, w$  denote batch size, frequency bands, channels, height, and width respectively. Both inputs are flattened, linearly projected into token sequences [23], and positional encoding is added. The tokens are then fed to their respective transformer

branches:  $T_{\text{topo}}$  (frequency domain encoding) and  $T_{\text{spectro}}$  (time-frequency domain encoding). Producing embeddings  $E_{\text{freq}} \in \mathbb{R}^{M_1}$  (frequency domain) and  $E_{\text{time-freq}} \in \mathbb{R}^{M_2}$  (time-frequency domain), where  $M_1$  and  $M_2$  denote the size of the embeddings (see Fig. 1 (a)).

**Graph-based embedding fusion.** We fuse the embeddings  $E_{\text{freq}}$  and  $E_{\text{time-freq}}$  using a GCN module  $\Phi$ . The concatenated embedding  $E_{\text{concat}} \in \mathbb{R}^{B \times G_1}$  is projected to  $\tilde{E}_{\text{concat}} \in \mathbb{R}^{B \times G_2}$  where  $B$  is the batch size,  $G_1$  is the initial embedding dimension, and  $G_2$  is dimension of the higher-dimensional space, which is defined as  $G_2 = N \times F$ , where  $N$  is the number of nodes in the graph with  $F$  being the feature dimension of each node.  $\tilde{E}_{\text{concat}}$  is reshaped into  $\tilde{E}_{\text{node}} \in \mathbb{R}^{B \times N \times F}$  to form a graph structure.

In the first GCN layer, the learnable weight matrix  $W_1 \in \mathbb{R}^{F \times F}$  transforms the node features as:

$$O_{\text{GCN}_1} = \tilde{E}_{\text{node}} W_1, \quad O_{\text{GCN}_1} \in \mathbb{R}^{B \times N \times F}, \quad (1)$$

where  $O_{\text{GCN}_1}$  is the output from the first GCN layer. Node features are updated by aggregating neighboring information via adjacency matrix  $A \in \mathbb{R}^{N \times N}$ , forming a fully connected graph in our case as:

$$\tilde{E}_{\text{node-update}} = A \cdot O_{\text{GCN}_1}. \quad (2)$$

This process is repeated for the second GCN layer, followed by flattening and a linear transformation to produce the final feature vector of size  $H$ . A ReLU activation function is applied after each GCN layer to introduce non-linearity.

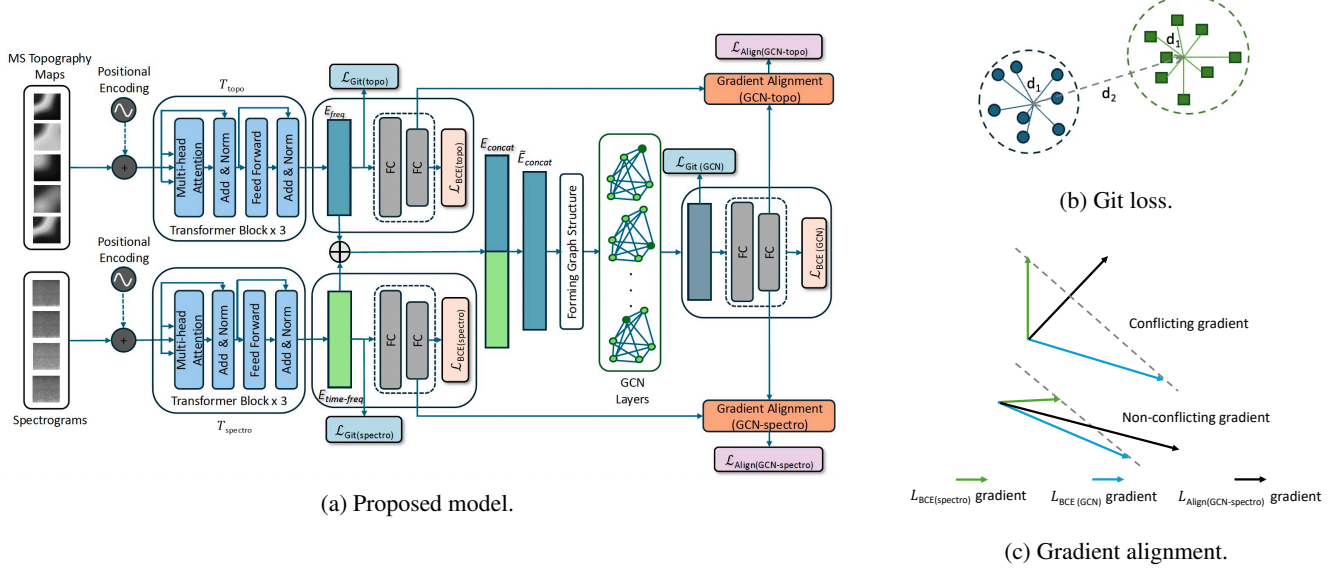
To train  $T_{\text{topo}}$ ,  $T_{\text{spectro}}$ , and the GCN, we use binary cross-entropy loss  $\mathcal{L}_{\text{BCE}}$  and Git loss [24]. Git loss is defined as:

$$\mathcal{L}_{\text{Git}} = \frac{1}{2} \sum_{i=1}^n \|E^i - c_y^i\|_2^2 + \sum_{i,j=1, i \neq j}^m \frac{1}{1 + \|E^i - c_y^j\|_2^2}, \quad (3)$$

where  $E^i$  is the feature vector of the  $i^{\text{th}}$  sample, and  $c_y^i$  is the center of the class to which  $E^i$  belongs.  $n$  and  $m$  are the total number of samples for the two classes, respectively. This loss combines center loss (first term of the equation) which reduces intra-class distances with pairwise difference loss (second part of the equation) which increases inter-class distances to enhance class separability as shown in Fig.1 (b).

**Gradient alignment.** Multiple domains in a single latent space can face the *gradient conflict* problem where the gradients from different domains may point at conflicting directions [8, 25] (see Fig.1(c)). This can result in sub-optimal training of the model and degrading of downstream performance. We align the two domains with respect to the fused domain rather than directly aligning the individual domains with each other as non-linear fusion reveals complex cross-domain interactions that remain hidden when domains are considered in isolation [26].

We define the gradients of losses computed over a mini-batch  $B$  as  $\nabla_B \mathcal{L}_{\text{BCE}(\text{topo})}$ ,  $\nabla_B \mathcal{L}_{\text{BCE}(\text{spectro})}$ , and  $\nabla_B \mathcal{L}_{\text{BCE}(\text{GCN})}$ .



**Fig. 1:** (a) The overview of our proposed network is depicted. (b) The concept of the Git loss is presented where we aim to minimize intra-class distances  $d_1$  and maximize inter-class distances  $d_2$ . (c) The concept of gradient alignment is presented.

When cosine similarity between gradients is negative ( $\cos \beta \leq 0$ , where  $\beta$  represents the angle between the gradients from different domains), conflicts exist. To resolve this, we use the Pareto optimization method that assigns weights  $\alpha^{\text{topo}}$ ,  $\alpha^{\text{spectro}}$ , and  $\alpha^{\text{GCN}}$  via a closed-form solution. The optimization problem for aligned gradient  $\mathcal{L}_{\text{Align(GCN-topo)}}$  is

$$\min_{\alpha^{\text{GCN}}, \alpha^{\text{topo}} \in \mathbb{R}} \left\| \alpha^{\text{GCN}} \nabla_B \mathcal{L}_{\text{GCN}} + \alpha^{\text{topo}} \nabla_B \mathcal{L}_{\text{topo}} \right\|^2, \quad (4)$$

subject to the constraints that  $\alpha^{\text{GCN}}, \alpha^{\text{topo}} \geq 0$  and  $\alpha^{\text{GCN}} + \alpha^{\text{topo}} = 1$ . Here, Eq. 4 minimizes the  $L_2$ -norm of the gradients within the convex hull of the gradient vectors  $\{\nabla_B \mathcal{L}_i\}_{i \in \{\text{GCN}, \text{topo}\}}$  [27]. The aligned gradient is:

$$h_{\text{GCN-topo}}^{\text{align}}(\theta) = 2\alpha^{\text{GCN}} \nabla_B \mathcal{L}_{\text{GCN}}(\theta) + 2\alpha^{\text{topo}} \nabla_B \mathcal{L}_{\text{topo}}(\theta), \quad (5)$$

where the resulting weights  $2\alpha^{\text{GCN}}$  and  $2\alpha^{\text{topo}}$  maintain the same weight summation (i.e.,  $2\alpha^{\text{GCN}} + 2\alpha^{\text{topo}} = 2$ ) and the model parameters  $\theta$  are updated as

$$\theta(t+1) = \theta(t) - \eta h_{\text{GCN-topo}}^{\text{align}}(\theta(t)). \quad (6)$$

Similar operations are performed for  $\mathcal{L}_{\text{Align(GCN-spectro)}}$  to align GCN and spectrogram gradients.

Finally, we define the total loss of GEEGA as:

$$\begin{aligned} \mathcal{L}_{\text{Total}} = & \mathcal{L}_{\text{Git(topo)}} + \mathcal{L}_{\text{Git(spectro)}} + \mathcal{L}_{\text{Git(GCN)}} \\ & + \mathcal{L}_{\text{BCE(topo)}} + \mathcal{L}_{\text{BCE(spectro)}} + \mathcal{L}_{\text{BCE(GCN)}} \\ & + \mathcal{L}_{\text{Align(GCN-topo)}} + \mathcal{L}_{\text{Align(GCN-spectro)}}. \end{aligned} \quad (7)$$

#### 4. EXPERIMENT SETUP

**Datasets.** We use three publicly available EEG datasets, namely BCI-2a [11], CL-Drive [10] and CLARE [9] for our

work. We use leave-one-subject-out (LOSO) evaluation. For BCI-2a, feet and tongue movement are used for binary classification, while for CL-Drive and CLARE, the subjective scores are binarized into low (1-5) and high (6-9) categories.

**Data preprocessing.** For BCI-2a, we use pre-processed data with each trial as an individual segment. For the other two datasets, we apply Butterworth bandpass filtering (1-75 Hz) and notch filtering following [10], then segment the signals into 10-second intervals. We generate multi-spectral topography maps and spectrograms from the segmented data.

**Multi-spectral topography maps.** To generate multi-spectral topography maps, we compute power spectral density (PSD) for each channel and five frequency bands: Delta, Theta, Alpha, Beta, and Gamma, following standard EEG practice [10, 28, 29]. Using Simpson's rule [30], we compute each band's power across all channels. These values are spatially mapped onto 2D grids using the international 10-20 electrode system with radial basis function (RBF) interpolation [31], creating multi-spectral topography maps of dimensions  $32 \times 32 \times 1$  for all datasets.

**Spectrograms.** While PSD captures power distribution across frequency bands, it fails to capture temporal dependencies. We address this using spectrograms containing time-frequency information. We compute Fast Fourier Transform (FFT) with non-overlapping 256-point windows, creating matrices where columns represent frequencies and rows represent time intervals. Spectrograms are generated for 4 channels (cognitive load datasets) or 22 channels (motor imagery dataset), each with dimensions  $32 \times 32 \times 1$ .

**Implementation details.** We use a batch size  $B$  of 32 and the Adam optimizer [34] (learning rate 0.0001, weight decay 0.00001). Training employs a Plateau scheduler (decay factor

**Table 1:** Performance compared to state-of-the-art solutions.

| Model              | BCI-2a             |                    | CL-Drive           |                    | CLARE               |                     |
|--------------------|--------------------|--------------------|--------------------|--------------------|---------------------|---------------------|
|                    | Accuracy           | F1                 | Accuracy           | F1                 | Accuracy            | F1                  |
| DGCNN [28]         | 65.29(9.26)        | 64.74(11.82)       | 65.77(4.71)        | 57.06(5.30)        | 61.84(3.96)         | 51.05(7.70)         |
| BiHDM [32]         | 67.86(9.29)        | 67.27(10.57)       | 62.01(15.57)       | 57.92(11.66)       | 68.14(16.43)        | 52.17(16.54)        |
| Conformer [33]     | 68.12(9.43)        | 67.53(11.25)       | 69.38(8.72)        | 63.29(9.29)        | 70.42(16.02)        | 58.28(12.00)        |
| MAE [15]           | 65.76(10.24)       | 65.98(10.92)       | 67.88(14.67)       | 61.25(13.18)       | 62.48(10.71)        | 57.51(7.29)         |
| VGG-style [10]     | 69.48(10.67)       | 69.73(10.24)       | 70.28(10.87)       | 63.12(9.39)        | 70.29(16.03)        | 60.24(13.16)        |
| DMMR [29]          | 65.57(10.23)       | 64.97(10.20)       | 61.15(13.74)       | 52.40(8.28)        | 69.02(22.07)        | 52.95(14.71)        |
| <b>GEEGA (our)</b> | <b>73.54(8.66)</b> | <b>72.86(8.04)</b> | <b>74.64(7.56)</b> | <b>64.53(8.24)</b> | <b>73.29(16.23)</b> | <b>60.68(14.42)</b> |

**Table 2:** Ablation experiments demonstrating the impact of each module within our proposed model. MS: multi-spectral topography maps, S: spectrograms, A: alignment.

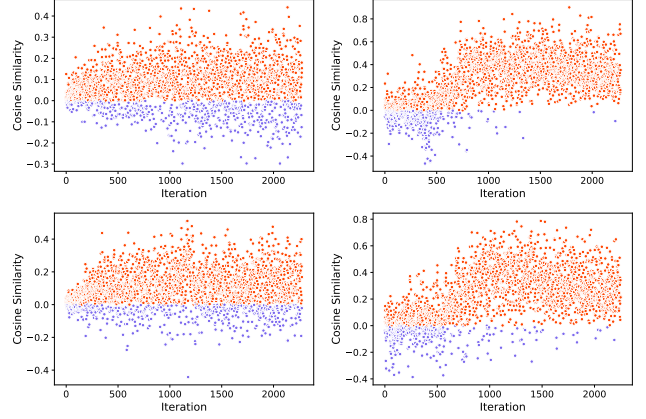
| MS | S | $\mathcal{L}_{Git}$ | A | BCI-2a             |                    | CL-Drive           |                    | CLARE               |                     |
|----|---|---------------------|---|--------------------|--------------------|--------------------|--------------------|---------------------|---------------------|
|    |   |                     |   | Accuracy           | F1                 | Accuracy           | F1                 | Accuracy            | F1                  |
| ✓  | ✓ | ✓                   | ✓ | <b>73.54(8.66)</b> | <b>72.86(8.04)</b> | <b>74.64(7.56)</b> | <b>64.53(8.24)</b> | <b>73.29(16.23)</b> | <b>60.68(14.42)</b> |
| ✓  | ✓ | ×                   | ✓ | 70.85(9.24)        | 69.20(9.83)        | 69.30(10.38)       | 60.07(7.72)        | 69.41(15.84)        | 54.28(12.30)        |
| ✓  | ✓ | ✓                   | × | 70.90(9.45)        | 69.87(9.73)        | 72.70(8.47)        | 62.65(6.81)        | 71.05(16.50)        | 56.39(14.00)        |
| ✓  | ✓ | ×                   | × | 69.48(8.29)        | 68.21(8.84)        | 70.20(8.99)        | 60.29(6.43)        | 70.07(16.50)        | 56.40(14.00)        |
| ✓  | × | ×                   | × | 66.00(9.41)        | 65.15(9.38)        | 67.52(9.29)        | 60.23(7.25)        | 66.82(17.24)        | 54.52(14.69)        |
| ×  | ✓ | ×                   | × | 66.43(9.36)        | 64.92(8.25)        | 67.46(8.57)        | 59.43(7.34)        | 70.50(15.38)        | 52.47(17.16)        |

0.1, patience 5) and warmup LambdaLR for first 5 epochs. Model is trained for 25 epochs on NVIDIA 2080 Ti using PyTorch. Both encoders  $T_{topo}$  and  $T_{spectro}$  use 3 transformer blocks with 8 attention heads, embedding dimension 512, and MLP hidden dimension 1024. The GCN module parameters are:  $G_1 = 1024$ ,  $G_2 = 1536$ ,  $N = 6$  nodes,  $F = 256$  and  $H = 512$ . This connects to FC layers (128, 1) with ReLU activation and 0.25 dropout. The FC layers after  $E_{freq}$  and  $E_{time-freq}$  use identical configurations. Dropout rates of 0.1 and 0.25 are applied to transformer/GCN blocks and FC layers respectively for regularization.

**Baseline methods.** We compare our proposed method with other popular and state-of-the-art recent works in EEG-based classification, and exclude methods requiring large-scale pre-training (EEGPT [35], BENDR [36]) following [21, 37].

## 5. RESULTS

**Performance.** We present the overall performance of our method in comparison to prior works in Table 1, where we observe that GEEGA achieves the best result across all three datasets. Notably, we observe that our method achieves higher accuracy and F1 scores than the two competing methods, the VGG-style CNN [10] and Conformer [33], by considerable margins. For instance, GEEGA outperforms the VGG by accuracy and F1 values of 4.06% and 3.13% respectively on the BCI-2a dataset, 4.36% and 1.41% on the CL-Drive dataset, and 3.00% and 0.44% on the CLARE dataset. Similarly, our method outperforms the widely used Conformer model by accuracy and F1 values of 5.42% and 5.33% on BCI-2a dataset, 5.26% and 1.24% on CL-Drive, and 2.87%



(a) w/o gradient alignment

(b) w/ gradient alignment

**Fig. 2:** The first row shows cosine similarities between multi-spectral topography maps and the fused domain, while the second row shows the same for spectrograms. Blue (values  $< 0$ ) indicates gradient conflicts, while red (values  $> 0$ ) indicates no conflict.

and 2.40% on CLARE. The results show that performance does not always correlate with the number of parameters, for instance, simpler models like VGG can still perform well.

**Gradient alignment.** Our key contribution is gradient alignment across domains to minimize conflicts and improve training. Fig. 2(a) shows misaligned gradients (positive (red) and negative (blue)) throughout training w/o our alignment process. Fig. 2(b) demonstrates reduced misaligned gradients as training progresses, confirming our alignment strategy’s effectiveness for both frequency and time-frequency domains.

**Ablation.** In Table 2, we present the results of detailed ablation experiments conducted to evaluate the impact of individual components in our method. We remove key components, including multi-spectral topography maps, spectrograms, the git loss ( $\mathcal{L}_{Git}$ ), and the alignment mechanism, and compare the results. We observe that our proposed GEEGA method with all the components achieves the best results compared to the other ablated combinations. Specifically, we observe that removing the git loss or the alignment step individually results in considerable drops in performance.

## 6. CONCLUSION

We propose GEEGA for EEG representation learning by integrating frequency and time-frequency domains using parallel transformer encoders and graph-based fusion. Our method addresses gradient conflicts through alignment strategies and enhances class separability using center loss with pairwise difference loss. Results on three benchmark datasets demonstrate superior performance over existing methods. In the future cross-task transferability and real-time applications can be explored.

## 7. REFERENCES

- [1] Shitao Zheng and Dongrui Wu, “Semi-supervised domain adaptation for eeg-based sleep stage classification,” in *ICASSP*, 2024, pp. 1776–1780.
- [2] Yi Ding, Yong Li, Hao Sun, Rui Liu, Chengxuan Tong, Chenyu Liu, Xinliang Zhou, and Cuntai Guan, “Eeg-deformer: A dense convolutional transformer for brain-computer interfaces,” *IEEE JBHI*, 2024.
- [3] Shivam Grover, Amin Jalali, and Ali Etemad, “Segment, shuffle, and stitch: A simple layer for improving time-series representations,” in *NeurIPS*, 2024.
- [4] He He and Dongrui Wu, “Transfer learning for brain–computer interfaces: A euclidean space data alignment approach,” *IEEE TBME*, vol. 67, no. 2, pp. 399–410, 2019.
- [5] Rui Li, Yiting Wang, and Bao-Liang Lu, “A multi-domain adaptive graph convolutional network for eeg-based emotion recognition,” in *ACMMM*, 2021, pp. 5565–5573.
- [6] Xiuzhen Yao, Tianwen Li, Peng Ding, Fan Wang, Lei Zhao, Anmin Gong, Wenya Nan, and Yunfa Fu, “Emotion classification based on transformer and cnn for eeg spatial–temporal feature learning,” *Brain sciences*, vol. 14, no. 3, pp. 268, 2024.
- [7] Maria Sayu Yamamoto, Khadijeh Sadatnejad, Toshihisa Tanaka, Md Rabiul Islam, Frédéric Dehais, Yuichi Tanaka, and Fabien Lotte, “Modeling complex eeg data distribution on the riemannian manifold toward outlier detection and multimodal classification,” *IEEE TBME*, 2023.
- [8] Yake Wei and Di Hu, “Mmpareto: Boosting multimodal learning with innocent unimodal assistance,” *ICML*, 2024.
- [9] Anubhav Bhatti, Prithila Angkan, Behnam Behinaein, Zunayed Mahmud, Dirk Rodenburg, Heather Braund, P. James Mclellan, Aaron Ruberto, Geoffrey Harrison, Daryl Wilson, Adam Szulewski, Dan Howes, Ali Etemad, and Paul Hungler, “Clare: Cognitive load assessment in realtime with multimodal data,” 2024.
- [10] Prithila Angkan, Behnam Behinaein, Zunayed Mahmud, Anubhav Bhatti, Dirk Rodenburg, Paul Hungler, and Ali Etemad, “Multimodal brain–computer interface for in-vehicle driver cognitive load measurement: Dataset and baselines,” *IEEE T-ITS*, 2024.
- [11] Clemens Brunner, Robert Leeb, Gernot Müller-Putz, Alois Schlögl, and Gert Pfurtscheller, “Bci competition 2008–graz data set a,” *Institute for knowledge discovery (laboratory of brain-computer interfaces), Graz University of Technology*, vol. 16, pp. 1–6, 2008.
- [12] Zhijiang Wan, Manyu Li, Shichang Liu, Jiajin Huang, Hai Tan, and Wenfeng Duan, “Eegformer: A transformer–based brain activity classification method using eeg signal,” *Front. Neurosci.*, vol. 17, pp. 1148855, 2023.
- [13] Han Wang, Lei Cao, Chenxi Huang, Jie Jia, Yilin Dong, Chunjiang Fan, and Victor Hugo C De Albuquerque, “A novel algorithmic structure of eeg channel attention combined with swin transformer for motor patterns classification,” *IEEE TNSRE*, 2023.
- [14] Yongling Xu, Yang Du, Ling Li, Honghao Lai, Jing Zou, Tianying Zhou, Lushan Xiao, Li Liu, and Pengcheng Ma, “Amdet: Attention based multiple dimensions eeg transformer for emotion recognition,” *IEEE Trans. Affect. Comput.*, 2023.
- [15] Dustin Pulver, Prithila Angkan, Paul Hungler, and Ali Etemad, “Eeg-based cognitive load classification using feature masked autoencoding and emotion transfer learning,” in *ICMI*, 2023, pp. 190–197.
- [16] Yuchen Liu and Ziyu Jia, “Bstt: A bayesian spatial-temporal transformer for sleep staging,” in *ICLR*, 2023.
- [17] Ziyu Jia, Youfang Lin, Jing Wang, Ronghao Zhou, Xiaojun Ning, Yuanlai He, and Yaoshuai Zhao, “Graphsleepnet: Adaptive spatial-temporal graph convolutional networks for sleep stage classification,” in *IJCAI*, 2020, vol. 2021, pp. 1324–1330.
- [18] Xuefen Lin, Jielin Chen, Weifeng Ma, Wei Tang, and Yuchen Wang, “Eeg emotion recognition using improved graph neural network with channel selection,” *Comput. Methods Programs Biomed.*, vol. 231, pp. 107380, 2023.
- [19] Yun Gu, Xinyue Zhong, Cheng Qu, Chuanjun Liu, and Bin Chen, “A domain generative graph network for eeg-based emotion recognition,” *IEEE JBHI*, vol. 27, no. 5, pp. 2377–2386, 2023.
- [20] Ming Jin, Changde Du, Huiguang He, Ting Cai, and Jinpeng Li, “Pgc: Pyramidal graph convolutional network for eeg emotion recognition,” *IEEE TMM*, 2024.
- [21] Chenyu Liu, Xinliang Zhou, Jiaping Xiao, Zhengri Zhu, Liming Zhai, Ziyu Jia, and Yang Liu, “Vsgt: variational spatial and gaussian temporal graph models for eeg-based emotion recognition,” in *IJCAI*, 2024, pp. 3078–3086.
- [22] Tengfei Song, Suyuan Liu, Wenming Zheng, Yuan Zong, and Zhen Cui, “Instance-adaptive graph for eeg emotion recognition,” in *AAAI*, 2020, vol. 34, pp. 2701–2708.
- [23] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *ICLR*, 2021.
- [24] Alessandro Calefati, Muhammad Kamran Janjua, Shah Nawaz, and Ignazio Gallo, “Git loss for deep face recognition,” *arXiv preprint arXiv:1807.08512*, 2018.
- [25] Xiaokang Peng, Yake Wei, Andong Deng, Dong Wang, and Di Hu, “Balanced multimodal learning via on-the-fly gradient modulation,” in *CVPR*, 2022, pp. 8238–8247.
- [26] Marcus Ghosh, Gabriel Béna, Volker Bormuth, and Dan FM Goodman, “Nonlinear fusion is optimal for a wide class of multisensory tasks,” *PLoS Comput. Biol.*, vol. 20, no. 7, pp. e1012246, 2024.
- [27] Jean-Antoine Désidéri, “Multiple-gradient descent algorithm (mgda) for multiobjective optimization,” *C. R. Math.*, vol. 350, no. 5-6, pp. 313–318, 2012.
- [28] Tengfei Song, Wenming Zheng, Peng Song, and Zhen Cui, “Eeg emotion recognition using dynamical graph convolutional neural networks,” *IEEE Trans. Affect. Comput.*, vol. 11, no. 3, pp. 532–541, 2018.
- [29] Yiming Wang, Bin Zhang, and Yujiao Tang, “Dmmr: Cross-subject domain generalization for eeg-based emotion recognition via denoising mixed mutual reconstruction,” in *AAAI*, 2024, vol. 38, pp. 628–636.
- [30] Daniel J Velleman, “The generalized simpson’s rule,” *Am. Math. Mon.*, vol. 112, no. 4, pp. 342–350, 2005.
- [31] Felix Havugimana, Kazi Ashraf Moinudin, and Mohammed Yeasin, “Deep learning framework for modeling cognitive load from small and noisy eeg data,” *IEEE TCDS*, 2023.
- [32] Yang Li, Lei Wang, Wenming Zheng, Yuan Zong, Lei Qi, Zhen Cui, Tong Zhang, and Tengfei Song, “A novel bi-hemispheric discrepancy model for eeg emotion recognition,” *IEEE TCDS*, vol. 13, no. 2, pp. 354–367, 2020.
- [33] Yonghao Song, Qingqing Zheng, Bingchuan Liu, and Xiaorong Gao, “Eeg conformer: Convolutional transformer for eeg decoding and visualization,” *IEEE TNSRE*, vol. 31, pp. 710–719, 2022.
- [34] Diederik P. Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *ICLR*, vol. abs/1412.6980, 2014.
- [35] Guangyu Wang, Wenchao Liu, Yuhong He, Cong Xu, Lin Ma, and Haifeng Li, “EEGPT: Pretrained transformer for universal and reliable representation of EEG signals,” in *NeurIPS*, 2024.
- [36] Demetres Kostas, Stephane Aroca-Ouellette, and Frank Rudzicz, “Bendr: Using transformers and a contrastive self-supervised learning task to learn from massive amounts of eeg data,” *Front. Hum. Neurosci.*, vol. 15, pp. 653659, 2021.
- [37] Qinke Ni, Hongyu Zhang, Cunhang Fan, Shengbing Pei, Chang Zhou, and Zhao Lv, “Dbpnet: Dual-branch parallel network with temporal-frequency fusion for auditory attention detection,” in *IJCAI*, 2024.