

MTEEG: A MULTI-TASK LEARNING FRAMEWORK FOR ENHANCED ELECTROENCEPHALOGRAPHY ANALYSIS USING LOW-RANK ADAPTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Electroencephalography (EEG) analysis using deep learning has traditionally placed a strong emphasis on models that are custom-built and optimized for specific datasets. Several recent research utilize self-supervised learning to extract generic representations from massive amounts of unlabeled EEG data. The pre-trained models are then fine-tuned on each downstream dataset independently, demonstrating promising results. However, in practical applications involving multiple tasks, utilizing a separate model for each is not ideal regarding computational and spatial cost. In this study, we go one step further and explore the simultaneous adaptation of a pre-trained model to multiple different tasks. The EEG signals exhibit significant heterogeneity due to their collection from various subjects using diverse devices and experimental setups, resulting in potential conflicts among different tasks that impede joint optimization. To tackle this challenge, we propose MTEEG, a multi-task EEG recognition framework which incorporates a task-agnostic temporal encoder and task-specific low-rank adaptation modules to disentangle the parameter space, facilitating both task interaction and specification. Experiments show that MTEEG surpasses other multi-task methods and performs on par with state-of-the-art single-task methods on abnormal detection, event type classification, emotion recognition, seizure detection, sleep stage classification and motor imagery classification after being tuned jointly on six publicly available datasets. MTEEG shows the potential of multi-task EEG recognition and promotes the development of general-purpose brain-computer interfaces in the future. The source code will be released.

1 INTRODUCTION

Electroencephalography (EEG) is a widely used neuroimaging technique that captures electrical activity of the brain through non-invasive scalp electrodes. In recent years, deep learning models, such as convolutional neural networks (CNNs) and transformers, have demonstrated remarkable success in extracting meaningful patterns from EEG data, leading to significant improvements in various applications including emotion recognition (Li et al., 2022b), motor imagery classification (Li et al., 2022b) and seizure detection (Boonyakitanont et al., 2020). However, despite their power, these models are typically customized for specific tasks and input formats, which causes them to overfit and become ungeneralizable.

Drawing inspirations from the advancements of large language models (Devlin, 2018; Achiam et al., 2023), some researchers (Yang et al., 2023a; Yi et al., 2024; Jiang et al., 2024) employ self-supervised learning to extract generic representations from large amounts of unlabeled EEG data, significantly improving the model’s generalizability. Despite their remarkable performance, these models necessitate individual fine-tuning for each downstream dataset, thereby constraining their versatility and applicability in practical scenarios involving multiple tasks. For example, an EEG-based health monitoring system may need to perform and switch between seizure detection, emotion recognition and sleep stage classification per demand to have a comprehensive evaluation of the patient’s condition, both physically and mentally. In this case, a pre-trained model must be replicated and fine-tuned three times, once for each task, resulting in significant computational and spatial

054 overhead. Therefore, it would be beneficial to have a unified system that is capable of handling
 055 different tasks simultaneously.
 056

057 Despite the promise, challenges persist to build an efficient multi-task model for EEG processing.
 058 The EEG signals, collected from various subjects utilizing different devices and experimental con-
 059 figurations, exhibit markedly distinct intrinsic characteristics. This variability can mislead the model
 060 with conflicting parameter update directions, leading to a substantial decrease in learning efficacy.
 061 Similar heterogeneity-induced issues have also been noted in other domains (Yu et al., 2020; Zhou
 062 et al., 2024b), and many methods have been proposed to tackle them; some incorporate separate
 063 modules for specific tasks (Liu et al., 2022b; Mahabadi et al., 2021), while others use soft-gating
 064 mechanisms to flexibly assign modules for different tasks (Ma et al., 2018; Cheng et al., 2016). Nev-
 065 ertheless, the majority of these studies focus on the analysis of image, text and audio data, raising
 doubts about the applicability of their findings to EEG.

066 In this study, we propose MTEEG, a novel
 067 EEG recognition framework which exploits a
 068 pre-trained LaBraM (Jiang et al., 2024) along
 069 with task-specific modules to facilitate effi-
 070 cient multi-task joint training. It consists of
 071 three major components: 1) a temporal en-
 072 coder that’s shared across all the tasks; 2)
 073 a transformer encoder with a frozen shared
 074 backbone and multiple task-specific low-rank
 075 adapters; 3) task-specific classification heads
 076 that output the final predictions. During train-
 077 ing, the task-agnostic temporal encoder pro-
 078 motes interaction among different tasks and the
 079 reuse of global knowledge, whereas the trans-
 080 former encoder allocates specialized low-rank
 081 adapters to each task, explicitly isolating the
 082 parameters. Thus, the disentanglement of task-
 083 specific knowledge towards their correspond-
 084 ing adapters effectively reduces conflicts aris-
 085 ing from heterogeneity. Furthermore, since
 086 the task-specific modules are implemented with
 087 low-rank adapters, the computational and spa-
 088 tial overhead they incur is significantly lower
 than that of fully fine-tuning a pre-trained model. In summary, our contributions are as follows:

- 089 • We investigate multi-task EEG recognition, which is a crucial yet underexplored aspect
 090 in the practical application of brain-computer interfaces. Concurring with prior research
 091 on other data types, we observe that joint training on heterogeneous EEG datasets also
 092 presents the issue of conflicts between different tasks, leading to substantial performance
 093 deterioration of the model.
- 094 • We present the MTEEG framework, which enhances a pre-trained model by incorporating
 095 task-specific modules to achieve parameter isolation across different tasks. This isolation
 096 allows for the separation of gradients to prevent conflicts, hence facilitating efficient multi-
 097 task joint training.
- 098 • Through extensive experiments, we demonstrate that after joint optimization on six pub-
 099 licly available datasets, MTEEG can handle abnormal detection, event type classification,
 100 emotion recognition, seizure detection, sleep stage classification and motor imagery simul-
 101 taneously, achieving performance superior than other multi-task methods and on par with
 102 state-of-the-art single-task methods.

104 **2 RELATED WORK**

105 **Self-supervised EEG pre-training.** Despite the scarcity of annotated EEG data, there is a substan-
 106 tial volume of unlabeled EEG data collected from various sources. Consequently, there has been
 107 a growing interest in adopting self-supervised methods to learn generic representations from these

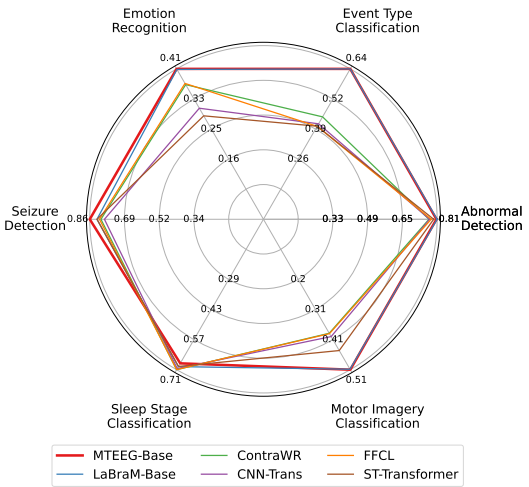


Figure 1: Overview of MTEEG’s performance (balanced accuracy) on downstream datasets.

unlabeled data to improve the model’s performance and generalizability. BENDR (Kostas et al., 2021) utilizes a contrastive learning model, wav2vec 2.0 (Baevski et al., 2020), to learn compressed representations of raw EEG signals. Neuro-GPT (Cui et al., 2024) masks random parts of the input and lets the model learn to reproduce the original signal. Brant-2 incorporates both mask-prediction and forecasting pretext tasks to enhance the model’s robustness and scalability. EEG2Rep (Mohammadi Foumani et al., 2024) reconstructs the masked samples in an abstract representation space to enhance the semantic quality of EEG representations. MMM (Yi et al., 2024) spatially divides the scalp into 17 regions and allocate a learnable token to each of them, enabling a unified topology for cross-dataset pre-training. LaBraM (Jiang et al., 2024) learns common spatial embeddings based on the 10-20 international system to be compatible with different electrode configurations, and adopts a two-stage pre-training paradigm to facilitate representation learning from noisy EEG signals.

Multi-task learning. Multi-task learning (MTL) aims to develop a model capable of handling various tasks simultaneously. The existing methods for MTL differ in how and where different tasks interact with each other. Hard parameter sharing (HPS) methods (Long et al., 2017; Lu et al., 2017) employ a single encoder for all tasks, resulting in exceptional scalability but limitations in their ability to deal with the conflicts between different tasks. The cross-stitch network (Misra et al., 2016) introduces a sharing unit to linearly combine the activation values at each layer. MTAN (Liu et al., 2019) uses attention modules to compute attention masks, thereby controlling the parameters involved in processing each task. MMoE (Ma et al., 2018) proposes to share multiple experts among different tasks with weights computed by task-specific gates, thus enabling the model to automatically learn how to balance the experts given specific inputs. PLE (Tang et al., 2020) explicitly divides experts into shared and task-specific ones, further improving the model’s robustness. In addition to the aforementioned methods that specifically target image processing, the concept of MTL has also been incorporated into EEG analysis. MIN2Net (Autthasan et al., 2021) and ERPENet (Ditthapron et al., 2019) utilize multi-task autoencoder to achieve good performance on motor imagery and P300 classification, respectively. GMSS (Li et al., 2022c) constructs different pretext tasks for a graph-based self-supervised learning model to reduce the chance of overfitting. These methods are fundamentally different from MTEEG in that they hand-craft tasks to serve for better optimization on a single dataset, while MTEEG is designed to be jointly optimized on heterogeneous datasets.

Low-rank adaptation. Low-Rank Adaptation (LoRA) (Hu et al., 2021) is a parameter-efficient fine-tuning method, which aims at reducing space and computation cost without sacrificing the model’s expressiveness. It has been widely used for adapting large foundation models to specific domains (Zhang et al., 2023; Zhou et al., 2024a). In the context of MTL, LoRA has also shown great potential because of its high level of flexibility. LoraHub (Huang et al., 2023) combines multiple LoRA modules to enhance cross-task generalization in few-shot scenarios. MOELoRA (Liu et al., 2023) integrates LoRA into a Mixture-of-Experts (MOE) framework and demonstrates superior performance. LoRAMOE (Dou et al., 2024) utilizes LoRA as an MOE-style plugin to alleviate the world knowledge forgetting problem in large language models. MoLA (Zhou et al., 2024b) includes LoRA during the training procedure and verifies their method on multiple types of heterogeneous data. However, unlike MTEEG which targets a cross-dataset setting, these methods are still limited to tasks within the same dataset.

3 METHOD

3.1 PROBLEM FORMULATION

Assume there are a total of P datasets. For $p \in \{1, 2, \dots, P\}$, given any multi-channel EEG signal $X \in \mathbb{R}^{C_p \times T_p}$ in the p -th dataset, where C_p and T_p represent the number of channels and the input duration respectively, the model aims to predict the corresponding label $y \in \mathcal{Y}_p$, where \mathcal{Y}_p represents the set of all possible outputs.

3.2 MODEL ARCHITECTURE

The architecture of MTEEG is built upon that of LaBraM. An input EEG sample $X \in \mathbb{R}^{C_p \times T_p}$ is first segmented in the temporal dimension with a non-overlapping window of length w , resulting in patches $\mathbf{x} = \{x_{i,j} | i = 1, 2, \dots, C_p, j = 1, 2, \dots, \lfloor \frac{T_p}{w} \rfloor\}$. The patches are then processed se-

162 quentially by the temporal encoder, transformer encoder and classification head to produce the final
163 output.

164 **Temporal Encoder.** The temporal encoder takes the segmented input patches and encode them
165 into embeddings, serving to capture the intricate temporal features in the signal. It consists of
166 multiple temporal convolution blocks, each of which is composed of a 1-D convolution layer, a
167 group normalization layer, and a GELU activation function. Formally, given a set of input patches
168 \mathbf{x} from dataset p , the output can be denoted as

$$169 \{e_{i,j} = TE(x_{i,j}) \in \mathbb{R}^d | x_{i,j} \in \mathbf{x}, i = 1, 2, \dots, C_k, j = 1, 2, \dots, \lfloor \frac{T_p}{w} \rfloor\},$$

170 where TE represents the temporal encoder and d is the dimension of the embeddings.
171

172 **Transformer Encoder.** To take account of the global features in the signal, we add the patch
173 embeddings with temporal and spatial embeddings based on the 10-20 international system, then
174 feed them into the transformer encoder to be processed with the attention mechanism. The attention
175 function can be formulated as

$$176 \text{Attention}(Q, K, V) = \text{softmax}\left(\frac{\text{LN}(Q)\text{LN}(K)^T}{\sqrt{d_p}}\right)V,$$

177 where d_p is the dimension of the key and query, and LN stands for layer normalization, which are
178 added to stabilize training by avoiding overly large values in the attention logits.

179 Following common practice, we employ multi-head attention to let the model attend to information
180 from different representational subspaces:

$$181 \text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$182 \text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

183 where h is the number of heads, $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$, $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$,
184 $W^O \in \mathbb{R}^{h d_v \times d_{\text{model}}}$ are the linear projection matrices.

185 3.3 TRAINING PROCEDURE

186 The training of MTEEG entails a two-stage process. In the first stage, a LaBraM model is pre-
187 trained on unlabeled data to provide a solid foundation for extracting useful information raw EEG
188 signals. Specifically, we start by training a neural tokenizer which is inspired by VQ-VAE (Van
189 Den Oord et al., 2017). The tokenizer employs the architecture outlined in Section 3.2 and is fol-
190 lowed by a neural codebook which quantizes the continuous representations into discrete tokens.
191 The learning process is then guided by the reconstruction of the amplitude and phase from these
192 discrete tokens. After the tokenizer is sufficiently trained, we train the LaBraM model by randomly
193 masking a proportion of the input patches and letting the model predict their corresponding indices
194 in the codebook. Some technical details are omitted here since the pre-training stage is not the main
195 focus of this work.

196 In the second stage, the pre-trained model is adapted to downstream datasets via a fine-tuning pro-
197 cess, in which we incorporate two major designs. Firstly, the parameters of the temporal encoder
198 are shared across and updated by all the tasks to promote the reuse of global knowledge. Secondly,
199 in the transformer encoder, we allocate specialized low-rank adapters to each task to achieve param-
200 eter isolation. An overview of the fine-tuning stage is shown in Figure 2. For any linear layer f
201 with weight matrix $W_0 \in \mathbb{R}^{m \times n}$ and bias b_0 , we define a set of low-rank decomposition matrices
202 $\Delta W = \{\Delta W_p = B_p A_p | B_p \in \mathbb{R}^{m \times r}, A_p \in \mathbb{R}^{r \times n}, p = 1, 2, \dots, P\}$ where r is the rank and P
203 is the total number of tasks. When the model performs the p -th task, the corresponding adapter is
204 injected into the layer and the original linear operation is transformed into

$$205 f(x) = W_0 x + \Delta W_p x + b_0$$

$$206 = (W_0 + B_p A_p) x + b_0$$

207 We apply this transformation to the linear projections of query, key, value and output matrices, as
208 well as the fully connected feed-forward network that follows the attention layers. Formally, for task
209 p , the output of a single attention head is

$$210 \text{head}_i = \text{Attention}(Q(W_i^Q + B_{i,p}^Q A_{i,p}^Q), K(W_i^K + B_{i,p}^K A_{i,p}^K), V(W_i^V + B_{i,p}^V A_{i,p}^V))$$

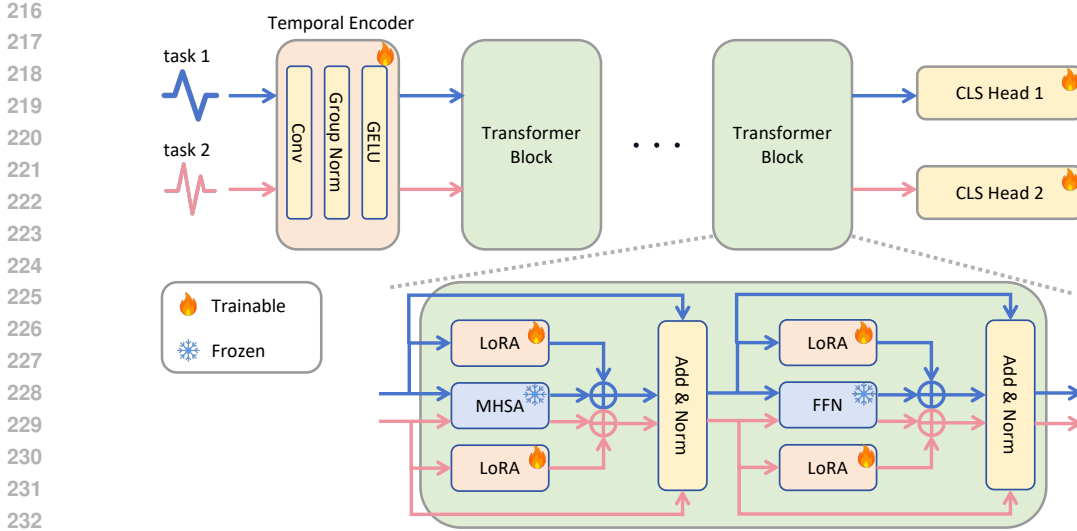


Figure 2: Overview of the fine-tuning stage. The temporal encoder, task-specific low-rank adapters and classification heads are trainable, while the pre-trained weights in the transformer encoder remain frozen.

and the full multi-head attention can be rewritten as

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)(W^O + B_p^O A_p^O)$$

where h is the number of heads, W_i^Q, W_i^K, W_i^V, W_O are the pre-trained weights for linear projections and $B_{i,p}^Q A_{i,p}^Q, B_{i,p}^K A_{i,p}^K, B_{i,p}^V A_{i,p}^V, B_p^O A_p^O$ are the corresponding task-specific low-rank adapters.

Throughout the fine-tuning stage, all the pre-trained weights in the transformer encoder are kept frozen and only the low-rank adapters are trainable. In this way, the gradients from different tasks are distinctly separated and confined within different modules, thereby alleviating the heterogeneous conflict issue.

4 EXPERIMENTS

4.1 DOWNSTREAM DATASETS

After pre-training, we fine-tune and evaluate our MTEEG jointly on the following six datasets, the statistics of which are detailed in Table 1.

TUAB (abnormal detection) (Obeid & Picone, 2016): A corpus of EEGs that have been annotated as normal or abnormal.

TUEV (event type classification) (Obeid & Picone, 2016): A subset of TUEG that contains annotations of EEG segments as one of six classes: (1) spike and sharp wave (SPSW), (2) generalized periodic epileptiform discharges (GPED), (3) periodic lateralized epileptiform discharges (PLED), (4) eye movement (EYEM), (5) artifact (ARTF) and (6) background (BCKG).

SEED-V (emotion recognition) (Liu et al., 2021): An emotion EEG dataset collected while 16 subjects watched video clips corresponding to five emotion categories (happy, sad, neutral, disgust, and fear).

CHB-MIT (seizure detection) (Shoeb, 2009): A database from Children’s Hospital Boston consisting of EEG recordings from 22 pediatric subjects with intractable seizures. Signals are sampled with 23 bipolar channels and we select the 16 standard montages in the experiments. Since the dataset is highly imbalanced (about 0.3% positive ratio), we segment the seizure regions with a 1-second stride to generate overlapping samples. In addition, we follow common practices (Lee et al., 2024; Chung et al., 2024) to randomly select 10% of the negative samples during training.

Sleep-EDF (sleep stage classification) (Goldberger et al., 2000): A database containing 197 whole-night PolySomnoGraphic sleep recordings, among which we use the 153 recordings from the study of age effects in healthy subjects (SC) in the experiments. Samples are manually annotated as one of the eight classes (W, N1, N2, N3, N4, REM, MOVEMENT, UNKNOWN). Following previous works (Supratak et al., 2017; Supratak & Guo, 2020), we exclude movement artifacts at the beginning and the end of each sleep data that was labeled as MOVEMENT or UNKNOWN, as they do not belong to the five sleep stages. In addition, we merge the N3 and N4 stages into a single stage N3 to stick to the AASM manual (Berry, 2012).

PhysioNet (motor imagery classification) (Goldberger et al., 2000): A dataset containing EEG recordings from 109 participants, with trials that belong to 5 classes: left hand, right hand, both hands, both feet, as well as rest. Following previous works (Barmpas et al., 2023; Zoumpourlis & Patras, 2024), we discard data from 6 participants (S088, S090, S092, S100, S104, S106) that have inconsistent sampling frequencies or trial lengths.

Table 1: Downstream dataset statistics

Dataset	# Channel	Sampling Rate (Hz)	Duration (seconds)	# Sample	Task
TUAB	23	256	10	409,455	Binary classification
TUEV	23	256	5	112,491	6-class classification
SEED-V	62	1000	1	148,694	5-class classification
CHB-MIT	16	256	10	26,483	Binary classification
Sleep-EDF	2	100	30	195,479	6-class classification
PhysioNet	64	160	4	18,540	5-class classification

4.2 EXPERIMENTAL SETUP

Preprocessing. We first filter the EEG signals within the range of 0.1 Hz to 75 Hz to eliminate low-frequency noise. A 50 Hz notch filter is subsequently employed to eliminate power-line interference. After that, all EEG signals are resampled to a frequency of 200 Hz. The typical range of EEG values is between -0.1 mV and 0.1 mV, which we normalize by setting the unit to 0.1 mV to ensure the values predominantly fall between -1 and 1.

Pre-training & Fine-tuning. We construct MTEEG utilizing two different configurations of LaBraM, specifically LaBraM-Base and LaBraM-Large, yielding MTEEG-Base and MTEEG-Large correspondingly. For the pre-training of LaBraM, We use the default hyperparameters outlined in the original paper. The pre-training data comprises nine public datasets, detailed in Appendix A, with a total duration of approximately 2000 hours. In the fine-tuning stage, the datasets are first split into training, validation and test subsets as outlined in Appendix B. Subsequently, we train the models using binary cross-entropy loss for binary classification tasks and cross-entropy loss for multi-class classification tasks. Due to the significantly larger data volume of TUAB compared to other datasets, which leads to early convergence and overfitting, we randomly sample 10% of the data points in TUAB for each training epoch to balance the optimization. All the experiments are conducted on Linux servers equipped with NVIDIA A100 GPUs and Python 3.10.14 + PyTorch 2.2.2 + CUDA 12.1 environment. The optimal models are trained on the training set, selected from the validation set, and finally evaluated on the test set. We report the average and standard deviation values on three different random seeds to obtain comparable results.

Baselines. For single-task baselines, we consider both self-supervised and supervised methods. Self-supervised baselines include LaBraM and BIOT (Yang et al., 2023a). Supervised baselines include SPaRCNet (Jing et al., 2023), ContraWR (Yang et al., 2021), CNN-Transformer (Peh et al., 2022), FFCL (Li et al., 2022a) and ST-Transformer (Song et al., 2021). LaBraM and BIOT are publicly accessible in their official repositories, with the supervised methods implemented by BIOT. We use the default hyperparameters for fair comparison.

Given that multi-task learning in EEG processing is underexplored and there is currently no public method for comparison, we integrate a pre-trained LaBraM-Base as the backbone network within three established multi-task learning frameworks to set up the multi-task baselines. These frame-

works include: (1) HPS (Long et al., 2017; Lu et al., 2017) where different tasks share the same expert (backbone network), except for the classification heads, (2) MMoE (Ma et al., 2018) where multiple experts are shared among different tasks with weights controlled by task-specific gates, (3) CGC (Cheng et al., 2016) where both shared and task-specific experts are included to enhance the extraction of heterogeneous features. The implementation is based on LibMTL (Lin & Zhang, 2022). Following common practice, we set the number of shared experts in MMoE and CGC to match the number of tasks, which is six in our case, and we designate one task-specific expert per task in CGC.

Metrics. We use the following metrics for evaluating the models: (1) Balanced Accuracy: the average of recall (sensitivity) on each class. (2) AUC-PR: area under the precision-recall curve, which summarizes the trade-off between precision and recall at different classification thresholds. This metric is used for binary classification. (3) AUROC: area under the receiver operating characteristic curve, which summarizes the trade-off between the true positive rate (sensitivity) and the false positive rate (1-specificity) at different classification thresholds. This metric is used for binary classification. (4) Cohen’s Kappa: an assessment of the agreement between two classifiers on a categorical scale, taking into account the possibility of agreement occurring by chance. This metric is used for multi-class classification. (5) Weighted F1: a weighted average of individual F1-scores for each class. This metric is used for multi-class classification. AUROC and Cohen’s Kappa are used as the monitoring metrics for binary and multi-class classifications respectively. For multi-task methods, we monitor the average values of these metrics across all tasks. We use PyHealth (Yang et al., 2023b) for the implementation of all the metrics.

4.3 COMPARISON WITH OTHER METHODS

The main results are summarized in Table 2, 3 and 4. The best results of multi-task and single-task methods in each column are highlighted in bold and underlined, respectively. Based on these results, we make the following observations.

Firstly, there exists a significant performance gap between HPS and LaBraM-Base across all tasks and metrics, despite their architectural similarities. This suggests that, similar to other data types, EEG signals from diverse sources can also confuse the model due to conflicting optimization directions, resulting in substantial performance degradation. Although multi-task methods such as MMoE and CGC have demonstrated efficacy in addressing this issue in other domains, their effectiveness in EEG processing remains limited. This may result from the gating mechanism in these methods being implemented with basic linear layers, which may be inadequate for differentiating the intricate intrinsic properties of highly noisy EEG signals. Secondly, in comparison to its multi-task counterparts, our proposed MTEEG-Base exhibits comparable performance on SEED-V and significantly outperforms them across all other datasets, thereby demonstrating the efficacy of gradient separation with task-specific low-rank adapters. Moreover, MTEEG even performs on par with the state-of-the-art single-task method. Comparing to LaBraM-Base, MTEEG-Base performs better on TUEV, SEED-V, CHB-MIT, and PhysioNet and slightly worse on TUAB and Sleep-EDF. The same phenomenon is also evident in the large variant of the model, confirming the scalability of our approach. Thirdly, MTEEG has the advantage of being lightweight. The base and large variants have only 1.8M and 7.4M trainable parameters fine-tuning respectively, compared to 5.8M and 46M for LaBraM-Base and LaBraM-Large. The time and space efficiency associated with this lightweight design would be beneficial in practical applications, particularly when computational resources are constrained or latency is critical.

4.4 ABLATION STUDIES

Ablation studies were performed on all six datasets; however, results are only presented for TUAB, TUEV, and SEED-V in the main paper to conserve space. For additional results on the other datasets, please refer to Appendix C.

Impact of adapter rank r . We assign different values to r , ranging from 4 to 32 to examine its impact on the model’s downstream performance. As illustrated in Figure 3, the base variant consistently achieves its maximum performance at $r = 8$ across all datasets, whereas the large variant reaches peak performance at $r = 16$ on TUAB and $r = 8$ on the remaining datasets. This indicates that a higher rank does not necessarily yield better performance, likely due to over-fitting

Impact of adapter locations. The selection of locations for applying low-rank adapters is known to significantly influence the model’s performance (Hu et al., 2021). Thus, we evaluate three different configurations of adapter locations: (1) only in multi-head self-attention modules (MHSA), (2) only in the feed-forward networks (FFN) that follow MHSA, (3) in both MHSA and FFN. As shown in Figure 4, the adaptations of both MHSA and FFN are crucial, as the elimination of either leads to a significant decline in performance.

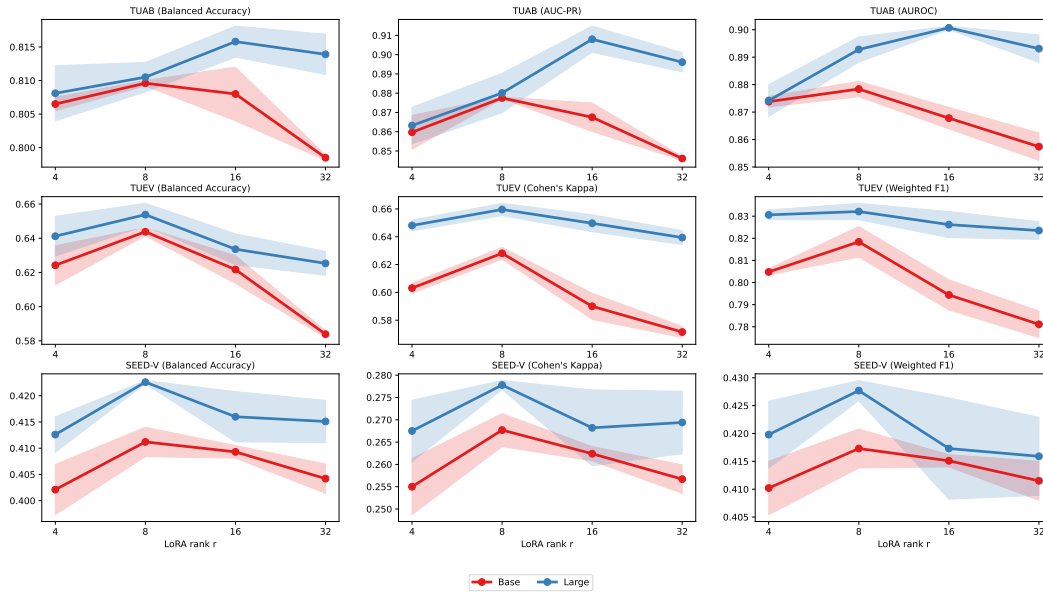


Figure 3: Ablation study on the impact of adapter rank r .

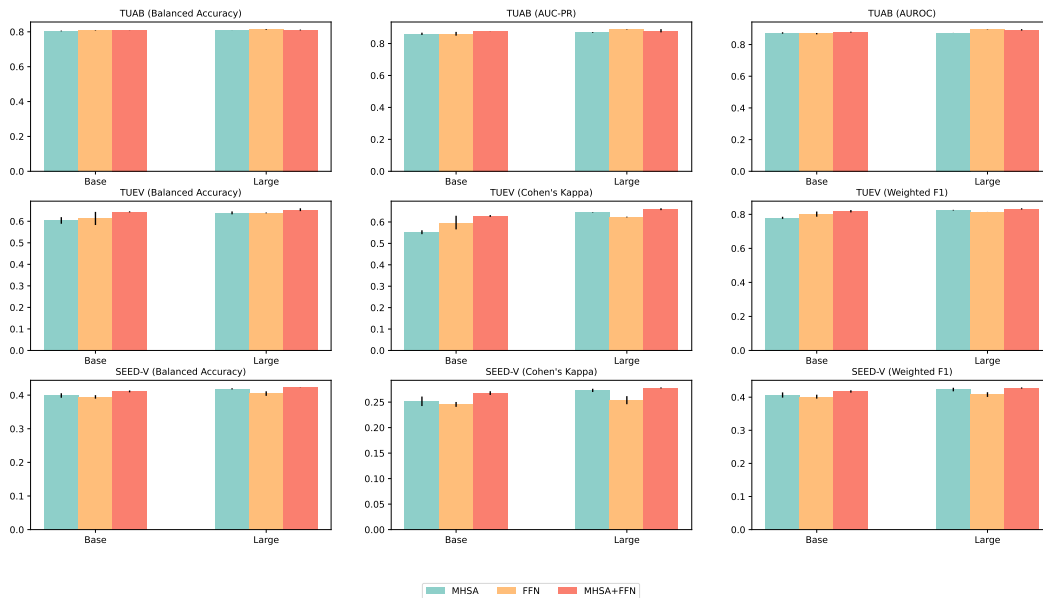


Figure 4: Ablation study on the impact of adapter locations.

Contribution of temporal encoder. The task-agnostic temporal encoder is designed to promote interaction among different tasks. To examine its actual contribution to the model’s downstream performance, we freeze it during fine-tuning and observe the resultant impact. As shown in Figure 5, freezing the temporal encoder leads to a notable decline in performance across all the tasks and metrics, with a more pronounced decrease observed in the more challenging multi-class classifica-

tion tasks. This suggests that the temporal encoder manages to capture global knowledge that helps with reducing overfitting and enhancing the generalizability of the model.

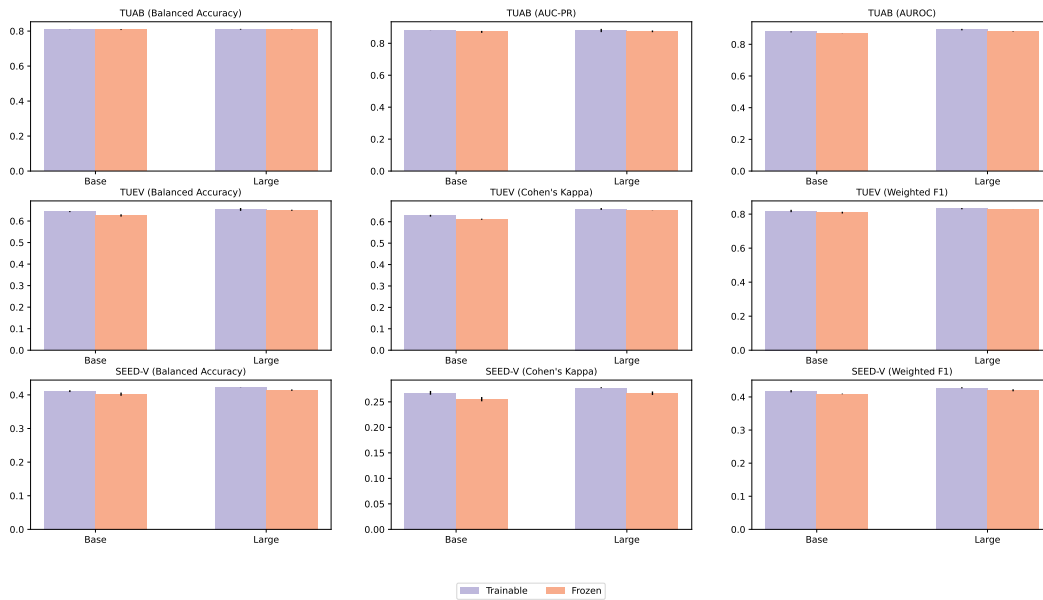


Figure 5: Ablation study on the contribution of temporal encoder.

5 CONCLUSION

This paper introduces MTEEG, an innovative multi-task EEG recognition framework. Utilizing a powerful pre-trained model, MTEEG incorporates a task-agnostic temporal encoder to capture global knowledge, along with task-specific low-rank adaptation modules to disentangle the parameter spaces for different tasks, thereby alleviating the conflicts stemming from the heterogeneity of EEG signals. We validate the effectiveness of MTEEG by fine-tuning it jointly on six publicly available datasets. Experiments show that MTEEG can simultaneously manage abnormal detection, event type classification, emotion recognition, seizure detection, sleep stage classification and motor imagery classification, outperforming other multi-task methods and matching the performance of state-of-the-art single-task methods. The adaptability and applicability of MTEEG demonstrate the significant potential of multi-task EEG recognition and promote the advancement of general-purpose brain-computer interfaces in the future.

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Phairot Autthasan, Rattanaphon Chaisaen, Thapanun Sudhawiyangkul, Phurin Rangpong, Suk-
546 tipol Kiatthaveephong, Nat Dilokthanakul, Gun Bhakdisongkhram, Huy Phan, Cuntai Guan, and
547 Theerawit Wilaiprasitporn. Min2net: End-to-end multi-task learning for subject-independent mo-
548 tor imagery eeg classification. *IEEE Transactions on Biomedical Engineering*, 69(6):2105–2118,
549 2021.
- 550 Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A frame-
551 work for self-supervised learning of speech representations. *Advances in neural information*
552 *processing systems*, 33:12449–12460, 2020.
- 553 Konstantinos Barmpas, Yannis Panagakis, Stylianos Bakas, Dimitrios A Adamos, Nikolaos
554 Laskaris, and Stefanos Zafeiriou. Improving generalization of cnn-based motor-imagery eeg
555 decoders via dynamic convolutions. *IEEE Transactions on Neural Systems and Rehabilitation*
556 *Engineering*, 31:1997–2005, 2023.
- 557 RB Berry. The aasm manual for the scoring of sleep and associated events. *Rules, Terminology and*
558 *Technical Specifications. Version, 2*, 2012.
- 559 Poomipat Boonyakitanont, Apiwat Lek-Uthai, Krisnachai Chomtho, and Jitkomut Songsiri. A re-
560 view of feature extraction and performance evaluation in epileptic seizure detection using eeg.
561 *Biomedical Signal Processing and Control*, 57:101702, 2020.
- 562 G Buckwalter, S Chhin, S Rahman, I Obeid, and J Picone. Recent advances in the tuh eeg corpus:
563 improving the interrater agreement for artifacts and epileptiform events. In *2021 IEEE Signal*
564 *Processing in Medicine and Biology Symposium (SPMB)*, pp. 1–3. IEEE, 2021.
- 565 Wei Cheng, Zhishan Guo, Xiang Zhang, and Wei Wang. Cgc: A flexible and robust approach to
566 integrating co-regularized multi-domain graph for clustering. *ACM Transactions on Knowledge*
567 *Discovery from Data (TKDD)*, 10(4):1–27, 2016.
- 568 Yoon Gi Chung, Anna Cho, Hunmin Kim, and Ki Joong Kim. Single-channel seizure detection
569 with clinical confirmation of seizure locations using chb-mit dataset. *Frontiers in Neurology*, 15:
570 1389731, 2024.
- 571 Wenhui Cui, Woojae Jeong, Philipp Thölke, Takfarinas Medani, Karim Jerbi, Anand A Joshi, and
572 Richard M Leahy. Neuro-gpt: Towards a foundation model for eeg. In *2024 IEEE International*
573 *Symposium on Biomedical Imaging (ISBI)*, pp. 1–5. IEEE, 2024.
- 574 Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding.
575 *arXiv preprint arXiv:1810.04805*, 2018.
- 576 Apiwat Dittaporn, Nannapas Banluesombatkul, Sombat Kettrat, Ekapol Chuangsuwanich, and
577 Theerawit Wilaiprasitporn. Universal joint feature extraction for p300 eeg classification using
578 multi-task autoencoder. *IEEE access*, 7:68415–68428, 2019.
- 579 Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Wei Shen, Limao Xiong, Yuhao Zhou, Xiao
580 Wang, Zhiheng Xi, Xiaoran Fan, et al. Loramoe: Alleviating world knowledge forgetting in
581 large language models via moe-style plugin. In *Proceedings of the 62nd Annual Meeting of the*
582 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1932–1945, 2024.
- 583 Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G
584 Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank,
585 physiotoolkit, and physionet: components of a new research resource for complex physiologic
586 signals. *circulation*, 101(23):e215–e220, 2000.
- 587 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
588 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
589 *arXiv:2106.09685*, 2021.

- 594 Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub: Effi-
595 cient cross-task generalization via dynamic lora composition. *arXiv preprint arXiv:2307.13269*,
596 2023.
- 597 Wei-Bang Jiang, Li-Ming Zhao, and Bao-Liang Lu. Large brain model for learning generic repre-
598 sentations with tremendous eeg data in bci. *arXiv preprint arXiv:2405.18765*, 2024.
- 600 Jin Jing, Wendong Ge, Shenda Hong, Marta Bento Fernandes, Zhen Lin, Chaoqi Yang, Sungtae An,
601 Aaron F Struck, Aline Herlopian, Ioannis Karakis, et al. Development of expert-level classifica-
602 tion of seizures and rhythmic and periodic patterns during eeg interpretation. *Neurology*, 100(17):
603 e1750–e1762, 2023.
- 604 Demetres Kostas, Stephane Aroca-Ouellette, and Frank Rudzicz. Bendr: Using transformers and a
605 contrastive self-supervised learning task to learn from massive amounts of eeg data. *Frontiers in*
606 *Human Neuroscience*, 15:653659, 2021.
- 608 Dohyun Lee, Byunghyun Kim, Taejoon Kim, Inwhee Joe, Jongwha Chong, Kyeongyuk Min, and
609 Kiyoungh Jung. A resnet- lstm hybrid model for predicting epileptic seizures using a pretrained
610 model with supervised contrastive learning. *Scientific Reports*, 14(1):1319, 2024.
- 612 Hongli Li, Man Ding, Ronghua Zhang, and Chunbo Xiu. Motor imagery eeg classification algorithm
613 based on cnn- lstm feature fusion network. *Biomedical signal processing and control*, 72:103342,
614 2022a.
- 615 Xiang Li, Yazhou Zhang, Prayag Tiwari, Dawei Song, Bin Hu, Meihong Yang, Zhigang Zhao,
616 Neeraj Kumar, and Pekka Marttinen. Eeg based emotion recognition: A tutorial and review. *ACM*
617 *Computing Surveys*, 55(4):1–57, 2022b.
- 618 Yang Li, Ji Chen, Fu Li, Boxun Fu, Hao Wu, Youshuo Ji, Yijin Zhou, Yi Niu, Guangming Shi,
619 and Wenming Zheng. Gmss: Graph-based multi-task self-supervised learning for eeg emotion
620 recognition. *IEEE Transactions on Affective Computing*, 14(3):2512–2525, 2022c.
- 622 Baijiong Lin and Yu Zhang. Libmtl: A python library for multi-task learning. *arXiv preprint*
623 *arXiv:2203.14338*, 2022.
- 624 Qidong Liu, Xian Wu, Xiangyu Zhao, Yuanshao Zhu, Derong Xu, Feng Tian, and Yefeng Zheng.
625 Moelora: An moe-based parameter efficient fine-tuning method for multi-task medical applica-
626 tions. *arXiv preprint arXiv:2310.18339*, 2023.
- 628 Shikun Liu, Edward Johns, and Andrew J Davison. End-to-end multi-task learning with attention.
629 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.
630 1871–1880, 2019.
- 632 Wei Liu, Jie-Lin Qiu, Wei-Long Zheng, and Bao-Liang Lu. Comparing recognition performance
633 and robustness of multimodal deep learning models for multimodal emotion recognition. *IEEE*
634 *Transactions on Cognitive and Developmental Systems*, 14(2):715–729, 2021.
- 635 Wei Liu, Wei-Long Zheng, Ziyi Li, Si-Yuan Wu, Lu Gan, and Bao-Liang Lu. Identifying similarities
636 and differences in emotion recognition with eeg and eye movements among chinese, german, and
637 french people. *Journal of Neural Engineering*, 19(2):026012, 2022a.
- 638 Yen-Cheng Liu, Chih-Yao Ma, Junjiao Tian, Zijian He, and Zsolt Kira. Polyhistor: Parameter-
639 efficient multi-task adaptation for dense vision tasks. *Advances in Neural Information Processing*
640 *Systems*, 35:36889–36901, 2022b.
- 642 Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Philip S Yu. Learning multiple tasks with
643 multilinear relationship networks. *Advances in neural information processing systems*, 30, 2017.
- 644 Yongxi Lu, Abhishek Kumar, Shuangfei Zhai, Yu Cheng, Tara Javidi, and Rogerio Feris. Fully-
645 adaptive feature sharing in multi-task networks with applications in person attribute classification.
646 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5334–
647 5343, 2017.

- 648 Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H Chi. Modeling task relation-
649 ships in multi-task learning with multi-gate mixture-of-experts. In *Proceedings of the 24th ACM*
650 *SIGKDD international conference on knowledge discovery & data mining*, pp. 1930–1939, 2018.
651
- 652 Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameter-
653 efficient multi-task fine-tuning for transformers via shared hypernetworks. *arXiv preprint*
654 *arXiv:2106.04489*, 2021.
- 655 Ishan Misra, Abhinav Shrivastava, Abhinav Gupta, and Martial Hebert. Cross-stitch networks for
656 multi-task learning. In *Proceedings of the IEEE conference on computer vision and pattern recog-*
657 *niton*, pp. 3994–4003, 2016.
658
- 659 Navid Mohammadi Foumani, Geoffrey Mackellar, Soheila Ghane, Saad Irtza, Nam Nguyen, and
660 Mahsa Salehi. Eeg2rep: Enhancing self-supervised eeg representation through informative
661 masked inputs. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discov-*
662 *ery and Data Mining*, pp. 5544–5555, 2024.
- 663 Iyad Obeid and Joseph Picone. The temple university hospital eeg data corpus. *Frontiers in neuro-*
664 *science*, 10:196, 2016.
665
- 666 Wei Yan Peh, Yuanyuan Yao, and Justin Dauwels. Transformer convolutional neural networks for
667 automated artifact detection in scalp eeg. In *2022 44th Annual International Conference of the*
668 *IEEE Engineering in Medicine & Biology Society (EMBC)*, pp. 3599–3602. IEEE, 2022.
- 669 Vinit Shah, Eva Von Weltin, Silvia Lopez, James Riley McHugh, Lillian Veloso, Meysam Gol-
670 mohammadi, Iyad Obeid, and Joseph Picone. The temple university hospital seizure detection
671 corpus. *Frontiers in neuroinformatics*, 12:83, 2018.
672
- 673 Ali Hossam Shoeb. *Application of machine learning to epileptic seizure onset detection and treat-*
674 *ment*. PhD thesis, Massachusetts Institute of Technology, 2009.
- 675 Yonghao Song, Xueyu Jia, Lie Yang, and Longhan Xie. Transformer-based spatial-temporal feature
676 learning for eeg decoding. *arXiv preprint arXiv:2106.11170*, 2021.
677
- 678 Akara Supratak and Yike Guo. Tinsleepnet: An efficient deep learning model for sleep stage
679 scoring based on raw single-channel eeg. In *2020 42nd Annual International Conference of the*
680 *IEEE Engineering in Medicine & Biology Society (EMBC)*, pp. 641–644. IEEE, 2020.
- 681 Akara Supratak, Hao Dong, Chao Wu, and Yike Guo. Deepsleepnet: A model for automatic sleep
682 stage scoring based on raw single-channel eeg. *IEEE transactions on neural systems and rehabil-*
683 *itation engineering*, 25(11):1998–2008, 2017.
684
- 685 Hongyan Tang, Junning Liu, Ming Zhao, and Xudong Gong. Progressive layered extraction (ple):
686 A novel multi-task learning (mtl) model for personalized recommendations. In *Proceedings of*
687 *the 14th ACM Conference on Recommender Systems*, pp. 269–278, 2020.
- 688 Logan Trujillo. Raw EEG Data. 2020. doi: 10.18738/T8/SS2NHB. URL [https://doi.org/](https://doi.org/10.18738/T8/SS2NHB)
689 [10.18738/T8/SS2NHB](https://doi.org/10.18738/T8/SS2NHB).
- 690 Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in*
691 *neural information processing systems*, 30, 2017.
692
- 693 L Veloso, J McHugh, E Von Weltin, S Lopez, I Obeid, and J Picone. Big data resources for eegs:
694 Enabling deep learning research. In *2017 IEEE Signal Processing in Medicine and Biology Sym-*
695 *posium (SPMB)*, pp. 1–3. IEEE, 2017.
- 696 Eva von Weltin, Tameem Ahsan, Vinit Shah, Dawer Jamshed, Meysam Golmohammadi, Iyad Obeid,
697 and Joseph Picone. Electroencephalographic slowing: A primary source of error in automatic
698 seizure detection. In *2017 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*,
699 pp. 1–5. IEEE, 2017.
700
- 701 Chaoqi Yang, Danica Xiao, M Brandon Westover, and Jimeng Sun. Self-supervised eeg representa-
tion learning for automatic sleep staging. *arXiv preprint arXiv:2110.15278*, 2021.

- 702 Chaoqi Yang, M Brandon Westover, and Jimeng Sun. Biot: Cross-data biosignal learning in the
703 wild. *arXiv preprint arXiv:2305.10351*, 2023a.
704
- 705 Chaoqi Yang, Zhenbang Wu, Patrick Jiang, Zhen Lin, Junyi Gao, Benjamin Danek, and Jimeng
706 Sun. PyHealth: A deep learning toolkit for healthcare predictive modeling. In *Proceedings of the*
707 *27th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*
708 *2023*, 2023b. URL <https://github.com/sunlabuiuc/PyHealth>.
- 709 Ke Yi, Yansen Wang, Kan Ren, and Dongsheng Li. Learning topology-agnostic eeg representations
710 with geometry-aware modeling. *Advances in Neural Information Processing Systems*, 36, 2024.
711
- 712 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.
713 Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems*,
714 33:5824–5836, 2020.
- 715 Longteng Zhang, Lin Zhang, Shaohuai Shi, Xiaowen Chu, and Bo Li. Lora-fa: Memory-efficient
716 low-rank adaptation for large language models fine-tuning. *arXiv preprint arXiv:2308.03303*,
717 2023.
- 718 Wei-Long Zheng and Bao-Liang Lu. Investigating critical frequency bands and channels for eeg-
719 based emotion recognition with deep neural networks. *IEEE Transactions on autonomous mental*
720 *development*, 7(3):162–175, 2015.
721
- 722 Wei-Long Zheng, Wei Liu, Yifei Lu, Bao-Liang Lu, and Andrzej Cichocki. Emotionmeter: A
723 multimodal framework for recognizing human emotions. *IEEE transactions on cybernetics*, 49
724 (3):1110–1122, 2018.
- 725 Yuhang Zhou, Haolin Li, Siyuan Du, Jiangchao Yao, Ya Zhang, and Yanfeng Wang. Low-rank
726 knowledge decomposition for medical foundation models. In *Proceedings of the IEEE/CVF Con-*
727 *ference on Computer Vision and Pattern Recognition*, pp. 11611–11620, 2024a.
728
- 729 Yuhang Zhou, Zihua Zhao, Haolin Li, Siyuan Du, Jiangchao Yao, Ya Zhang, and Yanfeng Wang.
730 Exploring training on heterogeneous data with mixture of low-rank adapters. *arXiv preprint*
731 *arXiv:2406.09679*, 2024b.
- 732 Georgios Zoumpourlis and Ioannis Patras. Motor imagery decoding using ensemble curriculum
733 learning and collaborative training. In *2024 12th International Winter Conference on Brain-*
734 *Computer Interface (BCI)*, pp. 1–8. IEEE, 2024.
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A PRE-TRAINING DATASETS

We use a selection of datasets from the original LaBraM paper, omitting the private ones, for pre-training. The overall duration is approximately 2000 hours.

Table 5: Information of datasets used for pre-training.

Dataset	#Channel	Rate (Hz)	Time (h)	Description
TUEP (Veloso et al., 2017)	19-23	256	591.22	A subset of TUEG that contains 100 subjects epilepsy and 100 subjects without epilepsy, as determined by a certified neurologist.
TUSL (von Weltin et al., 2017)	23	256	20.59	A subset of TUEG that contains annotations of slowing events.
TUSZ (Shah et al., 2018)	19-23	256	1138.53	A corpus containing EEG signals that have been manually annotated data for seizure events (start time, stop, channel and seizure type).
TUAR (Buckwalter et al., 2021)	23	256	92.22	A subset of TUEG that contains annotations of 5 different artifacts: (1) eye movement (EYEM), (2) chewing (CHEW), (3) shivering (SHIV), (4) electrode pop, electrode static, and lead artifacts (ELPP), and (5) muscle artifacts (MUSC).
SEED Series (Zheng & Lu, 2015; Zheng et al., 2018; Liu et al., 2022a)	62	1000	166.75	Emotional datasets collected when subjects watched videos. These datasets include SEED (15 subjects), SEED-IV (15 subjects), SEED-GER (8 subjects), and SEED-FRA (8 subjects).
Raw EEG Data (Trujillo, 2020)	64	256	34.35	A dataset containing EEG signals recorded during the reported Information-Integration categorization task and the reported multidimensional Rule-Based categorization task.

B ADDITIONAL DETAILS OF FINE-TUNING

B.1 DATA SPLIT

TUAB and **TUEV**: The training and test sets are provided by the original creator of the dataset. We adhere to BIOT and LaBraM to partition the training set into training and validation subsets at a ratio of 80% and 20%, respectively.

SEED-V: We divide the 15 trials of each session into three groups of five, then consolidate each group from all sessions to create the training, validation, and test sets.

CHB-MIT: There are a total of 23 cases collected from 22 subjects. Following BIOT, we use cases 1 to 19 for training, cases 20 and 21 for validation, and cases 22 and 23 for testing.

Sleep-EDF and **PhysioNet**: We partition the recordings by order into training, validation and test sets at a ratio of 64%, 16% and 20%, respectively.

B.2 HYPERPARAMETERS

Table 6: Hyperparameters for downstream fine-tuning.

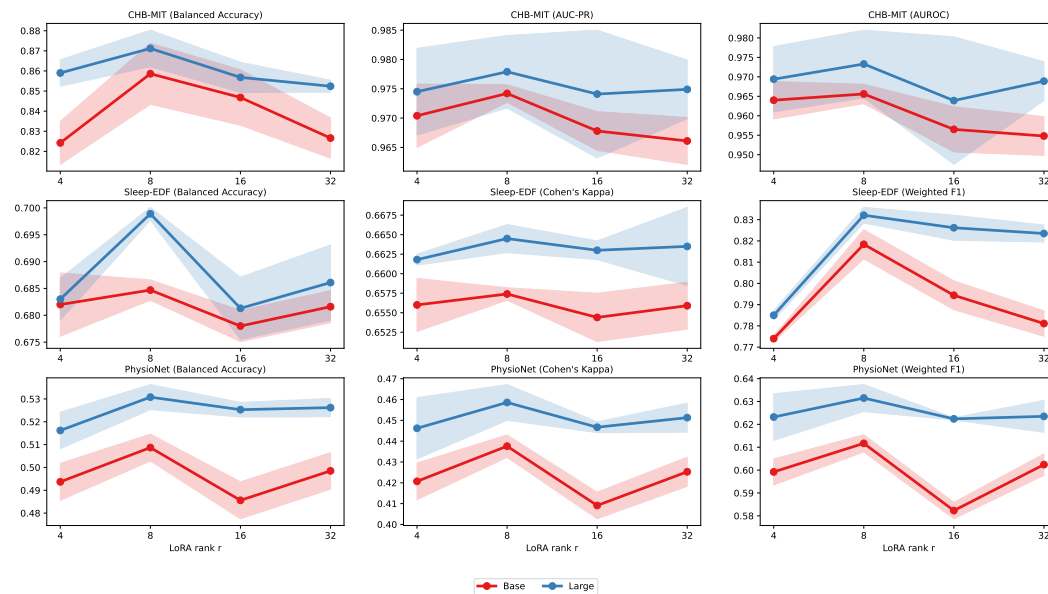
Hyperparameters	Values
Batch size	128
LoRA learning rate	5e-3
Temporal encoder learning rate	5e-4
Minimal learning rate	1e-6
Learning rate scheduler	Cosine
Optimizer	AdamW
Adam β	(0.9,0.999)
Weight decay	0.05
Total epochs	50
Warmup epochs	5
Drop path	0.1
Layer-wise learning rate decay	0.9
Label smoothing (multi-class classification)	0.1

C ADDITIONAL RESULTS OF ABLATION STUDIES

The results of ablation studies on CHB-MIT, Sleep-EDF and PhysioNet are shown in Figure 6, 7 and 8. We observe similar trends to those in Figure 3, 4 and 5, which are summarized as follows:

- MTEEG reaches peak performance when the rank of adapters is set to 8.
- Adaptations to both the MHSA and FFN modules in transformer encoder are crucial, as eliminating either of them results a significant decrease in the model’s downstream performance.
- The shared temporal encoder enables interaction between different tasks, thereby reducing overfitting and further boosting the performance.

These observations are consistent across all tasks and metrics, thereby affirming their validity.

Figure 6: Additional results of ablation study on the impact of adapter rank r .

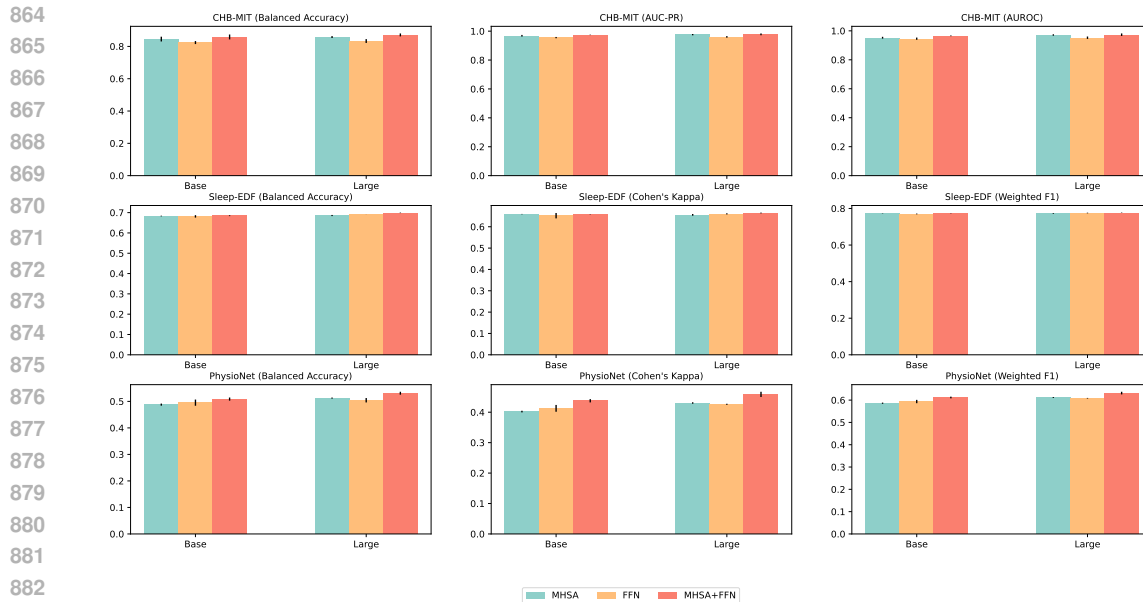


Figure 7: Additional results of ablation study on the impact of adapter locations.

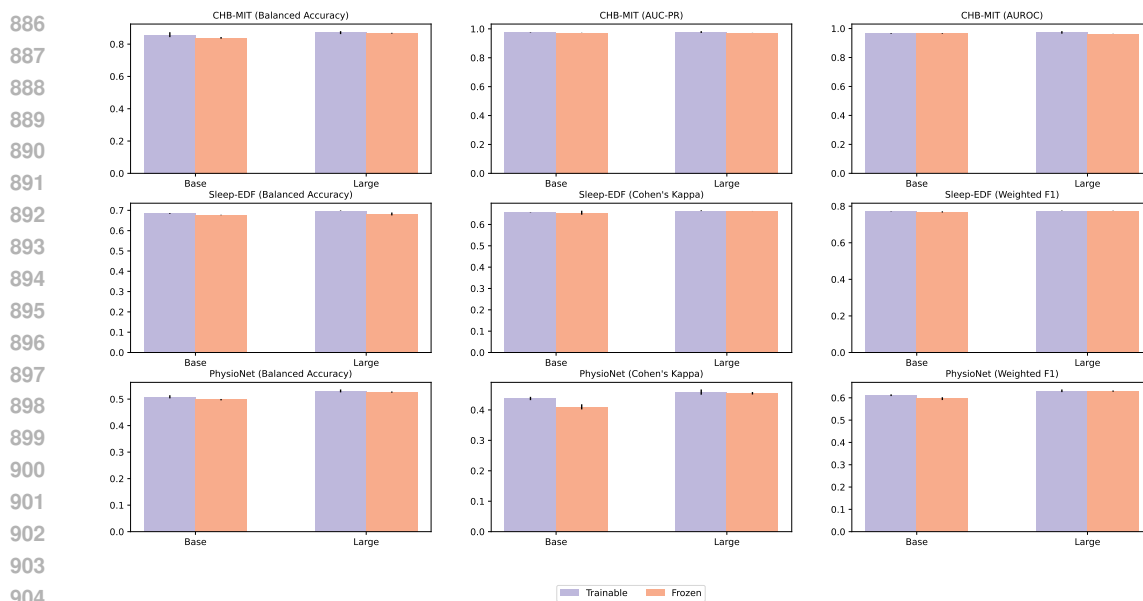


Figure 8: Additional results of ablation study on the contribution of temporal encoder.

D DISCUSSION

MTEEG represents a groundbreaking study in the joint optimization on heterogeneous EEG datasets to facilitate multi-task capability, yielding commendable results across diverse downstream tasks. Nonetheless, we note that it has the following limitations. Firstly, the representational ability of MTEEG is significantly influenced by the selection of the pre-trained model. The pre-training phase, although not the primary focus of this paper, is an essential element that establishes the upper limit of the model’s performance. Therefore, MTEEG would benefit from the future advancement of self-supervised EEG pre-training paradigms. Secondly, the EEG datasets exhibit significant variability in size and convergence speed, leading to challenges in balancing the optimization processes. In this study, we employ a rudimentary strategy to sample a subset of the data points in TUAB for each

918 training epoch, thereby decelerating convergence on this particular dataset; however, this approach
919 is suboptimal and presents significant opportunities for enhancement. Looking ahead, we believe
920 that adopting a more adaptive approach to handle the imbalance between different datasets would
921 greatly enhance multi-task joint training.
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971