

GEFM: Graph-Enhanced EEG Foundation Model

Limin Wang¹, Toyotaro Suzumura² and Hiroki Kanezashi³

Abstract—Electroencephalography (EEG) signals provide critical insights for applications in disease diagnosis and healthcare. However, the scarcity of labeled EEG data poses a significant challenge. Foundation models offer a promising solution by leveraging large-scale unlabeled data through pre-training, enabling strong performance across diverse tasks. While both temporal dynamics and inter-channel relationships are vital for understanding EEG signals, existing EEG foundation models primarily focus on the former, overlooking the latter. To address this limitation, we propose Graph-Enhanced EEG Foundation Model (GEFM), a novel foundation model for EEG that integrates both temporal and inter-channel information. Our architecture combines Graph Neural Networks (GNNs), which effectively capture relational structures, with a masked autoencoder to enable efficient pre-training. We evaluated our approach using three downstream tasks and experimented with various GNN architectures. The results demonstrate that our proposed model, particularly when employing the GCN architecture with optimized configurations, consistently outperformed baseline methods across all tasks. These findings suggest that our model serves as a robust foundation model for EEG analysis.

I. INTRODUCTION

Understanding brain signals is crucial for clinical diagnosis, neurological disorder prediction, and exploring human cognition. The analysis of these signals facilitates early disease detection and prevention, forming the foundation of critical healthcare applications. Electroencephalography (EEG) is a widely used method for brain activity measurement, with its analysis techniques undergoing continuous advancements. Recent progress in machine learning have significantly enhanced EEG analysis, employing models ranging from traditional approaches like support vector machines to advanced architectures such as Transformers and Graph Neural Networks (GNNs).

Nowadays, the success of large language models (LLMs) has driven the development and widespread adoption of foundation models. These models leverage extensive amounts of unlabeled data during pre-training, achieving strong performance across a wide range of downstream tasks with minimal reliance on labeled data. Furthermore, foundation models eliminate the need for training from scratch for each task, significantly reducing computational time and cost. A notable characteristic of foundation models is their “emergent ability,” where larger architectures trained on vast datasets exhibit enhanced generalization capabilities. As a result, foundation models have gained traction across various

domains, including language [1], [2], image [3], [4], and are now being explored for EEG analysis.

Foundation models offer significant advantages for EEG applications. The scarcity of labeled EEG data, due to the high cost and expertise required for manual annotation, contrasts with the abundance of unlabeled data. By leveraging this unlabeled data during pre-training, foundation models can be fine-tuned for specific tasks with minimal reliance on labeled samples.

Several foundation models for EEG analysis, such as BENDR [5] and Neuro-GPT [6], have been proposed, primarily treating EEG as a collection of time-series data for individual channels (electrodes). These models often rely on advanced sequence modeling techniques, including Transformers and autoencoders, to process temporal information. However, EEG signals inherently involve complex interactions between channels, making it essential to incorporate inter-channel relationships into the modeling process. For instance, abnormalities in these interactions can serve as indicators of neurological disorders, highlighting the potential diagnostic value of capturing such relationships. Despite this, existing research has not yet adequately addressed the importance of modeling contextual relationships across EEG channels.

In this study, we propose Graph-Enhanced EEG Foundation Model (GEFM), a novel approach that integrates GNNs into a Transformer-based EEG foundation model to effectively capture inter-channel relationships. GNNs are particularly well-suited for modeling and representing complex relationships between entities, making them an ideal choice for capturing the intricate interactions among EEG channels. While the integration of GNNs with sequence modeling techniques, such as combining knowledge graphs (KGs) with LLMs, has been extensively explored in other fields, to the best of our knowledge, this study is the first to apply this combination to EEG analysis. This approach aims to establish a versatile and task-agnostic foundation model, offering significant potential for advancing EEG research and its applications.

To realize our proposed approach, we build upon BENDR [5], a masked autoencoder-based model recognized for its strong performance and adaptability for customization. In our architecture, we integrate GNNs to enhance the model’s ability to capture inter-channel relationships by representing EEG data as a graph structure. Specifically, each EEG channel is treated as a node, the signals from the corresponding channels serve as node features, and the connections between channels are modeled as edges.

While this integration offers significant potential, it also

¹Limin Wang (wang-limin516@g.ecc.u-tokyo.ac.jp),

²Toyotaro Suzumura (suzumura@acm.org) and ³Hiroki Kanezashi (hkanezashi@acm.org) are with The University of Tokyo, 7-3-1 Hongo, Tokyo, Japan

introduces a key technical challenge. As a foundation model, the architecture must be applicable across diverse datasets. However, GNNs typically require fixed-length node features, whereas EEG datasets often vary in signal sequence lengths due to varying task objectives or experimental setups. This inconsistency makes it impractical to directly use raw EEG signals as node features across all datasets. To address this, we introduce a sequence length adjustment mechanism that standardizes EEG signal lengths to a predefined target length before they are fed into the GNNs. This ensures compatibility across datasets while maintaining model’s flexibility as a versatile foundation model for EEG analysis.

This study aims to address the following research questions:

- RQ1. Which GNN architectures demonstrate the best performance in the proposed framework?
- RQ2. For which types of downstream tasks does the integration of GNNs and sequence modeling techniques provide the most significant improvements?
- RQ3. What is the most effective strategy for adjusting sequence lengths of EEG signals within the proposed framework?
- RQ4. Which base model, when combined with GNNs, yields the most robust and versatile foundation model for EEG analysis?

To address the research questions, we conducted experiments using three GNN architectures: Graph Convolutional Networks (GCN) [7], Graph Attention Networks (GAT) [8], and GraphSAGE [9]. The models were pre-trained on a large-scale dataset and evaluated on three downstream tasks to assess their performance. To account for the variability in EEG sequence lengths, we examined two adjustment strategies. Furthermore, we tested of two configurations of BENDR [5], which have demonstrated high accuracy in prior studies, as base models integrated with GNNs to develop a robust foundation model for EEG analysis. The results show that our proposed model, particularly when employing the GCN architecture with optimized configurations, consistently outperformed baseline methods across all tasks.

The contributions of this paper are as follows:

- We propose a novel foundation model specifically designed for EEG analysis.
- We integrate a robust sequence modeling techniques with GNNs to effectively capture both temporal dynamics and inter-channel relationships in EEG data.
- We conduct extensive experimental evaluations across three downstream tasks using three GNN architectures and multiple configurations to identify the optimal setup, demonstrating the effectiveness and versatility of our model.

II. BACKGROUND AND RELATED WORK

A. Foundation Models

Foundation models are large-scale pre-trained models designed to serve as general-purpose frameworks across a

wide range of downstream tasks. These models are typically trained on extensive datasets using self-supervised learning, enabling the use of large amounts of unlabeled data. By leveraging their learned representations, foundation models can be adapted to specific tasks with minimal fine-tuning, demonstrating strong performance and versatility. This approach reduces the time and cost associated with task-specific model development. Foundation models have shown significant potential across various domains, including natural language processing (NLP) [1], [2] and computer vision (CV) [3], [4].

B. EEG Foundation Models

Foundation models have recently been introduced in the EEG field to address key challenges, such as the difficulty of large-scale data labeling. Among these, BENDR [5] is a notable example, leveraging techniques inspired by BERT [10] and wav2vec [11]. First, raw multi-channel EEG signals are passed through six convolutional layers (referred to as the “BENDR Encoder”) to generate convolved features. A portion of the convolved features is then masked, and a Transformer encoder reconstructs the masked features using the unmasked features as context. The model is optimized by calculating the contrastive loss between the reconstructed features and the original convolved features before masking.

MAEEG [12] is an improved model that builds upon the BENDR [5] method, achieving enhanced performance. Specifically, a linear layer and a convolutional layer are added after the Transformer encoder to directly reconstruct the EEG signal. And then instead of computing the contrastive loss between the original features and the reconstructed features, MAEEG calculates the reconstruction loss between the input EEG signal and the reconstructed EEG signal.

Neuro-GPT [6] is another model that has achieved higher performance than BENDR [5]. This model is based on the Generative Pre-trained Transformer (GPT) [13] method and uses a decoder-only Transformer Encoder for the encoded features by a new EEG Encoder, which consists of a convolutional module and a Transformer Encoder.

However, to the best of our knowledge, existing EEG foundation models focus exclusively on learning time-series information, while neglecting inter-channel relationships, which are also critical for capturing the underlying dynamics of EEG signals.

C. EEG Models with Graphs

Several EEG models have been developed to capture inter-channel relationships, although these are not foundation models but are instead trained from scratch for specific tasks. These models represent EEG signals as graphs, allowing the study of network properties such as channel connectivity. To process these graph-structured data, they employ Graph Neural Networks (GNNs), which are specifically designed for such applications [14]. One notable example is EEG-GCNN [15], which utilizes a Graph Convolutional Network (GCN) [7] architecture with edge weights calculated based

on spatial distances. This model learns inter-channel relationships effectively by leveraging these edge weights.

D. Graph-Enhanced Models in Other Domains

While no foundation model in the EEG domain has been developed to learn both time-series information and inter-channel relationships, similar approaches have been explored in the domain of NLP.

One example is the pre-training language models (LM) with knowledge graphs (KGs), such as K-BERT [16], which incorporates knowledge graphs into the pre-training process to enhance the model’s understanding of entities and relationships. Some research has proposed the combination of LMs and GNNs to jointly learn on texts and KGs [17], [18], [19]. We can apply this approach of integrating GNNs with LMs to EEG data, where the GNNs can capture the relationships between EEG channels, while the LMs can be replaced with Transformer-based sequence models to learn the temporal dynamics of the signals.

III. METHODS

We propose a foundation model that learns both the inter-channel relationships and the time-series dynamics of EEG signals. Since existing EEG foundation models primarily focus on learning time-series information, we extend these models by integrating inter-channel relationship learning. Among the state-of-the-art EEG foundation models, BENDR [5] is notable for its strong performance and adaptability for customization. Therefore, we adopt BENDR as the base foundation model for our proposed approach. In this section, we first provide a brief overview of BENDR, followed by the description of our proposed method.

A. BENDR

The model architecture of BENDR is illustrated in (a), (c) and (e) in Figure 1.

During pre-training, raw multi-channel EEG signals are passed through six convolutional layers (referred to as the “BENDR Encoder”) to generate convolved features. The procedure is inspired by wav2vec [11]. A portion of the convolved features is then masked, and a Transformer encoder reconstructs the masked features using the unmasked features as context. The model is optimized by calculating the contrastive loss between the reconstructed features and the original convolved features before masking.

For downstream tasks, several model configurations have been proposed. Among these, we focus on two configurations (“BENDR” and “Linear”) that achieved the highest performance according to [5]. In both configurations, the process up to encoding by the “BENDR Encoder” is identical to that of the pre-training phase. However, the subsequent steps differ. In the “BENDR” configuration, the convolved features are passed to the pre-trained Transformer encoder, followed by a linear layer for classification. In contrast, in the “Linear” configuration, the pre-trained Transformer encoder is omitted entirely, and the convolved features are aggregated and passed to a linear classification layer.

B. Our Proposal: GEFM

We propose Graph-Enhanced EEG Foundation Model (GEFM), extending the BENDR architecture described above by incorporating the learning of inter-channel relationships. Inter-channel relationships in EEG signals can be naturally represented using a graph structure, where each channel corresponds to a node and the edges represent the connectivity between channels. Graph Neural Networks (GNNs) are particularly well-suited for modeling such relationships in graph-structured data. Therefore, we propose integrating GNNs with BENDR to enhance its capability of learning inter-channel relationships.

To achieve this, we employ a simple yet effective strategy by inserting a two-layer GNN directly before the inputs are processed by the BENDR Encoder, in both the pre-training and downstream tasks. The GNN assumes an input graph structure $G = (V, E)$ defined as follows:

- Each node in V represents an EEG recording channel (electrode), where $|V| = C$, corresponding to the total number of channels.
- The node features are the EEG recordings for each channel, with a feature length equal to the sequence length n of the recordings.
- Following EEG-GCNN [15], the graph is fully connected, and each edge $(u, v) \in E$ has a weight. This setup enables the GNN to capture all pairwise connectivity between channels. The edge weights are represented as a matrix $W \in \mathbf{R}^{C \times C}$.
- The edge weights are defined as the reciprocal of the geodesic distance between two channels on the spherical scalp model, based on the hypothesis that closer channels interact more strongly. The geodesic distance D_{ij} between channels i and j is computed using their 3D coordinates (x_i, y_i, z_i) and (x_j, y_j, z_j) , along with the sphere’s radius r , as:

$$D_{ij} = \arccos\left(\frac{x_i x_j + y_i y_j + z_i z_j}{r^2}\right)$$

The coordinates and radius values are adopted from EEG-GCNN. These edge weights are utilized by the GNN to learn inter-channel relationships effectively.

Our proposed architecture is illustrated in (b), (d) and (f) in Figure 1. The GNNs in this architecture are pre-trained and subsequently fine-tuned for downstream tasks. The process following the “BENDR Encoder” remains identical to that of the base model, including the use of two model configurations, “BENDR” and “Linear”, for downstream tasks.

As a foundation model, this architecture must be applicable across diverse datasets. However, a key challenge arises because GNNs require fixed-length node features, whereas EEG datasets often vary in signal sequence lengths due to differing task objectives. To address this, we introduce a sequence length adjustment mechanism, such as padding or utilizing a linear layer, that standardizes EEG signal sequence lengths to a predefined target. This ensures compatibility across datasets while maintaining the model’s versatility for

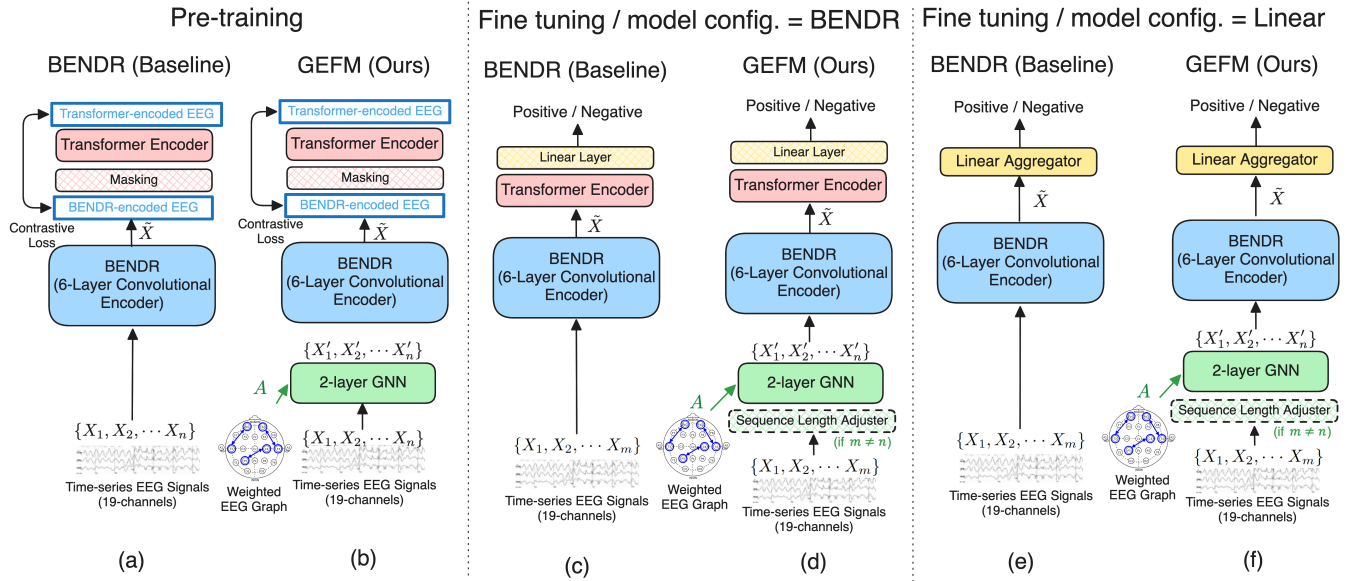


Fig. 1. Comparison between two model architectures during pre-training and fine tuning. Figures (a), (c) and (e) represent BENDR [5], while (b), (d) and (f) represent our proposed model, GEFM.

EEG analysis. If the original sequence length is already equal to the predefined target, no adjustment is applied to avoid unnecessary data manipulation.

C. Alternative Architectural Variants

Our proposed method aims to learn both temporal and spatial information in EEG signals as a foundation model. To achieve this, we leverage the strengths of existing EEG foundation models in learning temporal information and add a module to learn spatial relationships between channels.

There are three possible ways to combine these modules:

- 1) Learning spatial relationships first and then temporal information (This is the method we adopted, as described above).
- 2) Learning temporal information first and then spatial relationships.
- 3) Learning both in parallel and aggregating the results.

However, approach 2 is not suitable for our purpose because it would require placing the GNN after the BENDR Encoder and Transformer Encoder, which would integrate all channel information and make it difficult to learn spatial relationships.

Approach 3 is thought to be feasible, but we chose to focus on approach 1 in this study because it allows us to reuse the BENDR Encoder and pretext task with minimal modifications. Approach 3 is worth exploring in future research.

In this study, we adopt approach 1, where we first learn spatial relationships using a GNN and then learn temporal information using the BENDR Encoder and Transformer Encoder. This approach enables us to effectively learn both spatial and temporal information in EEG signals.

D. Incorporating Graphs Into Other EEG Foundation Models

In addition to BENDR [5], there also exist other EEG foundation models. Although we did not conduct experiments on them, we believe that a similar approach can be applied to them as well. Here we provide a brief overview of how our proposed method can be applied to some other EEG foundation models, such as MAEEG [12] and NeuroGPT [6].

Both models first use a Convolutional Encoder to process the input EEG signals and then use a Transformer Encoder to capture temporal information. We consider a simple yet effective method, inserting a GNN before the Convolutional Encoder to learn the relationships between channels, will be effective for these models. By incorporating graph structures into these EEG foundation models, we can potentially improve their performance and better capture the complex relationships between channels in EEG data.

IV. EXPERIMENTS

We conducted experiments, including both pre-training and downstream tasks, under various conditions to evaluate the performance of GEFM and gain insights into the research questions.

A. Dataset

1) *Pre-training Dataset*: Following BENDR [5], we utilized the Temple University Hospital EEG Corpus (TUEG) [20] for pre-training. TUEG provides a diverse range of subjects and includes recordings across multiple sessions over extended time periods, making it an ideal dataset for pre-training foundation models. For our study, we specifically focused on version 2 of this dataset, which consists of clinical recordings from over 10,000 individuals.

To accelerate experimentation and evaluation during development, we downsampled the dataset to one-tenth of its original size.

2) *Downstream Datasets*: We evaluated GEFM using the datasets from the following tasks.

- MMI [21], [22] This task involves predicting whether the participant is imagining the movement of the right (positive) or left (negative) hand.
- P300 [21], [23] This task involves predicting whether the participant focused on a flashing target letter (positive) or a non-target letter (negative).
- ERN [24] This task involves predicting whether the participant’s attempt to input a character using a P300 speller was recognized correctly (positive) or incorrectly (negative).

All three tasks are binary classification problems, which were previously used in BENDR [5]. These tasks involve EEG signals recorded from a sufficient number of channels in the 10/20 channel scheme [25], making them suitable for constructing graphs for GNN-based approaches. Detailed information about these datasets is provided in Table I.

The P300 and ERN datasets exhibit imbalanced class distributions, while the MMI dataset is balanced. To address class imbalance during fine tuning, we followed the methodology presented in BENDR [5] and applied the following steps:

- For imbalanced datasets, during training we performed undersampling of the majority class to equalize the number of samples across classes.
- During testing, we evaluated performance using metrics that account for class imbalance, specifically AUROC for P300 and ERN. For MMI, which has balanced classes, we used Accuracy to assess test performance.

3) *Preprocessing*: We applied the following preprocessing steps to both the pre-training and downstream datasets, following the methodology presented in BENDR [5]:

- To standardize the sampling frequency across datasets, we ensured that all recordings had a frequency of 256 Hz by applying over- or undersampling as necessary.
- We utilized 19 EEG channels from the 10/20 channel scheme [25] and ignored all other channels.
- For pre-training, 60-second sequences were extracted from the pre-training dataset for use in general training, while 20-second sequences were specifically used for the P300-related training, as described in BENDR [5]. For downstream tasks, the entire length of the sequences was used.

B. Setup

1) *GNN architectures*: We evaluated the performance of GEFM by individually incorporating three standard GNN architectures: Graph Convolutional Networks (GCN) [7], Graph Attention Networks (GAT) [8] and GraphSAGE [9]. Each architecture was used consistently throughout a single set of pre-training and downstream tasks. All GNN implementations were based on PyTorch 2.3.1 and PyTorch Geometric 2.5.3.

Our experiments included five configurations: GCN and GAT with and without edge weights, and GraphSAGE without edge weights. While GCN and GAT were configured to utilize edge weights, using the `edge_weight` and `edge_attr` inputs, respectively, GraphSAGE does not natively support edge weights in its standard implementations in PyTorch Geometric 2.5.3. Consequently, GraphSAGE was tested without edge weights. To ensure a fair comparison, we also tested GCN and GAT without edge weights.

2) *Sequence Length Adjusters*: As described in the previous section, all sequence lengths must be standardized to a specific value, n , to meet the requirements of a foundation model. To fully leverage the information available during pre-training, we fixed n to the sequence length of the pre-training dataset. Consequently, for downstream datasets with sequence lengths m , we adjusted m to match n .

Since m is typically smaller than n , we explored two simple yet effective adjustment methods. The first inserts a linear layer of size $m \times n$ immediately before the GNN. The second uses padding, where the last value of the original signal is repeated and appended to the sequence until the sequence reaches n .

C. Results

The results are presented in Tables II and III. The baseline corresponds to the original BENDR [5]. Note that the evaluation metric for MMI is Accuracy, while AUROC is used for P300 and ERN. Those with the statement “(with edge weights)” in the tables indicates that the GNN utilized edge weights. For the experiments with *padding* as the sequence length adjuster (shown in Table III), GCN and GAT without edge weights were excluded from these experiments due to their poor performance with the linear layer (see Table II). The following discussion addresses key points related to our research questions.

1) *Comparison of GNN Architectures (RQ1)*: Tables II and III indicate that our proposed approach performed better when using a linear layer for the sequence length adjustment compared to padding. Therefore, to analyze the performance variations introduced by different GNN architectures, the following discussion focuses on the results obtained with a linear layer.

As shown in Table II, among all the GNN architectures tested, only “GCN with edge weights” consistently outperformed the baseline across all three downstream tasks. Thus, “GCN with edge weights” emerges as the most suitable architecture for incorporation into foundation models for EEG analysis. The next best-performing architecture is “GraphSAGE”, which exceeded the baseline in two out of three tasks, making it another promising candidate. Investigating the underlying factors contributing to this behavior remains an open question and is expected to be addressed in future work.

2) *Differences of the Effect of GNNs Arising from Downstream Tasks (RQ2)*: This section examines how the effect of introducing GNNs varies across downstream tasks by comparing the performance differences between the baseline

Dataset	sfreq. (Hz)	Length (s)	Num of Ch.	Subjects	Folds
MMI [21], [22]	160	6	64	105	5
P300 [21], [23]	2,048	2	64	9	9
ERN [24]	200	2	56	26(10)	4

TABLE I

DOWNSTREAM DATASET BATTERY AND NUMBER OF CROSS-VALIDATION FOLDS USED, FOLLOWING BENDR [5].

Model	Config.	MMI	P300	ERN
Baseline	BENDR	0.646	0.577	0.522
	Linear	0.794	0.607	0.508
GraphSAGE	BENDR	<u>0.883</u>	<u>0.692</u>	0.501
	Linear	0.758	0.580	0.492
GCN	BENDR	0.514	<u>0.616</u>	<u>0.534</u>
	Linear	0.506	0.578	0.486
GCN (with edge weights)	BENDR	0.849	0.616	0.538
	Linear	0.508	0.574	0.504
GAT	BENDR	0.500	<u>0.618</u>	<u>0.551</u>
	Linear	0.509	0.578	0.496
GAT (with edge weights)	BENDR	0.509	<u>0.620</u>	<u>0.525</u>
	Linear	0.508	0.577	0.500

TABLE II

THE RESULTS FOR ALL DOWNSTREAM TASKS AND GNN ARCHITECTURES USING *a linear layer* AS THE SEQUENCE LENGTH ADJUSTER.

and GEFM. Specifically, we evaluate GEFM configured as ‘GCN with edge weights’, identified as the most suitable architecture in the previous section, with the base model configuration set to ‘BENDR’ and the sequence length adjuster implemented as *a linear layer*. For a fair comparison across tasks, we use the same configurations for the baseline. Hereafter, we refer to ‘GEFM configured as ‘GCN with edge weights’’ simply as ‘GEFM’.

As shown in Table II, on MMI, the baseline achieved a score of 0.646, while GEFM achieved 0.849, representing a 31.4% improvement. On P300, the baseline achieved 0.568, with GEFM improving performance by 8.53%. On ERN, the baseline achieved 0.522, and GEFM showed a 3.11% improvement. These results indicate that the higher the baseline performance on a task, the greater the relative improvement achieved by GEFM. Further analysis could investigate the relationship between model performance and task-specific characteristics, particularly the physiological features associated with each task, to gain deeper insights.

3) *Comparison of Sequence Length Adjusters (RQ3)*: Tables II, III show that all the models using a linear layer as the sequence length adjuster consistently outperformed those using padding. Consequently, more models outperformed the baseline when using a linear layer compared to padding. This indicates that adding a linear layer before GNNs is a more effective method for sequence length adjustment.

One possible explanation for this observation lies in the difference between the sequence lengths before and after adjustment. The sequence length before adjustment was less than half, and in the smallest cases as small as one-tenth, of that after adjustment. When using padding, a significant proportion of the adjusted sequence consisted of newly added

Model	Config.	MMI	P300	ERN
Baseline	BENDR	0.646	0.568	0.522
	Linear	0.794	0.608	0.508
GraphSAGE	BENDR	<u>0.874</u>	0.512	0.468
	Linear	0.538	0.503	0.475
GCN (with edge weights)	BENDR	0.521	0.504	0.481
	Linear	0.505	0.508	0.471
GAT (with edge weights)	BENDR	0.506	0.504	0.495
	Linear	0.502	0.493	0.490

TABLE III

THE RESULTS FOR ALL DOWNSTREAM TASKS AND THE GNN ARCHITECTURES WE EXPERIMENTED WITH USING *padding* AS THE SEQUENCE LENGTH ADJUSTER.

padding values, overshadowing the meaningful information from the original signals. Consequently, the model struggled to learn effective representations, leading to lower performance.

In contrast, when using a linear layer for sequence length adjustment, the original signal was distributed more sparsely across the adjusted sequence. While the data was ‘stretched’, the essential characteristics of the signals were preserved throughout the sequence. This allowed the model to capture the critical features of the original data more effectively, enabling the GNNs to leverage these features and achieve better performance.

4) *Comparison of Base Model Configurations (RQ4)*: This section discusses whether incorporating GNNs is more effective when the base model configuration is ‘BENDR’ or ‘Linear.’ As shown in Tables II and III, when comparing the performance of models employing GNNs across all task-model combinations, the ‘BENDR’ configuration outperformed the ‘Linear’ configuration in 21 out of 24 cases. In contrast, among the baseline models, ‘Linear’ exceeded ‘BENDR’ in 2 out of 3 cases. Notably, all GNN-based models that outperformed the baseline were configured with ‘BENDR.’ These results indicate that incorporating GNNs is more effective when the base model configuration is ‘BENDR.’

A possible explanation for this observation is as follows: Incorporating GNNs introduces additional information, such as inter-channel relationships, into the feature representation derived from raw signals. To process and utilize this enriched information effectively, the model requires a greater number of parameters after the GNN layers. Compared to ‘Linear’, the ‘BENDR’ configuration includes more parameters and, more importantly, employs a Transformer Encoder, which is a powerful mechanism for feature extraction. We hypothesize that these factors enable the model to better leverage the

information encoded by the GNNs, resulting in improved performance.

V. CONCLUSION

Foundation models for EEG analysis are particularly valuable due to the difficulty of collecting large amounts of labeled data and their ability to be applied to a wide range of EEG tasks with minimal computational time and cost. In this study, we propose Graph-Enhanced EEG Foundation Model (GEFM), a novel foundation model for EEG that leverages both inter-channel relationships and the temporal dynamics of EEG signals. The proposed architecture integrates GNNs, which are effective at learning relationships between entities, with a masked autoencoder-based framework. We evaluated the model using several GNN architectures across three downstream tasks. The results indicate that GEFM, when employing the GCN architecture [7] with specific configurations, consistently outperformed the baseline across all tasks. These findings demonstrate that incorporating inter-channel relationships learning through GNNs enhances the model's performance, establishing it as a more effective foundation model for EEG analysis.

As future work, we plan to evaluate GEFM on a broader range of tasks to further demonstrate its versatility. Additionally, we aim to investigate the mechanisms underlying the observed improvements achieved through integrating graph structures, using techniques such as GNNExplainer [26]. And another potential direction is to expand our approach, which integrates inter-channel relationships and time-series information, to other base models or various self-supervised learning approaches. Alternatively, we may design a novel architecture specifically optimized for EEG foundation models.

VI. ACKNOWLEDGMENTS

This work is partially supported by JSPS KAKENHI Grant JP21K17749 and JP23K28098.

REFERENCES

- [1] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [2] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.
- [3] W. Wang, J. Dai, Z. Chen, Z. Huang, Z. Li, X. Zhu, X. Hu, T. Lu, L. Lu, H. Li, *et al.*, "Internimage: Exploring large-scale vision foundation models with deformable convolutions," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 14408–14419, 2023.
- [4] L. Yuan, D. Chen, Y.-L. Chen, N. Codella, X. Dai, J. Gao, H. Hu, X. Huang, B. Li, C. Li, *et al.*, "Florence: A new foundation model for computer vision," *arXiv preprint arXiv:2111.11432*, 2021.
- [5] D. Kostas, S. Aroca-Ouellette, and F. Rudzicz, "Bendr: Using transformers and a contrastive self-supervised learning task to learn from massive amounts of eeg data," *Frontiers in Human Neuroscience*, vol. 15, p. 653659, 2021.
- [6] W. Cui, W. Jeong, P. Thölke, T. Medani, K. Jerbi, A. A. Joshi, and R. M. Leahy, "Neuro-GPT: Towards A Foundation Model For EEG," in *2024 IEEE International Symposium on Biomedical Imaging (ISBI)*, pp. 1–5, IEEE, 2024.
- [7] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [8] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [9] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," *Advances in neural information processing systems*, vol. 30, 2017.
- [10] J. Devlin, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [11] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in neural information processing systems*, vol. 33, pp. 12449–12460, 2020.
- [12] H.-Y. S. Chien, H. Goh, C. M. Sandino, and J. Y. Cheng, "Maeg: Masked auto-encoder for eeg representation learning," in *NeurIPS 2022 Workshop on Learning from Time Series for Health*, 2022.
- [13] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [14] D. Klepl, M. Wu, and F. He, "Graph neural network-based eeg classification: A survey," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 2024.
- [15] N. Wagh and Y. Varatharajah, "Eeg-gcnn: Augmenting electroencephalogram-based neurological disease diagnosis using a domain-guided graph convolutional neural network," in *Proceedings of the Machine Learning for Health NeurIPS Workshop* (E. Alsentzer, M. B. A. McDermott, F. Falck, S. K. Sarkar, S. Roy, and S. L. Hyland, eds.), vol. 136 of *Proceedings of Machine Learning Research*, pp. 367–378, PMLR, 11 Dec 2020.
- [16] W. Liu, P. Zhou, Z. Zhao, Z. Wang, Q. Ju, H. Deng, and P. Wang, "K-BERT: Enabling Language Representation with Knowledge Graph," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 2901–2908, 2020.
- [17] T. Sun, Y. Shao, X. Qiu, Q. Guo, Y. Hu, X.-J. Huang, and Z. Zhang, "CoLAKE: Contextualized Language and Knowledge Embedding," in *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 3660–3670, 2020.
- [18] M. Yasunaga, H. Ren, A. Bosselut, P. Liang, and J. Leskovec, "QA-GNN: Reasoning with Language Models and Knowledge Graphs for Question Answering," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 535–546, 2021.
- [19] X. Zhang, A. Bosselut, M. Yasunaga, H. Ren, P. Liang, C. D. Manning, and J. Leskovec, "GreaseLM: Graph REASoning enhanced language models," in *International Conference on Learning Representations*, 2022.
- [20] I. Obeid and J. Picone, "The temple university hospital eeg data corpus," *Frontiers in Neuroscience*, vol. 10, p. 196, 2016.
- [21] A. L. Goldberger, L. A. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, and H. E. Stanley, "Physiobank, physiotoolkit, and physionet: Components of a new research resource for complex physiologic signals," *circulation*, vol. 101, no. 23, pp. e215–e220, 2000.
- [22] G. Schalk, D. J. McFarland, T. Hinterberger, N. Birbaumer, and J. R. Wolpaw, "Bci2000: a general-purpose brain-computer interface (bci) system," *IEEE Transactions on biomedical engineering*, vol. 51, no. 6, pp. 1034–1043, 2004.
- [23] L. Citi, R. Poli, and C. Cinel, "Documenting, modelling and exploiting p300 amplitude changes due to variable target delays in donchin's speller," *Journal of Neural Engineering*, vol. 7, no. 5, p. 056006, 2010.
- [24] P. Margaux, M. Emmanuel, D. Sébastien, B. Olivier, and M. Jérémie, "Objective and subjective evaluation of online error correction during p300-based spelling," *Advances in Human-Computer Interaction*, vol. 2012, no. 1, p. 578295, 2012.
- [25] V. Jurcak, D. Tsuzuki, and I. Dan, "10/20, 10/10, and 10/5 systems revisited: Their validity as relative head-surface-based positioning systems," *Neuroimage*, vol. 34, no. 4, pp. 1600–1611, 2007.
- [26] Z. Ying, D. Bourgeois, J. You, M. Zitnik, and J. Leskovec, "Gnnexplainer: Generating explanations for graph neural networks," *Advances in neural information processing systems*, vol. 32, 2019.