

ALFEE: Adaptive Large Foundation Model for EEG Representation

Wei Xiong, Junming Lin, Jiangtong Li, Jie Li, Changjun Jiang
 Department of Computer Science and Technology, Tongji University
 Shanghai, China
 {xw1216,jiangtongli}@tongji.edu.cn

Abstract

While foundation models excel in text, image, and video domains, the critical biological signals, particularly electroencephalography (EEG), remain underexplored. EEG benefits neurological research with its high temporal resolution, operational practicality, and safety profile. However, low signal-to-noise ratio, inter-subject variability, and cross-paradigm differences hinder the generalization of current models. Existing methods often employ simplified strategies, such as a single loss function or a channel-temporal joint representation module, and suffer from a domain gap between pretraining and evaluation tasks that compromises efficiency and adaptability. To address these limitations, we propose the **Adaptive Large Foundation model for EEG signal representation (ALFEE)** framework, a novel hybrid transformer architecture with two learning stages for robust EEG representation learning. ALFEE employs a hybrid attention that separates channel-wise feature aggregation from temporal dynamics modeling, enabling robust EEG representation with variable channel configurations. A channel encoder adaptively compresses variable channel information, a temporal encoder captures task-guided evolution, and a hybrid decoder reconstructs signals in both temporal and frequency domains. During pretraining, ALFEE optimizes task prediction, channel and temporal mask reconstruction, and temporal forecasting to enhance multi-scale and multi-channel representation. During fine-tuning, a full-model adaptation with a task-specific token dictionary and a cross-attention layer boosts performance across multiple tasks. After 25,000 hours of pretraining, extensive experimental results on six downstream EEG tasks demonstrate the superior performance of ALFEE over existing models. Our ALFEE framework establishes a scalable foundation for biological signal analysis with implementation available at <https://github.com/xw1216/ALFEE>.

CCS Concepts

• **Computing methodologies** → **Artificial intelligence; Cognitive science**; • **Human-centered computing** → *HCI design and evaluation methods*.

Keywords

Foundation Model, Multi-task Learning, EEG, Pretraining

1 Introduction

The emergence of foundation models and their evolution into multimodal systems [1, 6, 13, 19, 51, 58] has transformed the artificial intelligence to process and integrate diverse information modalities. Specifically, multimodal foundation models exhibit strong

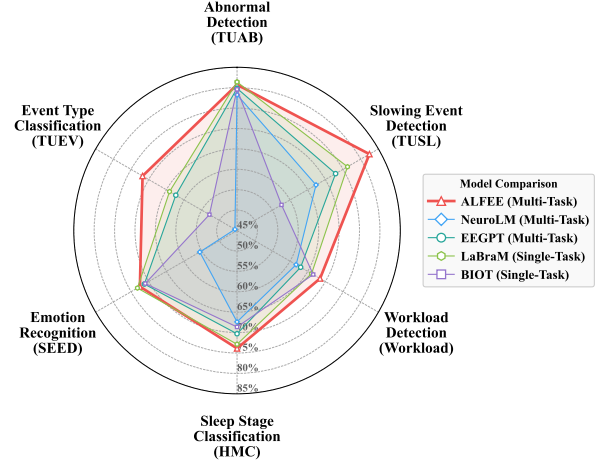


Figure 1: Balanced accuracy performance on six datasets.

capabilities in knowledge integration and data analysis while narrowing the gap between computational systems and real-world understanding, thereby driving advances in human productivity. Pioneering frameworks such as CLIP [54] and MAE [23] have set new benchmarks in visual-semantic understanding, while BEATs [10] has emerged as a paradigm-shifting approach for music audio feature extraction. Despite these achievements, current foundation models primarily focus on conventional modalities (e.g., text, image, video) [2, 15, 59, 79], with their architectural derivatives largely confined to these domains. This leaves critical biological signal modalities, i.e., electroencephalography (EEG), significantly underexplored in the era of foundation models. In particular, multi-channel time-series data like EEG offers vital insights into brain function and structure, which is essential in cognitive science.

EEG, as a widely adopted non-invasive neural recording modality, offers unique advantages for both neurological research and medical applications [11, 20, 36, 67, 80]. Specifically, EEG exhibits **high temporal resolution** (superior to Near-Infrared Spectroscopy), **operational practicality** (more cost-effective than functional Magnetic Resonance Imaging), and **safety profile** (non-invasive compared to Positron Emission Tomography). These advantages render EEG indispensable for detecting neurological anomalies, such as epileptic seizures, sleep disorders, and traumatic brain injuries [33].

Recent advances in EEG signal modeling have introduced several innovative pretraining frameworks with distinct methodological contributions. LaBraM [29] utilizes vector-quantized neural spectrum prediction to establish a comprehensive semantic tokenization system during the pretraining phase, while its neural transformer architecture facilitates simultaneous learning of spatiotemporal characteristics inherent in EEG signals. The EEGPT framework [77] introduces a dual self-supervised paradigm combining spatiotemporal representation alignment with mask-based

reconstruction. This hierarchical architecture initially derives stable spatial features from short-term EEG segments, subsequently modeling temporal dependencies across extended EEG signal sequences. Notably, NeuroLM [28] presents an integrated approach by combining the LaBraM encoder with the GPT-2 architecture [56], effectively harnessing the representational capacity of large-scale language models. This hybrid architecture demonstrates enhanced performance in multimodal learning scenarios, particularly in tasks requiring simultaneous processing of multiple cognitive objectives.

Despite significant advancements in EEG signal representation, EEG foundation models continue to encounter several critical challenges. First, variations in channel counts across different datasets pose challenges in efficiently learning channel representations, as they must adhere to varying standards within a unified framework. Second, EEG signals are characterized by two primary dimensions: the channel dimension, which reflects the activity of various brain regions at a specific moment, and the temporal dimension, which captures the evolution of these signals over time. Current methods typically rely on either a single loss function or a joint representation module, which limits their capability to fully capture the complex, dimension-specific features. Third, a significant domain gap persists between upstream pretraining tasks and downstream evaluation tasks. Existing methods either employ single-task fine-tuning for specific applications or freeze pretrained parameters while incorporating additional modules for downstream tasks, which reduces adaptability and impairs generalization and robustness.

In this work, we present the **Adaptive Large Foundation model for EEG signal representation (ALFEE)** framework, which is built on a hybrid attention architecture and employs two optimization stages to address the above issues. From the framework perspective, our ALFEE separates channel-wise feature aggregation from temporal dynamics modeling by first employing a channel encoder that adaptively compresses variable channel information to accommodate diverse channel counts. Subsequently, a temporal encoder captures task-guided temporal evolution. Finally, a decoder using hybrid attention with a learnable recovery query reconstructs the EEG signal in both the temporal domain and the frequency domain. From the training perspective, our framework optimizes four tasks during the pretraining stage: task prediction, channel mask reconstruction, temporal mask reconstruction, and temporal forecasting, effectively enhancing its multi-channel and multi-scale representation capabilities. In the fine-tuning stage, a full-model fine-tuning strategy is employed, utilizing a task-specific token dictionary augmented by a cross-attention layer, and comprehensively improving performance in multi-task scenarios.

To validate the effectiveness of our framework, we pretrain ALFEE with four model sizes on a large-scale EEG dataset, comprising over 25,000 hours and 15 datasets spanning 8 tasks [28]. Next, we finetune ALFEE with different model sizes on 6 downstream tasks, including TUAB [42], TUEV [22], TUSL [75], SEED [88], HMC [4], and Workload [89]. Extensive experimental results (Figure 1) demonstrate the superior performance of ALFEE over existing multi-task state-of-the-art (SOTA) EEG foundation models (*i.e.*, NeurLM and EEGPT), showcasing its advanced capability in robust and generalized feature representations from noisy EEG signals. Finally, further experiments confirm the contributions of each loss

and model size within the framework, validating the scaling-law in EEG signal representation. Our contributions are summarized as

- **Hybrid Attention Architecture:** An hybrid mechanism with self- and cross-attention that separates channel-wise feature aggregation from temporal dynamics modeling, enabling EEG representation with variable channel setting.
- **Multi-Task, Multi-Channel, Multi-Scale Pretraining:** A unified training framework that combines: a) temporal forecasting; b) channel masked reconstruction; c) temporal masked autoencoding; and d) task prediction, enhanced by Power Spectral Density (PSD) features for frequency domain.
- **Adaptive Representation Learning:** A modular architecture comprising a channel encoder, a temporal encoder, and an EEG decoder that adapts to heterogeneous EEG signals and enhances the capability in robust representations.
- **Extensive Experiments:** Extensive experimental results on a pretraining dataset comprising over 25,000 hours EEG signals and 6 evaluation datasets, demonstrating the superior performance of ALFEE over existing models and further validating the scaling law in EEG signal representation.

2 Methodology

In this section, we present a detailed description of the architecture and optimization of the ALFEE model. The model follows an encoder-decoder paradigm during pretraining, where an EEG feature encoder, composed of a composite feature extractor, a channel encoder, and a temporal transformer, first maps the input to a compact representation, and a decoder then reconstructs the input from this compact representation. In the finetuning stage, a finetuning head is attached to temporal encoder for downstream classification.

2.1 Preliminaries

To handle the issue that channel counts vary for different tasks and datasets, we first collect all existing electrode positions among different standards and predefine an electrode set $\mathbb{C} = c_1, c_2, \dots, c_{|\mathbb{C}|}$, where $|\mathbb{C}| = 90$, containing all standard positions in the international 10-10 system in addition to T1, T2, A1, and A2 [50]. For all involved datasets, the signals are resampled to $f_s = 256$ Hz. The learnable query is denoted as \mathbf{I}_* . The resampled EEG signal data is formulated as $\mathbf{x}_0 \in \mathbb{R}^{C_0 \times P_0}$, where C_0 denotes the number of original electrodes in the dataset and P_0 denotes the number of sample points. After slicing, the sample points are split into $n = C_0 \times \lfloor P_0 / f_s \rfloor$ non-overlapping one-second patches with patch size $P = f_s$. Subsequently, EEG patch samples are denoted as $\mathbf{x}_p \in \mathbb{R}^{B \times T \times C \times f_s}$, where B , T , C , and f_s denote the batch size, the number of time steps, the channel size, and the patch size after batching, respectively.

2.2 Model Architecture

ALFEE adopts a hierarchical architecture with the following components: 1) Feature Extractor; 2) Channel Encoder; 3) Temporal Encoder; 4) EEG Decoder; 5) Pretraining Head; and 6) Finetuning Head. As illustrated in Figure 2, ALFEE processes EEG signals with a time-frequency domain feature extractor, followed by encoder-decoder to capture EEG feature and reconstruct EEG signal.

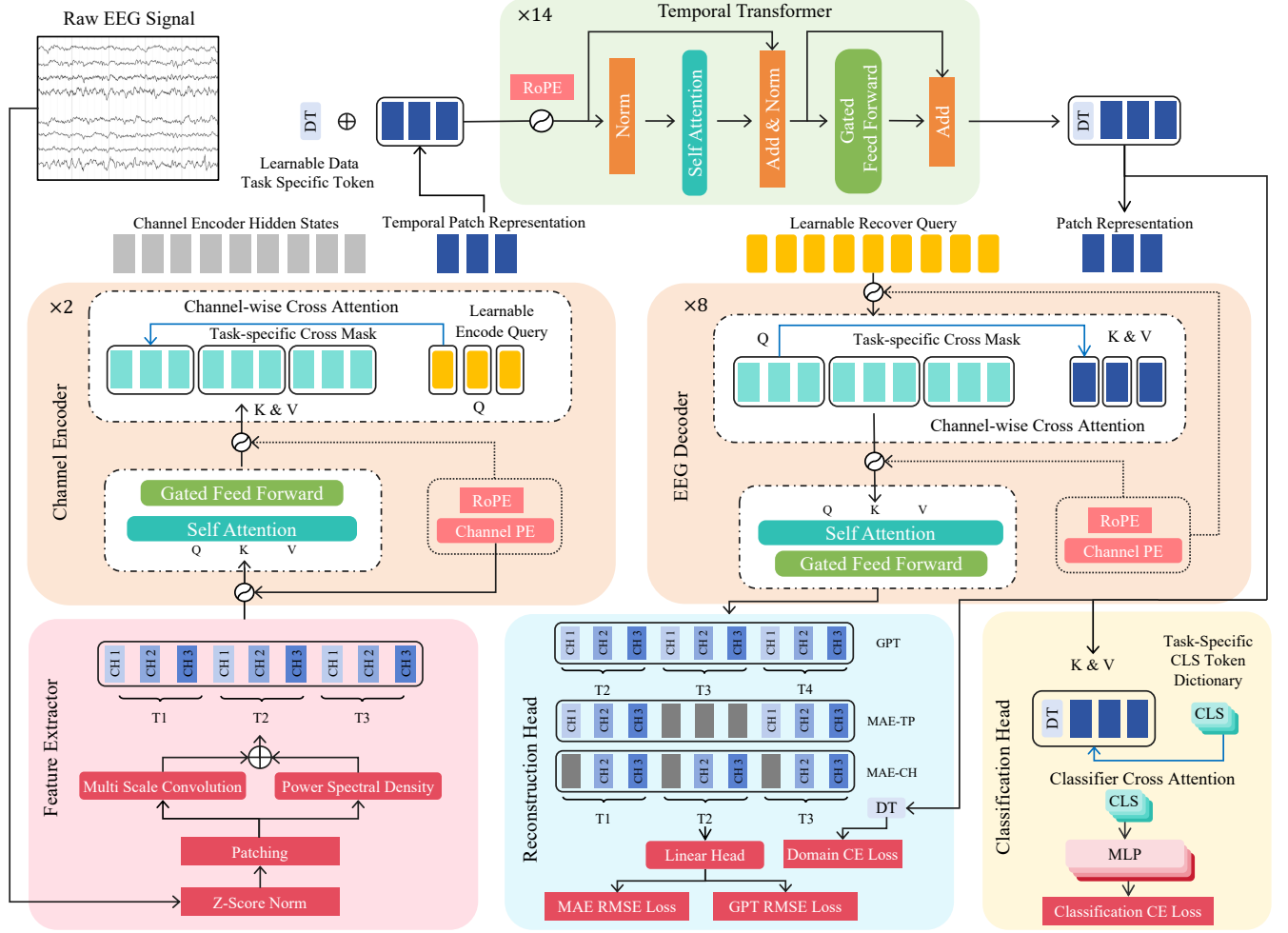


Figure 2: The overall architecture of ALFEE. (1) **Feature Extractor:** processes raw EEG signals via multi-scale convolution and PSD analysis; (2) **Channel Encoder:** captures channel-wise dependencies using cross-attention mechanisms; (3) **Temporal Encoder:** processes temporal feature through stacked transformer layers, guided by learnable task token; (4) **EEG Decoder:** reconstructs signals through attention-based refinement; (5) **Pretraining Head:** learns signal reconstruction via complicated loss minimization; and (6) **Finetuning Head:** gather information by task-specific token for downstream tasks. All the self- and cross-attention layers exploits the multi-head attention and different attention masking strategies are applied in each module.

2.2.1 Feature Extractor.

Input Normalization and Patching: To control the scale of the training loss and prevent gradient explosion in automatic mixed precision [46], we first apply z-score normalization to the resampled EEG signal data \mathbf{x}_0 along the temporal dimension:

$$\hat{\mathbf{x}}_0 = \frac{\mathbf{x}_0 - \mu_{\mathbf{x}_0}}{\sigma_{\mathbf{x}_0} + \epsilon} \quad (1)$$

where $\mu_{\mathbf{x}_0}$ and $\sigma_{\mathbf{x}_0}$ are calculated only on the temporal dimension, and $\epsilon = 1e^{-5}$ prevents division by zero. Then we divide the input samples as described in Section 2.1 and organize them into mini-batches, leading to $\mathbf{x}_p \in \mathbb{R}^{B \times T \times C \times f_s}$ for further operation.

Multi-Domain Feature Extraction: For the purpose of extracting informative EEG features in both temporal and frequency domains, we first implement multi-scale convolution using stacked Conv1D

layers [45, 69] and compute the PSD using the Fast Fourier Transform (FFT). Temporal domain features describe transient characteristics, while frequency domain features illustrate global characteristics. Given a normalized sample $\mathbf{x}_p \in \mathbb{R}^{B \times T \times C \times f_s}$, the tensor is first reshaped into $\mathbf{x}'_p \in \mathbb{R}^{(B \cdot T \cdot C) \times 1 \times f_s}$. The operation of the l -th Conv1D layer ($l \in 1, 2, \dots, L$) in the StackConv is defined as:

$$\mathbf{H}^{(l)} = \text{LayerNorm} \left(\text{GELU} \left(\text{Conv1D} \left(\mathbf{H}^{(l-1)} \right) \right) \right) \quad (2)$$

where $\mathbf{H}^{(l)}$ denotes the representation at the l -th layer with $\mathbf{H}^0 = \mathbf{x}'_p$. Then, the output ($\mathbf{H}^{(L)}$) is compressed along the patch dimension to a fixed length, followed by reshaping to restore the channel structure, leading to $\mathbf{Y} \in \mathbb{R}^{B \times T \times C \times D_{\text{Conv}}}$. The receptive field of

each StackConv is formulated by:

$$\mathcal{R}_l = \mathcal{R}_{l-1} + \left(K^l - 1\right) \cdot \prod_{i=1}^{l-1} s^i \quad (3)$$

where K^l and s^l are the kernel size and stride at the l -th layer, and \mathcal{R} is receptive field with $\mathcal{R}_0 = 1$. The multi-scale embedding module aggregates features from S independent StackConv blocks. Let $\mathbf{Y}_s \in \mathbb{R}^{B \times T \times C \times D_{Conv}}$ denote the output of the s -th stack.

$$\mathbf{x}_t = \text{LayerNorm} \left(\text{GELU} \left(\text{Concat}(\mathbf{Y}_1, \dots, \mathbf{Y}_S) \cdot \mathbf{W}_{\text{emb}}^\top \right) \right), \quad (4)$$

where $\mathbf{W}_{\text{emb}} \in \mathbb{R}^{(S \times D_{Conv}) \times D_{ti}}$ maps the concatenated features to the target temporal dimension $D_{ti} = 512$ for our base model.

Meanwhile, the frequency domain feature is extracted by $\mathbf{x}_{fq} \in \mathbb{C}^{B \times T \times C \times D_{fq}} = \mathcal{F}(\mathbf{x}_p)$, where \mathcal{F} is the FFT operator. Here $D_{fq} = \lfloor f_s \div 2 \rfloor + 1$ due to the Nyquist–Shannon sampling theorem [26], representing the frequency bins in one-sided spectrum. The single-sided PSD is formulated as $\mathbf{x}_{\text{PSD}} = |\mathbf{x}_{fq}|^2$ and normalized as

$$\mathbf{x}_{\text{PSD}}^{\text{norm}} = \frac{\mathbf{x}_{\text{PSD}}}{\sum_{p=0}^{P-1} w_{\text{Hann}}^2[p]} \quad (5)$$

where the Hanning window is defined as:

$$w_{\text{Hann}}[p] = 0.5 \left(1 - \cos \left(\frac{2\pi p}{P-1} \right) \right) \quad (6)$$

where P is patch size. Finally, $\mathbf{x}_{\text{PSD}}^{\text{norm}}$ is converted to decibels (dB):

$$\mathbf{x}_s = \log_{10}(\mathbf{x}_{\text{PSD}}^{\text{norm}}). \quad (7)$$

With z-score normalization, the extracted features is formulated as:

$$\mathbf{x} = \text{Linear}(\text{Z-Norm}(\mathbf{x}_t)) \oplus \text{Linear}(\text{Z-Norm}(\mathbf{x}_s)), \quad (8)$$

where \oplus denotes concatenation and the embedding dimension is $D = D_{ti} + D_{fq}$. The reference sample for reconstruction, combining temporal and frequency domain features, is generated as

$$\mathbf{x}_{\text{ref}} = \mathbf{x}_p \oplus \mathbf{x}_s. \quad (9)$$

Positional Encoding: The channel-wise positional information are encoded via learnable embeddings:

$$\mathbf{E}_c = \mathbf{W}_{pos}[c_1, c_2, \dots, c_{|C|}]^\top \in \mathbb{R}^{|C| \times d}, \quad (10)$$

where c_i denotes a channel in the predefined electrode set and the dimension d is given by $d = D/H_q$, where H_q denotes the number of query heads. \mathbf{W}_{pos} is initialized using a truncated normal distribution. We employ Rotary Positional Encoding (RoPE) [68] as the temporal relative position encoding \mathbf{E}_t to support EEG samples of various lengths. For query \mathbf{q} and key \mathbf{k} at position t :

$$\text{RoPE}(\mathbf{z}_m) = \begin{bmatrix} \mathbf{z}_{m,1} \odot \cos(m\theta) - \mathbf{z}_{m,2} \odot \sin(m\theta) \\ \mathbf{z}_{m,1} \odot \sin(m\theta) + \mathbf{z}_{m,2} \odot \cos(m\theta) \end{bmatrix} \quad (11)$$

where m is the position index, $\mathbf{z}_{m,1}$ denotes the first half of the vector \mathbf{z}_m , while $\mathbf{z}_{m,2}$ denotes the second half. During the forward pass of the model, the sample embeddings \mathbf{x} are augmented by adding \mathbf{E}_t along the sequence dimension and \mathbf{E}_c along the channel dimension in the channel encoder and decoder, while only \mathbf{E}_t is added to the sample embeddings in the temporal encoder.

2.2.2 Channel Encoder. Since electrode positions are directly correlated with brain regions, channels play an important role in the analysis of EEG. However, this property impedes the application of the transformer model to EEG as a type of multivariate data. The channel encoder unit solves the problem effectively and in parallel by employing a cross-attention mechanism that models inter-channel relationships while preserving temporal patterns. Built upon an interleaved transformer architecture with rotary positional encoding and learned channel positional encoding, this unit addresses two challenges: 1) Capturing long-range dependencies between non-adjacent EEG channels; and 2) Maintaining temporal coherence during channel-wise feature aggregation.

Given the input tensor $\mathbf{x} \in \mathbb{R}^{B \times T \times C \times D}$ and the channel positional encoding \mathbf{E}'_c selected based on the available channels in a specific sample, \mathbf{x} is first reshaped into $\mathbb{R}^{(B \times T) \times C \times D}$ and then the self-attention layer calculations are performed in the same manner as in Equations 16 and 17, which will be discussed in Section 2.2.3.

Then the output representation \mathbf{h} is fed to the cross-attention layer. We implement cross-attention by reshaping $\mathbf{x}' \in \mathbb{R}^{B \times (T \times C) \times D}$, then applying RMS Normalization, linear projection, and RoPE:

$$\mathbf{Q} = \mathbf{W}_q \cdot \mathbf{l}_c, \quad \mathbf{K} = \mathbf{W}_k \cdot \text{RMSNorm}(\mathbf{x}'), \quad \mathbf{V} = \mathbf{W}_v \cdot \text{RMSNorm}(\mathbf{x}'), \quad (12)$$

$$\mathbf{Q}' = \text{RoPE}(\mathbf{Q}), \quad \mathbf{K}' = \text{RoPE}(\mathbf{K}) + \mathbf{E}'_c, \quad \mathbf{V}' = \mathbf{V}, \quad (13)$$

where $\mathbf{l}_c \in \mathbb{R}^{B \times T \times D}$ is the learned query for the cross-attention layer. Afterwards, cross-attention is formulated as:

$$\mathbf{h} = \text{CrossAttn}(\mathbf{Q}', \mathbf{K}', \mathbf{V}') = \text{Softmax} \left(\frac{\mathbf{Q}' \mathbf{K}'^\top}{\sqrt{d_k}} \odot \mathcal{M}_c \right) \mathbf{V}' \quad (14)$$

where \mathcal{M}_c implements the cross-attention mask described in Section 2.3. As mentioned above, two distinct attention layers form a channel encoder layer, which can be stacked to enhance the representation capability [25, 51]. For each channel encoder layer l , we update the representations as follows:

$$\mathbf{h}_l = \text{SelfAttn}(\mathbf{h}_{l-1}); \quad \mathbf{q}_l = \text{CrossAttn}(\mathbf{Q} = \mathbf{q}_{l-1}, \mathbf{K} = \mathbf{V} = \mathbf{h}_l), \quad (15)$$

where $\mathbf{h}_0 = \mathbf{x}'$ and $\mathbf{q}_0 = \mathbf{l}_c$. Compared to the cross-attention layer, self-attention layers perform computations in a different manner: 1) the query is replaced by \mathbf{h}_l , 2) the query is augmented with both \mathbf{E}'_c and RoPE, and 3) task-specific attention masks, which are designed to enable interactions among channels at the same time step, are applied to learn robust channel representations.

2.2.3 Temporal Encoder. The temporal encoder captures the general EEG features through stacked self-attention layers, which 1) separate channel and temporal feature aggregation; 2) capture multi-scale temporal dependencies owing to the input embedding; and 3) learn robust representations in various pretraining tasks.

For the input feature, $\mathbf{h}_c \in \mathbb{R}^{B \times T \times D}$ produced by the channel encoder, we introduce a task-specific token $\mathbf{t}_{DT} \in \mathbb{R}^D$ for different datasets and attach it at the head of the input feature to guide the module in leveraging attention in different patterns, where $\mathbf{h} = \mathbf{t}_{DT} \oplus \mathbf{h}_c$. The self-attention is then formulated as

$$\mathbf{h}_{\text{attn}} = \text{SelfAttn}(\text{RMSNorm}(\text{RoPE}(\mathbf{h}))) \quad (16)$$

However, we adopt the innovative Gated Feed-Forward layer (GatedFFN) design in [51] which involves two residual connections and

a gated feed-forward layer. The GatedFFN is formulated as

$$\begin{aligned} \mathbf{h}^* &= \text{RoPE}(\mathbf{h}) + W_O \mathbf{h}_{\text{attn}}; \mathbf{h}_{\text{norm}} = \text{RMSNorm}(\mathbf{h}); \\ \mathbf{h}_{\text{gate}} &= \text{SiLU}(W_{\text{gate}} \mathbf{h}_{\text{norm}}); \mathbf{h}_{\text{act}} = \text{SiLU}(W_{\text{in}} \mathbf{h}_{\text{norm}}) \odot \mathbf{h}_{\text{gate}}; \\ \mathbf{h}_{\text{out}} &= \mathbf{h}^* + W_{\text{out}} \mathbf{h}_{\text{act}}, \end{aligned} \quad (17)$$

where all W_* denotes weights in linear layers. The GatedFFN implements parameter-efficient feature interaction, enabling dynamic feature selection through multiplicative interactions and enhancing non-linear representation capability. Dual residual pathways create a gradient highway connection for deep layer training.

2.2.4 EEG Decoder. To reconstruct the EEG signal from features, we introduce the EEG decoder includes a hybrid attention mechanism to preserve: 1) temporal coherence in signal dynamics, and 2) channel-specific neuro-physiological patterns.

For the input, $\mathbf{h}_e \in \mathbb{R}^{B \times T \times D}$, from the temporal encoder, and the learnable query $\mathbf{l}_d \in \mathbb{R}^{B \times T \times C \times D}$, RoPE is added to \mathbf{h}_e , and both RoPE and channel-specific positional encoding are added to \mathbf{l}_d . The positional encoding serves as a reconstruction guide to achieve a better correlation with the original input.

In contrast to the channel encoder, we utilize a cross-attention layer first, followed by a self-attention layer, to construct a decoder layer, which forms a symmetric decoding. For the l -th decoder layer,

$$\mathbf{q}'_l = \text{CrossAttn}(Q = \mathbf{q}_{l-1}, K = V = \mathbf{h}_e); \mathbf{q}_l = \text{SelfAttn}(\mathbf{q}'_l), \quad (18)$$

where $\mathbf{q}_0 = \mathbf{l}_d$. This symmetric design stabilizes encoder-decoder co-training and eliminates the need for an extra stabilization module.

2.2.5 Pretraining Head. Given $\mathbf{q}_L \in \mathbb{R}^{B \times T \times C \times D}$ as the output of the EEG decoder, we feed it into two linear projectors for dual-domain reconstruction, $\hat{\mathbf{x}} = [\Phi_t(\mathbf{q}_L); \Phi_f(\mathbf{q}_L)]$, where Φ_t and Φ_f aim to project the EEG feature to temporal domain $\mathbb{R}^D \rightarrow \mathbb{R}^{D_t}$ and frequency domain $\mathbb{R}^D \rightarrow \mathbb{R}^{D_f}$ on feature dimension respectively.

2.2.6 Finetuning Head. The finetuning head is designed to implement multi-task adaptive learning on different downstream tasks. Therefore, the finetuning head is composed of a cross-attention layer and a task-specific CLS token dictionary. For the EEG feature \mathbf{h}_e from the temporal encoder, the CLS token $\mathbf{t}_{\text{CLS}}^k$ for the k -th task is used as the query to extract task-specific features from the cross-attention layer, which is formulated as:

$$\mathbf{f}_{\text{CLS}}^k = \text{CrossAttn}(Q = \mathbf{t}_{\text{CLS}}^k, K = V = \mathbf{h}_e), \quad (19)$$

where the task-specific features are then projected into the corresponding probability space, $p_\theta(c|\mathbf{t}_{\text{CLS}}^k)$, where θ is the parameters of the cross-attention layer and the task-specific projector.

2.3 Training Objectives

2.3.1 Pretraining Stage. Our pretraining framework aims to learn robust EEG representations through a multi-task learning approach. These three complementary pretraining tasks, including 1) GPT for **EEG signal forecasting**, 2) MAE-TP for **temporal masked patch reconstruction**, and 3) MAE-CH for **channel masked patch reconstruction**, achieve temporal dynamics modeling and cross-domain adaptability. As illustrated in Figure 2, the data patch sequence is obtained and then fed into the pretraining reconstruction head to compute the auto-regressive loss.

Loss function: The losses for pretraining are defined as follows. Typically we adopt RMSE as reconstruction function. For GPT task:

$$\mathcal{L}_{\text{GPT}} = \frac{1}{BC(T-1)} \sum_{n=1}^B \sum_{t=1}^{T-1} \sum_{c=1}^C \|\hat{\mathbf{x}}_{n,t+1,c} - \mathbf{x}_{n,t,c}\|_2, \quad (20)$$

where B is the batch size, T is the number of time steps, and C is the number of channels. For two MAE pretraining tasks:

$$\mathcal{L}_{\text{MAE}} = \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} \|\hat{\mathbf{x}}_{i,j} - \mathbf{x}_{i,j}\|_2, \quad (21)$$

where Ω denotes the set of masked patches, and (i, j) indicates the index i along the time step axis and j along the channel axis, respectively. Moreover, we introduce a task token, \mathbf{t}_{DT} , to guide the model to further adapt to datasets collected in different experimental paradigms. The task token categorizes the datasets into several classes, including Emotional Recognition, Motor Imaginary, Motor Execution, Seizure Detection, Artifact Classification, Sleep Staging, Resting, Event-Related Potential (ERP), Visual Stimulus and Workload Estimation. Therefore the task classification loss is:

$$\mathcal{L}_{\text{DT}} = -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_{n,c} \log(p_\theta(c|\mathbf{t}_{\text{DT}}^n)), \quad (22)$$

where N is the batch size (equal to B), C is the number of dataset classes, $y_{n,c}$ is the one-hot ground truth label, and p_θ is the pre-decision probability obtained by applying a Softmax function to the logits. Finally, the pretrain loss is formulated as

$$\mathcal{L}_p = \lambda_1 \mathcal{L}_{\text{GPT}} + \lambda_2 \mathcal{L}_{\text{MAE}_{\text{TP}}} + \lambda_3 \mathcal{L}_{\text{MAE}_{\text{CH}}} + \lambda_4 \mathcal{L}_{\text{DT}}, \quad (23)$$

where λ_i is set based on the convergence speed of each loss.

Masking: Another important aspect is our masking strategy. In supervised objectives, we define three reconstruction losses within an encoder-decoder architecture for EEG representation learning. Unlike prior approaches that mask inputs directly or use causal masking in attention, our framework employs diverse masking patterns across modules. This strategy separates spatial and temporal operations into dedicated modules and dynamically adjusts masks based on task objectives and training phases, thereby addressing different patch reconstruction needs without modifying the raw input and preventing information leakage. For example, when the second time step is masked in MAE-TP and the first channel is masked in MAE-CH, we design an attention mask matrix for each task in the channel encoder as shown in Figure 3 (a). In the temporal encoder, corresponding time steps are masked for MAE-TP, a lower triangular matrix is used for GPT, and an all-ones matrix is applied for MAE-CH, while the task token \mathbf{t}_{DT} remains universally accessible. In the decoder, as shown in Figure 3 (b), full sequence masks allow masked patches to attend to others while preventing visible patches from attending to masked ones, reducing noise interference.

2.3.2 Finetuning Stage. We design a fine-tuning framework implements multi-task adaptive learning. Samples from multiple datasets are fine-tuned on the same foundation model equipped with different classification heads. Lightweight classification heads adapt to diverse downstream objectives, while the foundation model enables positive knowledge transfer between related tasks. In this method, the pretrained representation capabilities are fully exploited, and downstream datasets can be integrated easily and

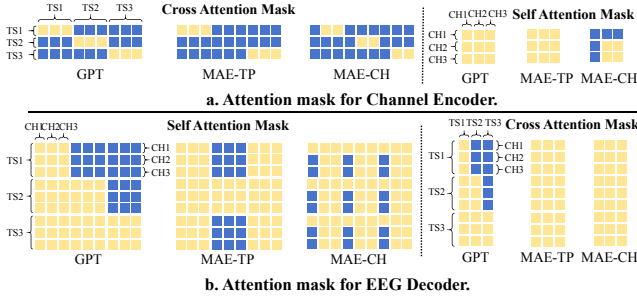


Figure 3: Attention masks in a) Channel Encoder and b) EEG Decoder. Attention is enabled as yellow. Rows and columns correspond to query (Q) and key (K) in attention calculation.

efficiently. Compared to SOTA multi-task EEG foundation models [28], there is no need to jointly train LLMs. As shown in Figure 2, we introduce a CLS token \mathbf{t}_{cls} to aggregate information from the compressed representations $\hat{\mathbf{x}}$ for classification.

Loss function: During the finetuning stage, K downstream tasks are trained simultaneously. The classification loss is formulated as:

$$\mathcal{L}_{cls}^k = -\frac{1}{N_k} \sum_{i=1}^{N_k} \sum_{c=1}^{C_k} w_c^k y_c^k[i] \log(p_{\theta}(c|\mathbf{t}_{cls}^k[i])), \quad (24)$$

where k is the task index, N_k is the number of samples in the mini-batch, C_k is the number of classes, w_c^k is the class weight computed as inversely proportional to the square root of class frequency, $y_c^k[i]$ is the one-hot indicator, and $p_{\theta}(c|\mathbf{t}_{cls}^k[i])$ is the probability that the sample CLS token $\mathbf{t}_{cls}^k[i]$ belongs to class c , with θ denoting the trainable model parameters. Besides, to preserve the ability of foundation model and prevent catastrophic interference, we introduce a reconstruction loss \mathcal{L}_{rec} calculated as in Equation 21 for all input patches. Finally, the fine-tuning loss is formulated as

$$\mathcal{L}_f = \sum_{k=1}^K (\alpha \mathcal{L}_{cls}^k + (1 - \alpha) \mathcal{L}_{rec}), \quad (25)$$

where α is weight for classification loss.

3 Experiments

3.1 Datasets

We follow the dataset composition in [28], which comprises 15 pretraining datasets and 6 evaluation datasets. The downstream evaluation dataset includes, TUAB [42], TUEV [22], TUSL [75], SEED [88], HMC [4], and Workload [89]. **More details about the dataset are in Supplementary Material.**

3.2 Experimental Setup

3.2.1 Data Preprocessing. An MNE [21] pipeline resamples, filters, aligns electrodes, converts units, compresses, and parallelizes EEG data, enabling reproducible large-scale feature extraction. **Detailed data process and evaluation are in Supplementary Material.**

3.2.2 Evaluation Metrics. To address the class imbalance observed in downstream EEG datasets, following the same metrics as in [28], we adopt Balanced Accuracy, Cohen’s Kappa, and Weighted F1 for multi-category classification, and Balanced Accuracy, AUC-PR, and

AUROC for binary classification. **Detailed descriptions of these metrics are in the Supplementary Material.**

3.2.3 Baselines. We compare our framework with three SOTA EEG foundation models (LaBraM [29], EEGPT [77], and NeuroLM [28]) and four supervised methods (SPaRCNet [30], ContraWR [82], FFCL [35], and BIOT [81]), where LaBraM, SPaRCNet, ContraWR, FFCL, and BIOT are designed for single-task scenarios and NeuroLM and EEGPT are multi-task methods. In the experiment, our ALFEE with its three variants (*i.e.*, Medium (M, 44.3M), Base (B, 120M), and Large (L, 300M)), is finetuned in two configurations: (1) a single-task setting, where the pretrained model is finetuned and evaluated on each of the six evaluation tasks separately; and (2) a multi-task setting, where the pretrained model is finetuned and evaluated on all six evaluation tasks jointly. **More details about the baselines, related work, model configurations, and training settings are provided in the Supplementary Material.**

3.3 Main Results

In Tables 1, 2, and 3, we present all the experimental results on both single-task and multi-task configuration. The best results on both settings are highlighted in bolded. All experiment are conducted for ten times, with mean \pm standard deviation reported. **Scaling law experiments are in Supplementary Material.**

Single-task Performance: In our single-task configuration, we evaluate ALFEE-M, ALFEE-B and ALFEE-L on each downstream dataset independently and compare their performance with several baselines, most notably LaBrAM-base. Overall, ALFEE-B achieves results comparable to LaBrAM-base on most datasets, while exhibiting higher Balanced Accuracy, Cohen’s Kappa, and Weighted F1 scores in TUEV and HMC. These results underscore ALFEE-B’s ability to capture essential EEG features with a relatively moderate model size. When we scale up to ALFEE-L, the performance improves further on most tasks, indicating that increasing model capacity confers distinct advantages in extracting nuanced EEG patterns. For instance, ALFEE-L exhibits superior sensitivity and specificity on TUAB, a dataset that contains more complex signals for medical diagnosis. Although LaBrAM-base remains competitive on certain metrics (*e.g.*, on SEED), our results show that ALFEE-B and ALFEE-L surpass it in most single-task settings. Beyond the numerical improvements, our analysis reveals two key characteristics of ALFEE in the single-task setting. First, the unified architecture of ALFEE-B achieves stable training dynamics and converges efficiently, suggesting that it can accommodate differences in input distributions. Second, as the model size varies from ALFEE-M to ALFEE-L, the larger model capacity of ALFEE-L allows it to capture subtle temporal and channel variations, leading to improved generalization on tasks requiring complex pattern recognition. However, this increased capacity also introduce overfitting risk, particularly in datasets with fewer samples (*e.g.*, Workload). Overall, our single-task evaluations confirm that ALFEE provides competitive and often superior performance compared to existing baselines.

Multi-task Performance: In the multi-task configuration, we evaluate ALFEE across four model sizes, ALFEE-M, ALFEE-B, ALFEE-L, on six datasets. Compared with NeuroLM and EEGPT, our models

Table 1: Results on TUAB and TUEV. Best performance on multi-task and single-task are highlight in bolded.

Methods	Multi-task	TUAB			TUEV		
		Balanced Acc.	AUC-PR	AUROC	Balanced Acc.	Cohen's Kappa	Weighted F1
SPaRCNet	✗	0.7896±0.0018	0.8414±0.0018	0.8676±0.0012	0.4161±0.0262	0.4233±0.0181	0.7024±0.0104
ContraWR	✗	0.7746±0.0041	0.8421±0.0104	0.8456±0.0074	0.4384±0.0349	0.3912±0.0237	0.6893±0.0136
FFCL	✗	0.7848±0.0038	0.8448±0.0065	0.8569±0.0051	0.3979±0.0104	0.3732±0.0188	0.6783±0.0120
BIOT	✗	0.7959±0.0057	0.8792±0.0023	0.8815±0.0043	0.5281±0.0225	0.5273±0.0249	0.7492±0.0082
LaBrAM-base	✗	0.8140±0.0019	0.8965±0.0016	0.9022±0.0009	0.6409±0.0065	0.6637±0.0093	0.8312±0.0052
ALFEE-M	✗	0.8069±0.0109	0.8835±0.0191	0.8828±0.0142	0.6190±0.0261	0.6612±0.0299	0.7953±0.0132
ALFEE-B	✗	0.8108±0.0072	0.8899±0.0043	0.8859±0.0019	0.6439±0.0052	0.7529±0.0073	0.8481±0.0093
ALFEE-L	✗	0.8239±0.0091	0.9042±0.0059	0.9015±0.0028	0.6525±0.0049	0.7737±0.0103	0.8626±0.0124
EEGPT	✓	0.7983±0.0030	-	0.8718±0.0050	0.6232±0.0114	0.6351±0.0134	0.8187±0.0063
NeuroLM-B	✓	0.7826±0.0065	0.6975±0.0081	0.7816±0.0079	0.4560±0.0048	0.4285±0.0048	0.7153±0.0028
ALFEE-M	✓	0.7951±0.0048	0.8515±0.0162	0.8589±0.0039	0.6553±0.0124	0.6607±0.0098	0.7959±0.0052
ALFEE-B	✓	0.8074±0.0082	0.8588±0.0171	0.8763±0.0091	0.7173±0.0102	0.7254±0.0222	0.8342±0.0093
ALFEE-L	✓	0.8090±0.0072	0.8861±0.0252	0.8812±0.0082	0.6987±0.0145	0.7863±0.0273	0.8683±0.0102

Table 2: Results on SEED and HMC. Results of EEGPT are reproduced based on public released code.

Methods	Multi-task	SEED			HMC		
		Balanced Acc.	Cohen's Kappa	Weighted F1	Balanced Acc.	Cohen's Kappa	Weighted F1
SPaRCNet	✗	0.5596±0.0244	0.3464±0.0372	0.5585±0.0297	0.4756±0.1109	0.3147±0.1315	0.4108±0.1310
ContraWR	✗	0.6106±0.0078	0.4220±0.0129	0.6137±0.0085	0.4242±0.0541	0.2340±0.0554	0.2987±0.0288
FFCL	✗	0.5808±0.0322	0.3732±0.0462	0.5743±0.0402	0.4427±0.0702	0.2542±0.0654	0.2902±0.0485
BIOT	✗	0.7097±0.0024	0.5682±0.0051	0.7134±0.0027	0.6862±0.0041	0.6295±0.0113	0.7091±0.0147
LaBrAM-base	✗	0.7318±0.0019	0.5994±0.0031	0.7354±0.0021	0.7286±0.0101	0.6812±0.0073	0.7554±0.0024
ALFEE-M	✗	0.6561±0.0088	0.4854±0.0049	0.6572±0.0093	0.7229±0.0105	0.6636±0.0091	0.7378±0.0066
ALFEE-B	✗	0.6963±0.0029	0.5461±0.0021	0.6924±0.0038	0.7490±0.0133	0.6922±0.0109	0.7627±0.0092
ALFEE-L	✗	0.6931±0.0035	0.5414±0.0052	0.6848±0.0045	0.7338±0.0162	0.7083±0.0134	0.7720±0.0112
EEGPT	✓	0.7122±0.0022	0.5734±0.0049	0.7099±0.0038	0.7029±0.0082	0.6584±0.0059	0.7323±0.0041
NeuroLM-B	✓	0.5554±0.0075	0.3393±0.0117	0.5599±0.0068	0.6737±0.0050	0.6188±0.0057	0.7126±0.0034
ALFEE-M	✓	0.6578±0.0034	0.4880±0.0062	0.6592±0.0082	0.7197±0.0073	0.6782±0.0083	0.7551±0.0039
ALFEE-B	✓	0.7248±0.0052	0.5887±0.0089	0.7221±0.0102	0.7388±0.0098	0.6815±0.0078	0.7584±0.0082
ALFEE-L	✓	0.7411±0.0072	0.6153±0.0078	0.7432±0.0079	0.7378±0.0074	0.6837±0.0068	0.7610±0.0083

Table 3: Results on Workload and TUSL. Results of EEGPT are reproduced based on public released code.

Methods	Multi-task	Workload			TUSL		
		Balanced Acc.	AUC-PR	AUROC	Balanced Acc.	Cohen's Kappa	Weighted F1
SPaRCNet	✗	0.5977±0.0071	0.6638±0.0314	0.6717±0.0172	0.4185±0.0452	0.1399±0.0799	0.3500±0.0968
ContraWR	✗	0.6966±0.0332	0.7668±0.0408	0.7685±0.0317	0.5857±0.0662	0.3567±0.0968	0.5458±0.0798
FFCL	✗	0.7069±0.0197	0.7823±0.0099	0.7857±0.0234	0.3819±0.0688	0.0628±0.0888	0.2120±0.0786
BIOT	✗	0.6655±0.0665	0.7189±0.0722	0.7342±0.0536	0.5758±0.0303	0.2012±0.0212	0.2394±0.0040
LaBrAM-base	✗	0.6609±0.0204	0.7174±0.0234	0.7272±0.0165	0.7625±0.0131	0.6407±0.0304	0.7614±0.0210
ALFEE-M	✗	0.6409±0.0298	0.3315±0.0272	0.7661±0.0379	0.7262±0.0201	0.5217±0.0322	0.6588±0.0208
ALFEE-B	✗	0.6518±0.0183	0.5922±0.0211	0.8004±0.0138	0.7535±0.0182	0.6417±0.0277	0.7573±0.0299
ALFEE-L	✗	0.7333±0.0201	0.6420±0.0232	0.8151±0.0172	0.7698±0.0129	0.6077±0.0329	0.7336±0.0310
EEGPT	✓	0.6299±0.0178	0.6792±0.0092	0.6928±0.0108	0.7288±0.0143	0.5972±0.0209	0.7233±0.0153
NeuroLM-XL	✓	0.6172±0.0113	0.5824±0.0080	0.6253±0.0160	0.6734±0.0436	0.5107±0.0617	0.6743±0.0394
ALFEE-M	✓	0.6285±0.0223	0.4268±0.0272	0.6339±0.0249	0.7292±0.0111	0.6283±0.0219	0.7490±0.0182
ALFEE-B	✓	0.6849±0.0239	0.6871±0.0282	0.8295±0.0292	0.8245±0.0172	0.7441±0.0282	0.8312±0.0143
ALFEE-L	✓	0.7333±0.0244	0.5516±0.0258	0.7992±0.0283	0.8680±0.0181	0.8158±0.0302	0.8799±0.0154

consistently achieve higher performance across all metrics. Specifically, on the SEED dataset, ALFEE-M attains a significant improvement (10%–15% across three metrics) over NeuroLM, while ALFEE-L pushes the performance even further (20%–26% across three metrics), surpassing single-task approaches and providing evidence of multi-task synergy. Notably, on TUEV, ALFEE-L exceeds the performance of LaBrAM, reaffirming that a well-designed multi-task approach can outperform single-task strategies. A similar trend is also observed on HMC, where ALFEE-B and ALFEE-L demonstrate

clear advantages over the baselines and maintain a robust margin of improvement in all metrics. Furthermore, on TUSL, despite its limited data size, our larger models still outperformed EEGPT and NeuroLM, emphasizing the flexibility of ALFEE in handling small-scale datasets without severe overfitting. One of the key observations is that scaling up our model size from ALFEE-M to ALFEE-L tends to achieve better performance in most tasks. At the same time, our results on tasks with fewer samples suggest that heavier models require careful early stopping strategies to avoid overfitting. The

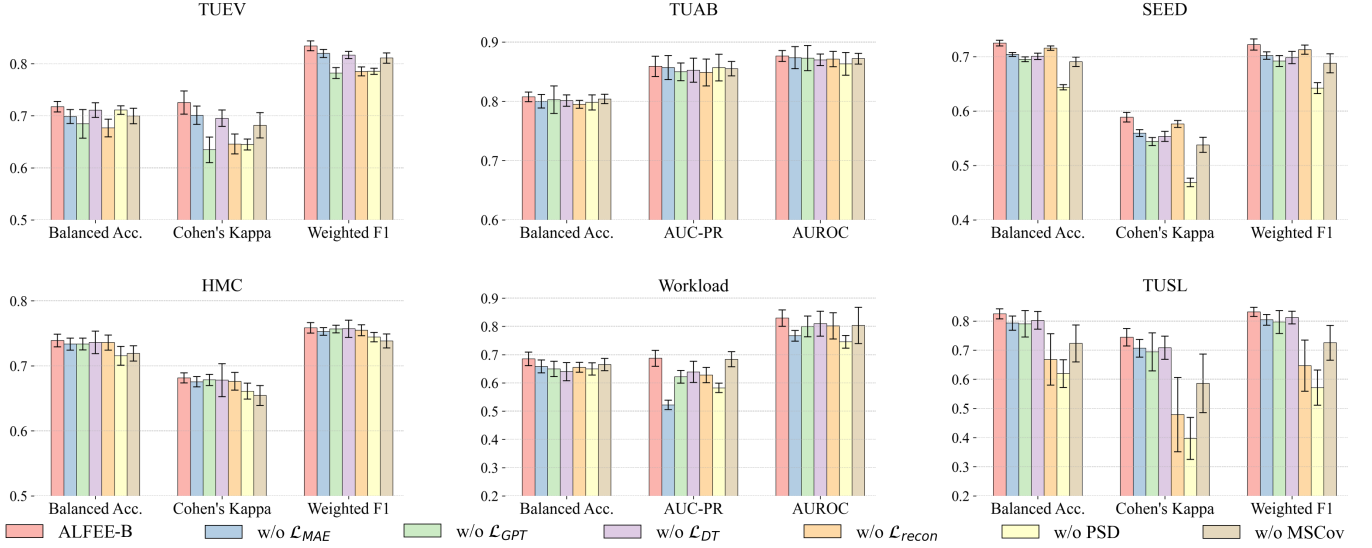


Figure 4: Ablation Study. We study the effect of four loss functions (w/o \mathcal{L}_{MAE} , w/o \mathcal{L}_{GPT} , and w/o \mathcal{L}_{DT} , w/o \mathcal{L}_{recon}), and two feature extraction module (the frequency domain feature, w/o PSD and the multi-scale convolution, w/o MSCov).

complementary information shared across tasks in the multi-task setting contributes to the improved performance, enabling ALFEE to acquire more robust representations of EEG signals. This effect is particularly meaningful because it provides a more efficient and unified framework for training, obviating the need for intensive task-specific fine-tuning. Overall, ALFEE outperform multi-task methods, demonstrating that a larger model capacity is beneficial for representation learning with appropriately dataset balance.

3.4 Ablation Study

In this section, we study the effect of loss functions (w/o \mathcal{L}_{MAE} , w/o \mathcal{L}_{GPT} , and w/o \mathcal{L}_{DT} for the pretraining stage, and w/o \mathcal{L}_{recon} for the finetuning stage) and two feature extraction modules (the frequency domain feature, w/o PSD, and the multi-scale convolution, w/o MSCov) using ALFEE-B as shown in Figure 4.

Effect of Different Losses: Comparing the performance among the three pretraining losses, we notice that \mathcal{L}_{GPT} and \mathcal{L}_{MAE} have a relatively large impact on downstream performance; without them, performance would degrade by approximately 2% to 10% across different metrics. Moreover, \mathcal{L}_{GPT} also has a large impact on model stability, as the standard deviation for “w/o \mathcal{L}_{GPT} ” in TUAB, TUEV, HMC, and TUSL is two times larger than that of ALFEE-B. Considering that temporal forecasting is a more advanced ability than mask reconstruction, such instability also indicates the effectiveness of \mathcal{L}_{GPT} for robust EEG representation. As for \mathcal{L}_{DT} , it has little impact on overall performance but is helpful for model stability because it aids the model in capturing task-specific characteristics beyond channel relationships and temporal dynamics during pretraining. In the finetuning stage, \mathcal{L}_{recon} is very useful, particularly when the downstream dataset is small, such as TUSL and Workload, thereby proving its effectiveness in preventing catastrophic interference during multi-task finetuning.

Effect of Different Feature Extraction Modules: For the feature extraction modules, we notice that both are very important for

downstream performance. The multi-scale convolution captures temporal dynamics at different granularities, thereby providing additional temporal patterns for downstream classification, especially in TUSL, HMC, and SEED datasets. The frequency domain feature is critical for TUEV, SEED, Workload, and TUSL datasets. For the Workload dataset, the γ -wave in the frequency domain is very important for attention assessment. For the TUSL dataset, the β -wave in the frequency domain helps detect seizures. For SEED and TUEV, the electromyographic signal in the frequency domain from the face can introduce noise during EEG collection, therefore, frequency domain feature can filter out such noise. These results suggest the effectiveness of our loss functions and feature extractors in the ALFEE for robust EEG representation.

4 Conclusion

In this paper, we propose the ALFEE framework, a EEG foundation model that incorporates a hybrid attention architecture to separate channel-wise feature aggregation from temporal dynamics modeling, enabling robust EEG representation with variable channel configurations. By separating EEG representations into adaptive channel-wise and task-guided temporal components, ALFEE leverages a hybrid attention architecture and a dual-stage optimization process, employing multi-task pretraining followed by multi-task fine-tuning with cross-attention layers, to reconstruct both temporal and frequency domain features and address the diverse challenges inherent in EEG signal representation. Our extensive experiments across six diverse EEG datasets demonstrate the superior performance of the ALFEE framework in EEG signal representation and multi-task learning. Overall, ALFEE represents a major advancement in brain-computer interfaces and healthcare, underscoring the significant potential of EEG foundation models in signal processing and multi-task learning. Moreover, scaling law analysis confirms that large-scale models such as ALFEE offer promising prospects for EEG applications and human-machine interactions.

In the Appendix, we provide related work in Appendix A to introduce the progress in the relevant research field. We also provide detailed descriptions of the datasets used in our two learning stages in Appendix B, and experimental settings in Appendix C. Moreover, we conduct an additional experiment to validate the scaling law in EEG signal representation in Appendix E. Finally, we conclude by discussing the limitations and future work in Appendix F.

A Related Work

A.1 EEG Modeling Approaches

Traditional EEG analysis initially relies on time-frequency features extracted through the discrete wavelet transform [8] or PSD [3] via Fourier transforms, combined with classifiers such as support vector machines or multilayer perceptrons. While these methods achieve reasonable performance in brain information decoding, their limited generalizability across experimental paradigms and subjects prompts the development of deep learning architectures. Early convolutional neural networks, including ConvNet [63] and EEGNet [34], establish end-to-end pipelines through spatiotemporal convolutions, though their capacity to capture long-range dependencies remained constrained. To address this limitation, recurrent neural networks, such as ChronoNet [60] and ATDD-LSTM [16], are introduced, thereby improving temporal modeling at the cost of training efficiency and gradient stability.

Graph neural networks like EEG-GNN [11] subsequently emerge, leveraging the electrode-channel relationships for neuroscientifically interpretable analysis. However, their fixed node representations often compromise temporal resolution and adaptability. The recent shift toward attention mechanisms and transformer architectures has yielded notable advancements in EEG analysis. ST-Transformer [65] pioneers transformer-based EEG analysis, while EEGConformer [66] hybridizes convolutional local feature extraction with global self-attention. Despite their potential, these models face challenges in cross-paradigm generalization and often require extensive datasets to mitigate overfitting. Self-supervised learning approaches have further addressed data scarcity. REMoNet [27] employs channel masking and spectral feature prediction for emotion recognition, while EEG2Rep [47] learns noise-invariant representations through abstract feature reconstruction. However, their reliance on task-specific feature designs limits broader applicability across diverse experimental paradigms.

A.2 Transformer in Foundation Models

The impact of self-supervised transformers in natural language processing (NLP), as exemplified by the masked language modeling of BERT [13] and the autoregressive pretraining of GPT [1, 43, 55, 56], has catalyzed cross-modal adaptations. Vision transformers [15] demonstrate that transformer architectures can surpass convolutional networks in image classification by processing sequences of image patches. CLIP [54] extends this paradigm through contrastive learning on massive image-text pairs, achieving unprecedented zero-shot generalization. MAE [23] further advances latent representation learning by reconstructing randomly masked

image patches with an asymmetric encoder-decoder design. Recently, BLT [51] introduces a novel architecture with alternating self-attention and cross-attention layers, enabling dynamic attention allocation across variable-length sequences, which is also a key inspiration for our hybrid attention mechanism.

While transformers dominate NLP and computer vision, their adoption in time series analysis remains nascent. For univariate forecasting, Lag-Llama [58] establishes a foundation model using autoregressive training on diverse datasets, demonstrating robust generalization to unseen domains. Multivariate approaches face greater complexity: PatchTST [48] addresses this through channel-independent patch embedding strategies, significantly improving long-term forecasting accuracy. Crossformer [86] enhances cross-variate dependency modeling via dimension-aware attention, while Pathformer [9] integrates multi-scale temporal resolutions through adaptive pathway selection. Morai [78] introduces innovations like multiple patch size projections and any-variate attention to handle heterogeneous sampling rates and variable input dimensions. Despite these advances, existing models primarily target structured numerical data (e.g., stock prices, statistical data), lacking mechanisms to handle the unique challenges of bioelectrical signals—non-stationarity, low signal-to-noise ratios, and inter-subject variability—which limit their direct applicability to EEG analysis.

A.3 Multi-Task Learning

The scarcity of labeled data for individual tasks motivates the adoption of multi-task learning to enhance data utilization efficiency by mining inter-task relationships and learning shared representations, thereby improving model generalization and transferability. MT-DNN [40] demonstrates robust domain adaptation capabilities by integrating four distinct task categories, also effectively mitigating overfitting while promoting universal feature learning. T5 [57] standardizes diverse NLP tasks into a unified text-to-text framework, revealing that multi-task pretraining with subsequent fine-tuning achieves performance parity with single-task pretraining approaches. UniLM [14] extends the architecture of BERT [13] through dynamic attention masking mechanisms, enabling a single model architecture to address heterogeneous language tasks without structural modifications. OmniNet [53] employs HogWild parallel training to facilitate cross-modal knowledge transfer in multi-task scenarios, while UniT [25] implements an end-to-end unified transformer that concurrently learns seven distinct tasks across eight datasets through joint training, achieving competitive performance while being parameter-efficient. For temporal data processing, UniTS [18] develops a domain-agnostic framework by integrating multi-domain datasets, demonstrating the superiority of specialized time series transformers over language-oriented LLMs in handling classification, forecasting, and imputation tasks. Nevertheless, these existing architectures require substantial adaptation to effectively process time-series neurophysiological data and achieve optimal performance in EEG applications.

A.4 Large Scale Pretrained EEG Models

The success of self-supervised learning in language modeling has spurred its adoption for EEG analysis, particularly given the scarcity of clinical annotations. BENDR [32] pioneers self-supervision to

Table 4: Detailed information about pretraining datasets.

Dataset	Category	#Channel	Duration	#Train	#Valid	Description
TUEG [49]	Clinical Recordings	60	60	1515391	75436	A rich archive of 26,846 clinical EEG recordings collected at Temple University Hospital.
SEED-IV [87]	Emotion Recognition	60	10	13678	977	A emotion EEG dataset that 15 subjects watch 72 film clips which have the tendency to induce happiness, sadness, fear or neutral emotions.
SEED-V [38]	Emotion Recognition	60	10	13246	2088	A multi-modal emotion dataset that 20 subjects watch video clips in happy, sad, disgust, fear and neutral.
SEED-GER [39]	Emotion Recognition	60	10	8440	919	A emotion dataset eight German subjects watch 20 film clips (positive, neutral, negative) as stimuli.
SEED-FRA [39]	Emotion Recognition	60	10	6846	978	Eight French subjects watch 21 film clips in French (positive, neutral, negative) as stimuli.
BCIC-1A [7]	Motor Imagery	43	8	3155	535	EEG recordings for motor imagery tasks, where subjects imagined moving either their left hand, right hand, or foot.
Emobrain [61]	Emotion Recognition	54	10	1370	405	Multimodal emotion detection dataset using brain signals (EEG, fNIRS) from 5 male subjects.
Grasp and Lift [41]	Motor Execution	32	5	7003	1390	Grasp and lift action from 12 subjects in total, 10 series of trials for each subject.
Inria BCI P300 [44]	ERP	56	5	13647	7901	A P300-based spelling dataset including 26 subjects.
Motor Movement Imagery [62]	Motor Imagery	64	12	13516	650	1500 EEG recordings dataset obtained from 109 volunteers perform opens four action in both in imaginary and reality.
Resting State (Trujillo 2017) [72]	Resting	64	10	981	100	A dataset comprising 22 subjects for a resting eyes closed and eyes open.
Raw EEG Data (Trujillo 2019) [71]	Visual Stimulus	64	30	4806	319	EEG was recorded during reported Information-Integration categorization and reported multidimensional Rule-Based categorization tasks.
Siena Scalp EEG Database [12]	Seizure Detection	27	40	8260	4408	EEG recordings of 14 patients that are labeled epilepsy and the classification of seizures is carefully reviewed.
SPIS Resting State [70]	Resting	64	10	270	30	A resting-state EEG from 10 subjects contains 2.5 minutes of eyes-open and 2.5 minutes of eyes-closed.
Target Versus Non-Target [31]	ERP	32	15	3403	333	dataset contains EEG recordings of 50 subjects playing to a visual P300 Brain-Computer Interface (BCI) video game.

learn general EEG data distributions, while EEG2Vec [5] introduces a conditional variational autoencoder framework for joint generative-discriminative representation learning. BIOT [81] addresses practical challenges such as variable electrode configurations and signal durations, demonstrating enhanced performance on clinical benchmarks. BrainBERT [76] adapts the methodology of BERT [13] to stereo-electroencephalography analysis through time-frequency representation learning. The emergence of large-scale pretrained EEG transformer models has significantly enhanced generalization capacity and analytical precision in biosignal processing. The Brant [83–85] series models, trained on terabyte-scale datasets, exhibit exceptional robustness against data variability

and excellent scalability, even demonstrating proficiency in modeling EEG-physiological signal correlations. LaBraM [29] currently achieves SOTA performance through its innovative integration of VQ-VAE [73] modules with dual-domain (frequency/phase) autoregressive learning. NeuroLM [28] bridges neurosignal-language modality gaps by embedding EEG signal into pretrained LLM frameworks, thereby balanced performance across unified downstream tasks. EEGPT [77] advances the field through the combination of dual self-supervised universal representation learning and stabilization mechanisms inspired by MoCo [24]. However, challenges such as variant channel counts, inadequate channel-temporal supervision, and upstream–downstream domain gaps still limit the performance of existing methods.

Table 5: Detailed information about evaluation datasets.

Dataset	Category	#Channel	Duration	#Train	#Valid	#Test	Task
TUAB	Clinical Recording	23	30	247728	12315	12277	Binary Classification
TUEV	Artifact Detection	21	5	87834	12473	13046	6-class Classification
TUSL	Seizure Detection	21,22	10	222	43	25	3-class Classification
SEED	Emotion Recognition	60	10	22455	7875	7560	3-class Classification
HMC	Sleep Staging	4	30	91681	22804	22440	5-class Classification
Workload	Workload Estimation	19	10	1537	300	297	Binary Classification

B Datasets Description

The detailed information about all 15 pretraining datasets are listed in Table 4. All data in the duration column is measured in seconds. Besides, the detailed information about the evaluation datasets is listed in Table 5: 1) The **TUAB** [42] dataset contains EEG records that are classified as clinically normal or abnormal; 2) The **TUEV** [22] dataset contains sessions that include events such as periodic lateralized epileptiform discharges, generalized periodic epileptiform discharges, spike and/or sharp wave discharges, artifacts, and eye movements; 3) The **TUSL** [75] dataset contains sessions that include seizure events, slowing events, and complex background events; 4) The **SEED** [88] dataset contains physiological signal data and corresponding emotion labels. Fifteen Chinese movie clips were carefully selected as stimuli to evoke different emotions; 5) The **HMC** [4] dataset collects 151 whole-night polysomnographic (PSG) sleep recordings as well as event annotations corresponding to the scoring of sleep patterns (hypnogram) performed by sleep technicians; 6) The **Workload** [89] dataset contains EEG recordings from 36 healthy volunteers during mental serial subtraction along with corresponding reference background EEGs. Based on task performance, subjects are divided into two groups: one for background and one for arithmetic.

C Experimental Setting

C.1 Data Preprocessing

The dataset preprocessing pipeline employs a structured approach for handling EEG data from various datasets using MNE-tools [21]. The pipeline initializes with a resampling operation that transforms the source sampling rate to $f_s = 256$ Hz, which facilitates patch division. An overlap-add Finite Impulse Response (FIR) high-pass filter is applied to remove low-frequency noise, otherwise the signal length is too short to meet the required filter length. Subsequently, a 50 Hz or 60 Hz notch filter is applied after human review. Afterwards, electrode configurations from the datasets are aligned with the predefined 10-10 channel set. Data unit conversion is performed from μV to Volts for MNE compatibility. For specific implementation, EEG data is serialized in Parquet format with Zstandard compression to expedite dataset loading. Remote storage is supported via the S3 protocol for distributed computing. This processing pipeline enables reproducible feature extraction while maintaining physiological signal fidelity. The implementation leverages parallel processing for efficient large-scale data handling.

C.2 Baselines

We mainly consider three SOTA EEG foundation models, namely LaBraM [29], EEGPT [77], and NeuroLM [28], as our baseline methods. LaBraM is pre-trained on 2,500 hours of data with integrated VQ-VAE [73] modules for dual-domain (frequency/phase) mask learning. EEGPT combines dual self-supervised universal representation learning and stabilization mechanisms, and NeuroLM bridges neurosignal-language modality gaps by embedding EEG signals into pretrained LLM frameworks using 25,000 hours of data. In addition, we select four supervised methods for comparison, including SPaRCNet [30], ContraWR [82], FFCL [35], and BIOT [81]. Among these methods, LaBraM, SPaRCNet, ContraWR, FFCL, and BIOT are designed for single-task scenarios and are also used as baselines in NeuroLM [28]. Therefore, we adopt the reported results from NeuroLM [28] in our experiments. Besides, since NeuroLM [28] and EEGPT [77] are multi-task methods, we also compare with them under multi-task settings. Since EEGPT [77] does not report its performance on the SEED, HMC, Workload, and TUSL datasets, we reproduce the results using its official, publicly available pretrained model weights and code. We follow the default settings in their linear_prob algorithm for training on these datasets.

C.3 Evaluation Strategy

We implement a greedy algorithm-based multi-label stratified splitting function to divide a dataset into training, validation, and test sets by subject, while ensuring that the label distribution in each split is balanced and aligns with predefined ratios. Subsequently, dataset partitioning is performed by splitting subjects into training, validation, and test sets as follows: 1) **TUAB**: The validation and test sets are obtained by equally splitting the original evaluation set by subject, while the training set remains unchanged; 2) **TUEV** and **TUSL**: Owing to the highly imbalanced label distribution in these datasets, the stratified splitting function is employed to create three splits from all the data, approximately aligning with predefined ratios of 0.8, 0.1, and 0.1; 3) **SEED**: Following prior research, the 15 trials are divided into three sets in a 9:3:3 ratio, and all sessions are merged together thereafter; 4) **HMC**: Subjects are randomly split into training, validation, and test sets at a ratio of 103:24:24; 5) **Workload**: Stratified splitting is employed to achieve approximate ratios of 0.72, 0.14, and 0.14 for training, validation, and test sets, respectively.

C.4 Evaluation Metrics

To address the class imbalance commonly observed in downstream EEG datasets, the following evaluation metrics are adopted for performance comparison:

- **Balanced Accuracy**: the arithmetic mean of recall (sensitivity) across all classes, mitigating the impact of imbalanced class distributions. It is particularly effective for evaluating classification models on datasets with significant disparities in class proportions.
- **Weighted F1**: a harmonic mean of precision and recall, weighted by the number of true instances in each class. This metric accounts for class imbalance by assigning higher importance to classes with larger sample sizes, ensuring a more representative evaluation of model effectiveness.
- **AUROC**: area under the ROC curve. It reflects the model’s ability to discriminate between classes across all possible decision boundaries.
- **AUC-PR**: area under the precision-recall curve. It provides a holistic evaluation of model performance under class imbalance.
- **Cohen’s Kappa**: the agreement level between predicted and true labels by comparing observed and expected frequencies along the diagonal of a confusion matrix. It is particularly suited for multi-class classification scenarios.

Among these metrics, AUROC and AUC-PR are used to evaluate binary classification tasks, while Cohen’s Kappa and Weighted F1 are applied to multi-category classification. Together, these metrics provide a robust evaluation framework under class imbalance.

C.5 Training Settings

To facilitate data loading, all samples in the datasets are transformed into an Arrow dataset after cleaning and preprocessing, thereby speeding up distributed computing and leveraging the GPU’s direct data access functionality. All experiments are conducted using Python 3.11.11, PyTorch 2.6.0, and CUDA 12.4 on eight A100 GPUs for pretraining and two A800 GPUs for finetuning. We enable autonomous mixed precision in the bfloat16 data type to improve GPU memory utilization and introduce GradScaler to prevent gradient explosion. We employ the AdamW optimizer and a two-phase learning rate scheduler, which combines linear warmup with cosine annealing. During the pretraining stage, the learning rate is set to 1×10^{-4} , and the hyperparameters λ_1 , λ_2 , λ_3 , and λ_4 are set to 0.4, 0.275, 0.275, and 0.05, respectively, to balance the different loss terms. During the finetuning stage, the learning rate is set to 5×10^{-5} , and α is set to 0.9 to prevent catastrophic interference. To compare with existing methods under different settings, the finetuning stage is conducted in two configurations: (1) a single-task setting, where the pretrained model is finetuned and evaluated on each of the six evaluation tasks separately; and (2) a multi-task setting, where the pretrained model is finetuned and evaluated on all six evaluation tasks jointly. The best models are trained on the training set, selected based on performance on the validation set, and finally evaluated on the test set.

C.6 Model Configurations

ALFEE has five variants: Small (S), Medium (M), Base (B), Large (L) and Extra Large (XL), with 16.3M, 44.3 M, 120 M, 300 M, and 540 M parameters, respectively. We can alter the parameter scale by changing the number of stacked transformer blocks and the embedding dimensionality. The maximum learnable query length and sequence length are set to 2048 to accommodate the longest input EEG patch sequences. The patch size is equal to $f_s = 256$, representing 1 second. 50% of channels and 40% of time steps are randomly masked. More details regarding the settings of the model variants are listed in Table 6.

D Visualization

D.1 Interpretable Decision Localization

In this section, we exploit Gradient-weighted Class Activation Mapping (Grad-CAM) [64] on the SEED dataset to reveal how our ALFEE framework captures emotional information from the spatial channels of EEG signals. Specifically, Grad-CAM is a visualization technique that highlights critical regions in input images that influence CNN predictions. By computing gradient-based weights from the target class score to the original EEG features, the resulting heatmaps provide intuitive explanations for model decisions. This technique has also been adopted to extract information about activated brain regions from model parameters and to perform visualization to enhance explainability in neuroscience [37]. In our implementation, we attach it to the channel encoder to visualize channel-wise decision-making evidence.

The attention of our ALFEE framework to different target classes may vary across individuals; however, the averaged Grad-CAM visualization on a standard brain model reflects the model’s overall attention patterns. As shown in Figure 5, this visualization confirms that our ALFEE framework focuses on relevant electrodes or brain regions for classification. Apparent lateralization can be observed, which is consistent with previous research on the functional networks of emotion [17]. In Figure 5a, our framework pays more attention to Fp1, Fpz, F4, F8, and T8 electrodes, suggesting that the right hemisphere may be dominant in negative emotion. For positive emotion, the left frontal lobe around F7 and FT7 receives strong attention from our framework. Additionally, other regions including Pz, TP7, and Oz can be viewed as the basis for classification in neutral emotion. All of these results are consistent with established scientific evidence [37, 52].

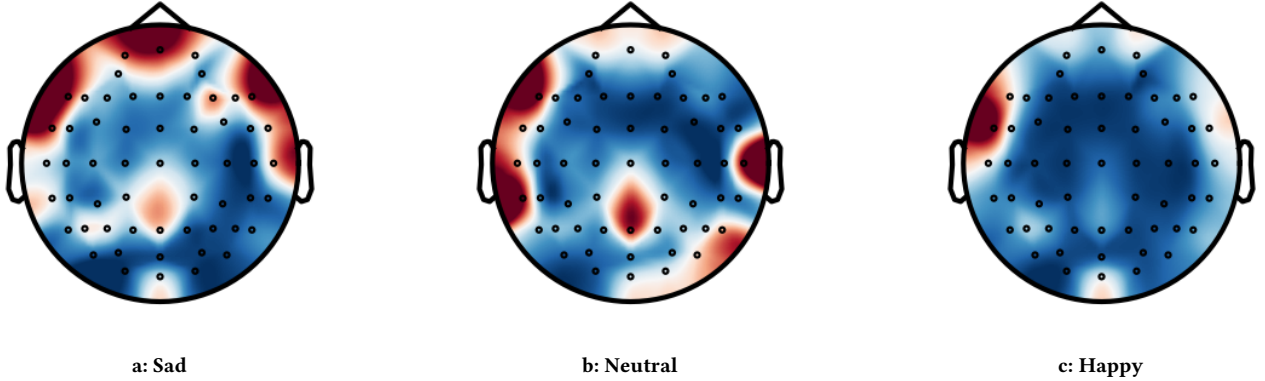
D.2 Latent Feature Clustering via t-SNE

t-SNE is a nonlinear dimensionality reduction technique designed to visualize high-dimensional data in a low-dimensional space (typically 2D or 3D) while preserving local structures and cluster relationships [74]. In Figure 6, we visualize the features of six datasets from the classification head of ALFEE to demonstrate the representational ability of our framework.

For the TUAB dataset, as shown in Figure 6a, two lateral clusters emerge with a broad transitional region, which are approximately linearly separable. For the TUEV dataset, as shown in Figure 6b, the six predefined classes show strong separation, except for the diffuse artifact category spanning the embedding space, consistent with

Table 6: Configurations for model variants.

	ALFEE-S	ALFEE-M	ALFEE-B	ALFEE-L	ALFEE-XL
#Parameters	16.3 M	44.3 M	120 M	300 M	540 M
Model Dim	384	512	640	896	1152
MLP Dim	256	512	512	768	768
Channel Block	1	1	2	2	3
Temporal Block	5	7	14	16	19
Decoder Block	2	4	8	10	12
Batch Size	256	256	128	96	64
Attn. Head	4	4	8	8	12

**Figure 5: Average Grad-CAM visualization showing the model’s region of interest for target class prediction. Warmer colors indicating higher relevance, generated by computing gradient flow of channel encoder on ALFEE-B.**

its miscellaneous definition during the annotation process. For the SEED dataset, as shown in Figure 6c, positive and negative emotions form distinct clusters, while neutral samples exhibit intermediate positioning with partial overlap, likely due to inter-individual variability in perception. For the HMC dataset, as shown in Figure 6d, the five sleep stages form tightly grouped clusters, reflecting high intra-class consistency aligned with physiological patterns. These results confirm the representation ability of ALFEE to disentangle discriminative features while tolerating noise.

E Scaling Laws

In this section, we study the relationship between model performance and training scale, including the model size and pretraining data size. For the model size, we design 5 variants, as described in Appendix C.6, to investigate the effects of model size on pretraining loss and downstream task accuracy. For the data size, we change the pretraining data from 2,500 hours to 25,000 hours on ALFEE-B to investigate the effects of pretraining dataset on pretraining loss and downstream task accuracy. The pretraining loss is calculated on a small partition of validation pretraining data, and the downstream

task accuracy is calculated based on the balanced accuracy on SEED and TUSL dataset.

E.1 For Parameter Size

In Figures 7a, 7b, 7c, we provide the results of the scale law experiments on the pretraining loss function (\mathcal{L}_p), SEED, and TUSL dataset. The results on the pretraining loss function show that the scaling law of the test \mathcal{L}_p with model size (N) is: $\mathcal{L}_p = -0.029 * \ln(N) + 1.173$, where R^2 is 0.972. The results on the SEED dataset show that the scaling law of the test balanced accuracy with model size (N) is: $B\text{Acc} = 0.036 * \ln(N) + 0.028$, where R^2 is 0.988. The results on the TUSL dataset show that the scaling law of the test balanced accuracy with model size (N) is: $B\text{Acc} = 0.056 * \ln(N) - 0.236$, where R^2 is 0.919. The results of the pretraining loss and downstream task balanced accuracy indicate that larger models generally achieve higher accuracy.

E.2 For Data Size

In Figures 8a, 8b, 8c, we provide the results of the scale law experiments on the pretraining loss function (\mathcal{L}_p), SEED, and TUSL dataset. The results on the pretraining loss function show that the

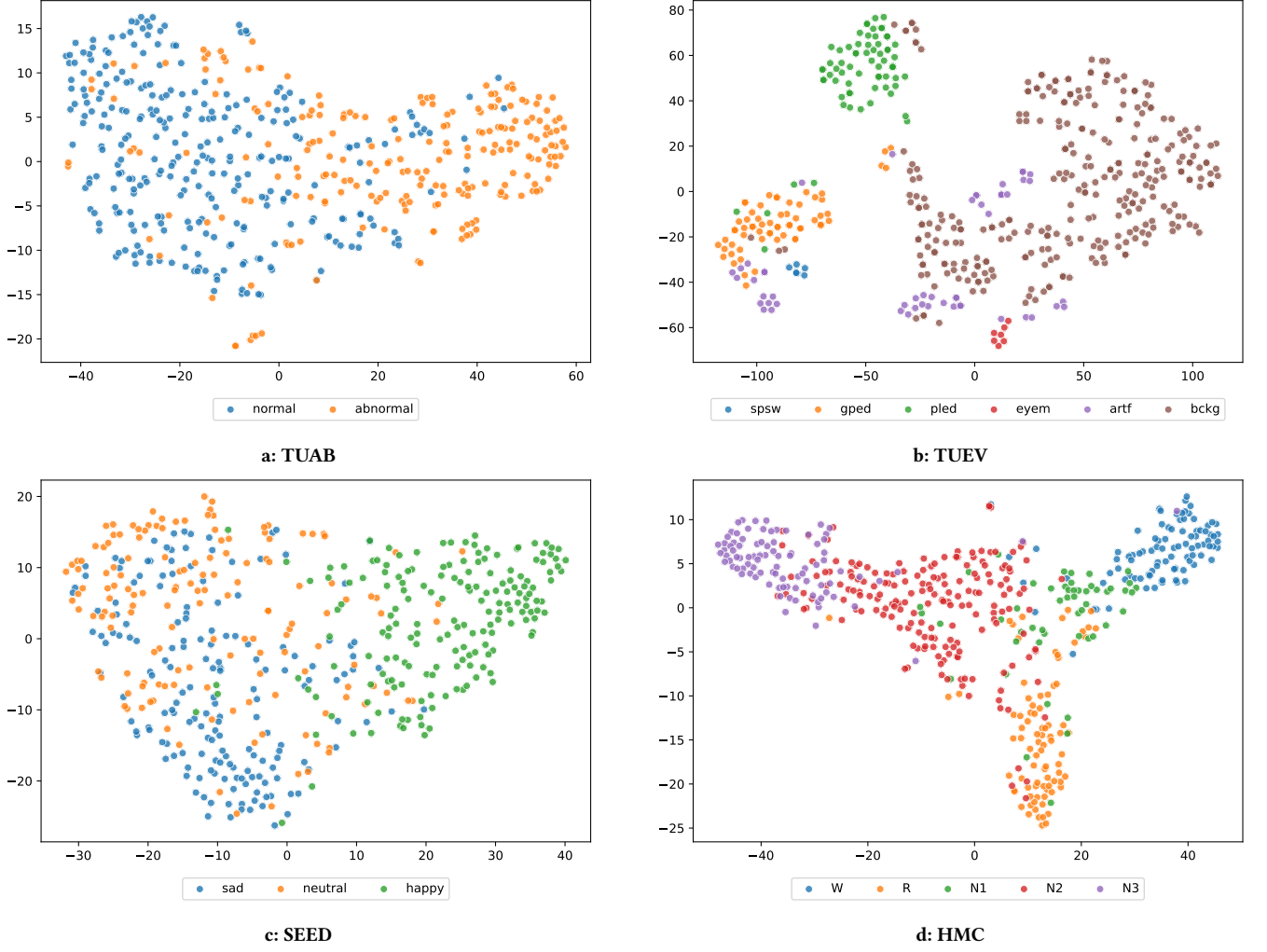


Figure 6: t-SNE visualizations of feature embeddings reduced to 40 dimensions by PCA on different datasets with ALFEE-B. t-SNE runs for 1500 iterations in which the perplexity is 30 and the number of samples is 500.

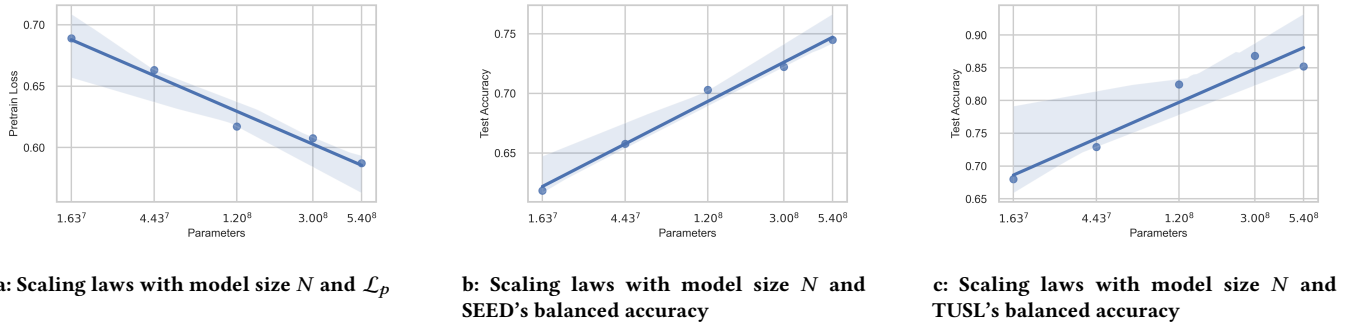


Figure 7: Scaling laws with model size N and a) \mathcal{L}_p ; b) SEED balanced accuracy; c) TUSL balanced accuracy. Axes are all on a logarithmic scale.

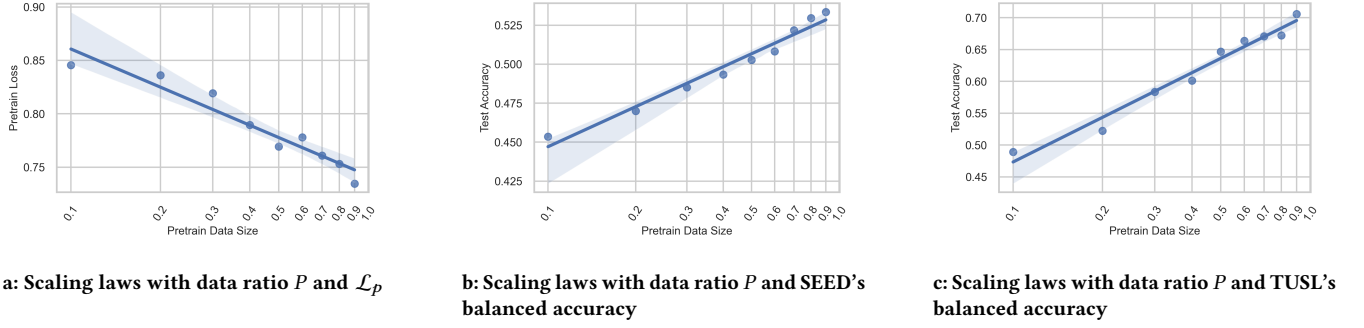


Figure 8: Scaling laws with data partition P and a) \mathcal{L}_p ; b) SEED balanced accuracy; c) TUSL balanced accuracy. Axes are all on a logarithmic scale.

scaling law of the test \mathcal{L}_p with data ratio (P) is: $\mathcal{L}_p = -0.051 * \ln(P) + 0.742$, where R^2 is 0.923. The results on the SEED dataset show that the scaling law of the test balanced accuracy with data ratio (P) is: $B\text{Acc} = 0.037 * \ln(P) + 0.532$, where R^2 is 0.968. The results on the TUSL dataset show that the scaling law of the test balanced accuracy with data ratio (P) is: $B\text{Acc} = 0.101 * \ln(P) + 0.706$, where R^2 is 0.971. The results of the pretraining loss and downstream task balanced accuracy indicate that more pretraining data generally achieve higher accuracy.

F Limitations and Future Work

First, although we pretrain the ALFEE with up to 540M parameters using 25,000 hours of EEG data, the overall data scale and model capacity remain limited compared to those of contemporary language and vision models, and the variable quality and low signal-to-noise ratio of EEG data necessitate new preprocessing techniques for more robust representation learning. Second, while the cross-attention layer combined with task-specific CLS tokens enables excellent multi-task performance under full-parameter fine-tuning, integrating new tasks still requires additional fine-tuning; developing a more comprehensive and universal model similar to those in natural language processing remains an open challenge. Finally, our current work primarily focuses on robust EEG representation; further research is needed to extend the model to multimodal scenarios such as Video-EEG, Image-EEG, Text-EEG, and fMRI-EEG for broader applicability.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. GPT-4 technical report. *OpenAI Blog* (2023).
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. In *NeurIPS*. 23716–23736.
- [3] Mashail Alsolamy and Anas Fattouh. 2016. Emotion estimation from EEG signals during listening to Quran using PSD features. In *CSIT*. 1–5.
- [4] Diego Alvarez-Estevéz and Roselyne M. Rijsman. 2021. Inter-database validation of a deep learning approach for automatic sleep scoring. *PLOS ONE* 16, 8 (08 2021), 1–27.
- [5] David Bethge, Philipp Hallgarten, Tobias Grosse-Puppenthal, Mohamed Kari, Lewis L Chuang, Ozan Özdenizci, and Albrecht Schmidt. 2022. EEG2Vec: Learning affective EEG representations via variational autoencoders. In *SMC*. 3150–3157.
- [6] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. 2023. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf> 2, 3 (2023), 8.
- [7] Benjamin Blankertz, Guido Dornhege, Matthias Krauledat, Klaus-Robert Müller, and Gabriel Curio. 2007. The non-invasive Berlin Brain–Computer Interface: Fast acquisition of effective performance in untrained subjects. *NeuroImage* 37, 2 (2007), 539–550.
- [8] Duo Chen, Suiren Wan, Jing Xiang, and Forrest Sheng Bao. 2017. A high-performance seizure detection algorithm based on Discrete Wavelet Transform (DWT) and EEG. *PLOS One* 12, 3 (2017).
- [9] Peng Chen, Yingying Zhang, Yunyao Cheng, Yang Shu, Yihang Wang, Qingsong Wen, Bin Yang, and Chenjuan Guo. 2024. Pathformer: Multi-scale transformers with adaptive pathways for time series forecasting.
- [10] Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, Wanxiang Che, Xiangzhan Yu, and Furu Wei. 2023. BEATS: audio pre-training with acoustic tokenizers. In *ICML*.
- [11] Andac Demir, Toshiaki Koike-Akino, Ye Wang, Masaki Haruna, and Deniz Erdogmus. 2021. EEG-GNN: Graph neural networks for classification of electroencephalogram (EEG) signals. In *EMBC*. 1061–1067.
- [12] Paolo Detti, Giampaolo Vatti, and Garazi Zabalo Manrique de Lara. 2020. EEG Synchronization Analysis for Seizure Prediction: A Study on Data of Noninvasive Recordings. *Processes* 8, 7 (2020).
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. 4171–4186.
- [14] Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. Unified language model pre-training for natural language understanding and generation. In *NeurIPS*.
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale.
- [16] Xiaobing Du, Cuixia Ma, Guanhua Zhang, Jinyao Li, Yu-Kun Lai, Guozhen Zhao, Xiaoming Deng, Yong-Jin Liu, and Hongan Wang. 2020. An efficient LSTM network for emotion recognition from multichannel EEG signals. *IEEE Transactions on Affective Computing* 13, 3 (2020), 1528–1540.
- [17] Guido Gainotti. 2019. The role of the right hemisphere in emotional and behavioral disorders of patients with frontotemporal lobar degeneration: an updated review. *Frontiers in aging neuroscience* 11 (2019), 55.
- [18] Shanghua Gao, Teddy Koker, Owen Queen, Thomas Hartvigsen, Theodoros Tsiligkaridis, and Marinka Zitnik. 2024. UniTS: Building a Unified Time Series Model. *NeurIPS*.
- [19] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. 2023. Imagebind: One embedding space to bind them all. In *CVPR*. 15180–15190.
- [20] Peiliang Gong, Ziyu Jia, Pengpai Wang, Yueying Zhou, and Daoqiang Zhang. 2023. ASTDF-net: attention-based spatial-temporal dual-stream fusion network for EEG-based emotion recognition. In *ACM MM*. 883–892.
- [21] Alexandre Gramfort, Martin Luessi, Eric Larson, Denis A. Engemann, Daniel Strohmeier, Christian Brodbeck, Roman Goj, Mainak Jas, Teon Brooks, Lauri Parkkonen, and Matti Hämäläinen. 2013. MEG and EEG data analysis with MNE-Python. *Frontiers in Neuroscience* 7 (2013).
- [22] A. Harati, M. Golmohammadi, S. Lopez, I. Obeid, and J. Picone. 2015. Improved EEG event classification using differential energy. In *SPMB*. 1–4.

- [23] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *CVPR*. 16000–16009.
- [24] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *CVPR*. 9729–9738.
- [25] Ronghang Hu and Amanpreet Singh. 2021. Unit: Multimodal multitask learning with a unified transformer. In *CVPR*. 1439–1449.
- [26] A.J. Jerri. 1977. The Shannon sampling theorem—Its various extensions and applications: A tutorial review. *Proc. IEEE* 65, 11 (1977), 1565–1596.
- [27] Wei-Bang Jiang, Yu-Ting Lan, and Bao-Liang Lu. 2024. REmoNet: Reducing Emotional Label Noise via Multi-regularized Self-supervision. In *ACM MM*. 2204–2213.
- [28] Wei-Bang Jiang, Yansen Wang, Bao-Liang Lu, and Dongsheng Li. 2025. NeuroLM: A Universal Multi-task Foundation Model for Bridging the Gap between Language and EEG Signals. In *ICLR*.
- [29] Wei-Bang Jiang, Li-Ming Zhao, and Bao-Liang Lu. 2024. Large brain model for learning generic representations with tremendous EEG data in BCI. In *ICLR*.
- [30] Jin Jing, Wendong Ge, Shenda Hong, Marta Bento Fernandes, Zhen Lin, Chaoqi Yang, Sungtae An, Aaron F Struck, Aline Herlopian, Ioannis Katakis, et al. 2023. Development of expert-level classification of seizures and rhythmic and periodic patterns during EEG interpretation. *Neurology* 100, 17 (2023), e1750–e1762.
- [31] Louis Korczowski, Martine Cederhout, Anton Andreev, Grégoire Cattani, Pedro Luiz Coelho Rodrigues, Violette Gautheret, and Marco Congedo. 2019. *Brain Invaders calibration-less P300-based BCI with modulation of flash duration Dataset (bi2015a)*. Technical Report.
- [32] Demetres Kostas, Stephane Aroca-Ouellette, and Frank Rudzicz. 2021. BENDR: Using transformers and a contrastive self-supervised learning task to learn from massive amounts of EEG data. *FHN* 15 (2021), 653–659.
- [33] Surbhi Kumari and Amit Kumar Dutta. 2024. Brain waves, neuroimaging (fMRI, EEG, MEG, PET, NIR). In *Computational Techniques in Neuroscience*. 67–82.
- [34] Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and Brent J Lance. 2018. EEGNet: a compact convolutional neural network for EEG-based brain–computer interfaces. *Journal of Neural Engineering* 15, 5 (2018).
- [35] Hongli Li, Man Ding, Ronghua Zhang, and Chunbo Xiu. 2022. Motor imagery EEG classification algorithm based on CNN-LSTM feature fusion network. *Biomedical Signal Processing and Control* 72 (2022), 103342.
- [36] Rui Li, Yiting Wang, Wei-Long Zheng, and Bao-Liang Lu. 2022. A multi-view spectral-spatial-temporal masked autoencoder for decoding emotions with self-supervised learning. In *ACM MM*. 6–14.
- [37] Bingxiu Liu, Jifeng Guo, CL Philip Chen, Xia Wu, and Tong Zhang. 2023. Fine-grained interpretability for EEG emotion recognition: Concat-aided grad-CAM and systematic brain functional network. *IEEE Transactions on Affective Computing* 15, 2 (2023), 671–684.
- [38] Wei Liu, Jie-Lin Qiu, Wei-Long Zheng, and Bao-Liang Lu. 2022. Comparing Recognition Performance and Robustness of Multimodal Deep Learning Models for Multimodal Emotion Recognition. *IEEE TCDS* 14, 2 (2022), 715–729.
- [39] Wei Liu, Wei-Long Zheng, Ziyi Li, Si-Yuan Wu, Lu Gan, and Bao-Liang Lu. 2022. Identifying similarities and differences in emotion recognition with EEG and eye movements among Chinese, German, and French People. *Journal of Neural Engineering* 19, 2 (2022), 026012.
- [40] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019. Multi-task deep neural networks for natural language understanding. In *ACL*.
- [41] Matthew D. Luciw, Ewa Jarocka, and Benoni B. Edin. 2014. Multi-channel EEG recordings during 3,936 grasp and lift trials with varying weight and friction. *Scientific Data* 1, 1 (25 Nov 2014), 140047.
- [42] S. López, G. Suarez, D. Jungreis, I. Obeid, and J. Picone. 2015. Automated identification of abnormal adult EEGs. In *SPMB*. 1–5.
- [43] Ben Mann, N Ryder, M Subbiah, J Kaplan, P Dhariwal, A Neelakantan, P Shyam, G Sastry, A Askell, S Agarwal, et al. 2020. Language models are few-shot learners. *OpenAI Blog* (2020).
- [44] Perrin Margaux, Maby Emmanuel, Daligault Sébastien, Bertrand Olivier, and Mattout Jérémie. 2012. Objective and Subjective Evaluation of Online Error Correction during P300-Based Spelling. *Advances in Human-Computer Interaction* 2012, 1 (2012), 578295.
- [45] Jonathan Masci, Ueli Meier, Dan Cireşan, and Jürgen Schmidhuber. 2011. Stacked convolutional auto-encoders for hierarchical feature extraction. In *ICANN*. 52–59.
- [46] Paulius Micekevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. 2017. Mixed Precision Training. *arXiv preprint arXiv:1710.03740* (2017).
- [47] Navid Mohammadi Foumani, Geoffrey Mackellar, Soheila Ghane, Saad Irtza, Nam Nguyen, and Mahsa Salehi. 2024. EEG2Rep: Enhancing Self-supervised EEG Representation Through Informative Masked Inputs. In *KDD*. 5544–5555.
- [48] Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. 2023. A time series is worth 64 words: Long-term forecasting with transformers.
- [49] Iyad Obeid and Joseph Picone. 2016. The temple university hospital EEG data corpus. *Frontiers in Neuroscience* 10 (2016), 196.
- [50] Robert Oostenveld, Pascal Fries, Eric Maris, and Jan-Mathijs Schoffelen. 2011. FieldTrip: open source software for advanced analysis of MEG, EEG, and invasive electrophysiological data. *Computational Intelligence and Neuroscience* 2011, 1 (2011), 156869.
- [51] Artidoro Pagnoni, Ram Pasunuru, Pedro Rodriguez, John Nguyen, Benjamin Muller, Margaret Li, Chunting Zhou, Lili Yu, Jason Weston, Luke Zettlemoyer, et al. 2024. Byte Latent Transformer: Patches scale better than tokens. *arXiv preprint arXiv:2412.09871* (2024).
- [52] Luiz Pessoa. 2017. A network model of the emotional brain. *Trends in cognitive sciences* 21, 5 (2017), 357–371.
- [53] Subhojeet Pramanik, Priyanka Agrawal, and Aman Hussain. 2019. Omninet: A unified architecture for multi-modal multi-task learning. In *ICLR*.
- [54] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*. 8748–8763.
- [55] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. *OpenAI Blog* (2018).
- [56] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI Blog* (2019).
- [57] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR* 21, 140 (2020), 1–67.
- [58] Kashif Rasul, Arjun Ashok, Andrew Robert Williams, Hena Ghonia, Rishika Bhagwatkar, Arian Khorasani, Mohammad Javad Darvishi Bayazi, George Adamopoulos, Roland Riachi, Nadhir Hassen, et al. 2023. Lag-LLaMA: Towards foundation models for probabilistic time series forecasting. In *NeurIPS Workshop*.
- [59] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *CVPR*. 10684–10695.
- [60] Subhrajit Roy, Isabell Kiral-Kornek, and Stefan Harrer. 2019. ChronoNet: A Deep Recurrent Neural Network for Abnormal EEG Identification. In *Artificial Intelligence in Medicine*, David Riaño, Szymon Wilk, and Annette ten Teije (Eds.). 47–56.
- [61] Arman Savran, Koray Çiftçi, Guillaume Chanel, Javier Mota, Luong Viet, Bulent Sankur, Lale Akarun, Alice Caplier, and Michèle Rombaut. 2006. Emotion Detection in the Loop from Brain Signals and Facial Images. (01 2006).
- [62] G. Schalk, D.J. McFarland, T. Hinterberger, N. Birbaumer, and J.R. Wolpaw. 2004. BCI2000: a general-purpose brain-computer interface (BCI) system. *IEEE TBE* 51, 6 (2004), 1034–1043.
- [63] Robin Tibor Schirrmeyer, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin Glasstetter, Katharina Eggensperger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. 2017. Deep learning with convolutional neural networks for EEG decoding and visualization. *Human Brain Mapping* 38, 11 (2017), 5391–5420.
- [64] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*. 618–626.
- [65] Yonghao Song, Xueyu Jia, Lie Yang, and Longhan Xie. 2021. Transformer-based spatial-temporal feature learning for EEG decoding. *arXiv preprint arXiv:2106.11170* (2021).
- [66] Yonghao Song, Qingqing Zheng, Bingchuan Liu, and Xiaorong Gao. 2022. EEG Conformer: Convolutional transformer for EEG decoding and visualization. *TNSRE* 31 (2022), 710–719.
- [67] Zhenxi Song, Ruihan Qin, Huixia Ren, Zhen Liang, Yi Guo, Min Zhang, and Zhiguo Zhang. 2024. EEG-MACS: Manifold Attention and Confidence Stratification for EEG-based Cross-Center Brain Disease Diagnosis under Unreliable Annotations. In *ACM MM*. 340–349.
- [68] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing* 568 (2024), 127063.
- [69] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander Alemi. 2017. Inception-v4, inception-resnet and the impact of residual connections on learning. In *AAAI*.
- [70] Mastaneh Torkamani-Azar, Sumeysra Demir Kanik, Serap Aydin, and Mujdat Cetin. 2020. Prediction of Reaction Time and Vigilance Variability From Spatio-Spectral Features of Resting-State EEG in a Long Sustained Attention Task. *IEEE JBHI* 24, 9 (2020), 2550–2558.
- [71] Logan T. Trujillo. 2019. Mental Effort and Information-Processing Costs Are Inversely Related to Global Brain Free Energy During Visual Categorization. *Frontiers in Neuroscience* 13 (2019).
- [72] Logan T. Trujillo, Candice T. Stanfield, and Ruben D. Vela. 2017. The Effect of Electroencephalogram (EEG) Reference Choice on Information-Theoretic Measures of the Complexity and Integration of EEG Signals. *Frontiers in Neuroscience*

- 11 (2017).
- [73] Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural discrete representation learning. In *NeurIPS*.
 - [74] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning Research* 9, 86 (2008), 2579–2605. <http://jmlr.org/papers/v9/vandermaaten08a.html>
 - [75] E. von Weltin, T. Ahsan, V. Shah, D. Jamshed, M. Golmohammadi, I. Obeid, and J. Picone. 2017. Electroencephalographic slowing: A primary source of error in automatic seizure detection. In *SPMB*. 1–5.
 - [76] Christopher Wang, Vighnesh Subramaniam, Adam Uri Yaari, Gabriel Kreiman, Boris Katz, Ignacio Cases, and Andrei Barbu. 2023. BrainBERT: Self-supervised representation learning for intracranial recordings. In *ICLR*.
 - [77] Guangyu Wang, Wenchao Liu, Yuhong He, Cong Xu, Lin Ma, and Haifeng Li. 2024. EEGPT: Pretrained transformer for universal and reliable representation of eeg signals. In *NeurIPS*. 39249–39280.
 - [78] Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sahoo. 2024. Unified training of universal time series forecasting transformers. In *ICML*.
 - [79] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2024. NExT-GPT: Any-to-Any Multimodal LLM. In *ICML*.
 - [80] Xueyuan Xu, Li Zhuo, Jinxin Lu, and Xia Wu. 2024. WSEL: EEG feature selection with weighted self-expression learning for incomplete multi-dimensional emotion recognition. In *ACM MM*. 350–359.
 - [81] Chaoqi Yang, M Westover, and Jimeng Sun. 2023. BIOT: Biosignal transformer for cross-data learning in the wild. *NeurIPS*, 78240–78260.
 - [82] Chaoqi Yang, Cao Xiao, M Brandon Westover, Jimeng Sun, et al. 2023. Self-supervised electroencephalogram representation learning for automatic sleep staging: model development and evaluation study. *JMIR AI* 2, 1 (2023), e46769.
 - [83] Zhizhang Yuan, Daoze Zhang, Junru Chen, Geifei Gu, and Yang Yang. 2024. Brant-2: Foundation model for brain signals. *arXiv e-prints* (2024), arXiv–2402.
 - [84] Daoze Zhang, Zhizhang Yuan, Junru Chen, Kerui Chen, and Yang Yang. 2024. Brant-X: A Unified Physiological Signal Alignment Framework. In *KDD*. 4155–4166.
 - [85] Daoze Zhang, Zhizhang Yuan, Yang Yang, Junru Chen, Jingjing Wang, and Yafeng Li. 2023. Brant: Foundation model for intracranial neural signal. In *NeurIPS*. 26304–26321.
 - [86] Yunhao Zhang and Junchi Yan. 2023. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In *ICLR*.
 - [87] W. Zheng, W. Liu, Y. Lu, B. Lu, and A. Cichocki. 2018. EmotionMeter: A Multimodal Framework for Recognizing Human Emotions. *IEEE TC* (2018), 1–13.
 - [88] Wei-Long Zheng and Bao-Liang Lu. 2015. Investigating Critical Frequency Bands and Channels for EEG-Based Emotion Recognition with Deep Neural Networks. *IEEE TAMD* 7, 3 (2015), 162–175.
 - [89] Igor Zyma, Sergii Tukaev, Ivan Seleznev, Ken Kiyono, Anton Popov, Mariia Chernykh, and Oleksii Shpenkov. 2019. Electroencephalograms during Mental Arithmetic Task Performance. *Data* 4, 1 (2019).