# CEReBrO: Compact Encoder for Representations of Brain Oscillations Using Efficient Alternating Attention

**Alexandru Dimofte**[1][*] **Glenn Anta Bucagu**[1] **Thorir Mar Ingolfsson**[2] **Xiaying Wang**[2]

**Andrea Cossettini**[2] **Luca Benini**[2] **Yawei Li**[2][†]

[1]University of Zurich    [2]Integrated Systems Laboratory, ETH Zürich, Switzerland

## Abstract

Electroencephalography (EEG) is indispensable in neurology and brain–computer interfaces, yet learning from it is hampered by scarce labels, diverse montages, and long, high-dimensional time series. Recent self-supervised, attention-based EEG models mitigate the first issue. However, two challenges remain: (i) electrode-specific geometry and referencing schemes, (ii) quadratic memory and time complexity of self-attention limiting the processing of long high channel-count signals.

To address these challenges, we introduce **CEReBrO** – a **C**ompact **E**ncoder for **Re**presentations of **Br**ain **O**scillations. CEReBrO augments each token with an electrode position-aware channel embeddings in addition to the standard learnable temporal embeddings, and processes the sequence with a novel alternating attention Transformer approach that interleaves intra-channel (temporal) and inter-channel (spatial) attention. This design cuts memory usage by up to $6\times$ and inference latency by $2\times$ compared with vanilla self-attention, without degrading performance.

We pre-train CEReBrO on 25,000 hours of public clinical and commercial EEG data from >12,000 subjects spanning 2-64 channels. On several publicly available downstream benchmarks covering clinical and BCI tasks, CEReBrO is competitive with existing state of the art EEG foundation models, but with a lower pre-training and fine-tuning cost.

## 1 Introduction

An Electroencephalogram (EEG) records the brain's spontaneous electrical activity. Its extensive applications cover disease diagnosis, medical monitoring, and brain-computer interfacing among many others [1, 2]. Modeling EEG is challenging: the signals are inherently non-linear, correlated, and non-stationary [3], making classical time series models inadequate [4]. Compounding this difficulty is the scarcity of labeled EEG data: obtaining expert annotations is time-consuming, costly, and prone to human errors [5, 6, 7, 8]. For example, continuous multi-day EEG monitoring of epileptic patients demands substantial manual effort and can yield high inter-rater variability. Additionally, artifacts (e.g. ocular, cardiac, and muscular activity) can mimic seizure patterns and lead to false alarms and misinterpretations [9, 10, 11, 12].

---

[*]Glenn Anta Bucagu, Alexandru Dimofte and Thorir Mar Ingolfsson contributed equally to this work.

[†]Correspondence to `yawli@iis.ee.ethz.ch`.

To alleviate the limitations of small labeled datasets, self-supervised EEG foundation models have recently emerged (see Section 2). They leverage large-scale, unlabeled EEG signals and can be fine-tuned on downstream tasks. While promising, they face two key challenges.

**1) Modeling of spatial and temporal characteristics of EEG signals**. EEG signals are distinct from text or images, exhibiting unique spatial and temporal dynamics. Preserving critical invariances while capturing these dynamics is non-trivial. Although self-attention has shown promise [13, 14, 15, 16, 17, 18], it can be both computationally expensive and, at times, suboptimal for EEG-specific structure [19].

**2) Handling variations in channel configurations**. EEG signals can be unipolar (each channel relative to a reference) or bipolar (difference between two neighboring electrodes), leading to diverse setups. Existing models often overlook inter-channel correlations or assume fixed channel indexing [17], with limited generalization to different montages. While CBraMod [19] uses asymmetric conditional positional encodings, it still relies on fixed electrode ordering and does not handle referencing differences.

To address these gaps, we present **CEReBrO**, a *Compact Encoder for Representations of Brain Oscillations* with an alternating attention strategy and position-aware channel embeddings. Our contributions are as follows:

- **Large-scale pre-training**: Leveraging over 25,000 hours of public EEG from 12,000+ subjects, we ensure coverage of diverse montages, pathologies, and hardware. CEReBrO is exposed to configurations ranging from 2 to 64 channels, enabling robust performance across both consumer-grade and high-density clinical setups.
- **Electrode Position-Aware Channel Embeddings**: To our knowledge, we are the first to explicitly encode the spatial location of an electrode on the scalp into the model. We embed 3D scalp coordinates via an MLP, achieving invariance to channel count, order, referencing scheme, or montage. This flexible design enables CEReBrO to handle both unipolar and bipolar configurations seamlessly.
- **Novel Alternating Attention Mechanism**: We introduce a dual strategy to separately model intra-channel (temporal) and inter-channel (spatial) dependencies in alternating layers. This reduces memory usage up to $6\times$ and runtime by $2\times$ compared to standard self-attention.
- **Comprehensive Evaluations:** CEReBrO achieves competitive performance on 11 public benchmarks, including anomaly classification, seizure detection, emotion recognition, and sleep staging, while drastically reducing computational costs.

## 2 Related Works

Recent advancements in EEG foundation models have leveraged Transformer architectures and self-supervised learning techniques to address the scarcity of labeled EEG data. These models aim to learn robust representations from raw EEG signals, facilitating various downstream tasks without extensive labeled datasets.

BENDR [20] is inspired by wav2vec [21] and employs a stack of short-receptive-field 1D convolutions to transform raw EEG waveforms into a sequence of embeddings, which are fed to a Transformer encoder with linear attention modules. Pre-training is done via a contrastive learning objective. However, without explicit channel-specific embeddings, BENDR may not optimally differentiate between channels, especially when the number of channels varies between training examples. BrainBERT [13] uses iEEG spectrogram patches as input tokens to a standard Transformer encoder model, pre-trained via MAE. By processing channels individually, BrainBERT fails to capture inter-channel correlations, missing crucial spatial relationships inherent in multi-channel EEG data. BIOT [16] frames EEG as a "sentence" of fixed-length channel tokens with index-based and positional embeddings to encode spatio-temporal context, enabling cross-dataset transfer despite mismatched channels or missing data. Its linear attention design trades off some accuracy for speed, and its index-only embeddings remain sensitive to channel ordering and reference variations. Neuro-GPT [14] adapts causal auto-regressive MAE for EEG waveform modeling. In this model, each token aggregates information from multiple channels, limiting the attention mechanism's ability to balance spatio-temporal characteristics. Neuro-GPT standardizes to 22 channels via nearest neighbor interpolation, which can degrade performance by losing detail from high channel counts and introducing artifacts when interpolating from fewer
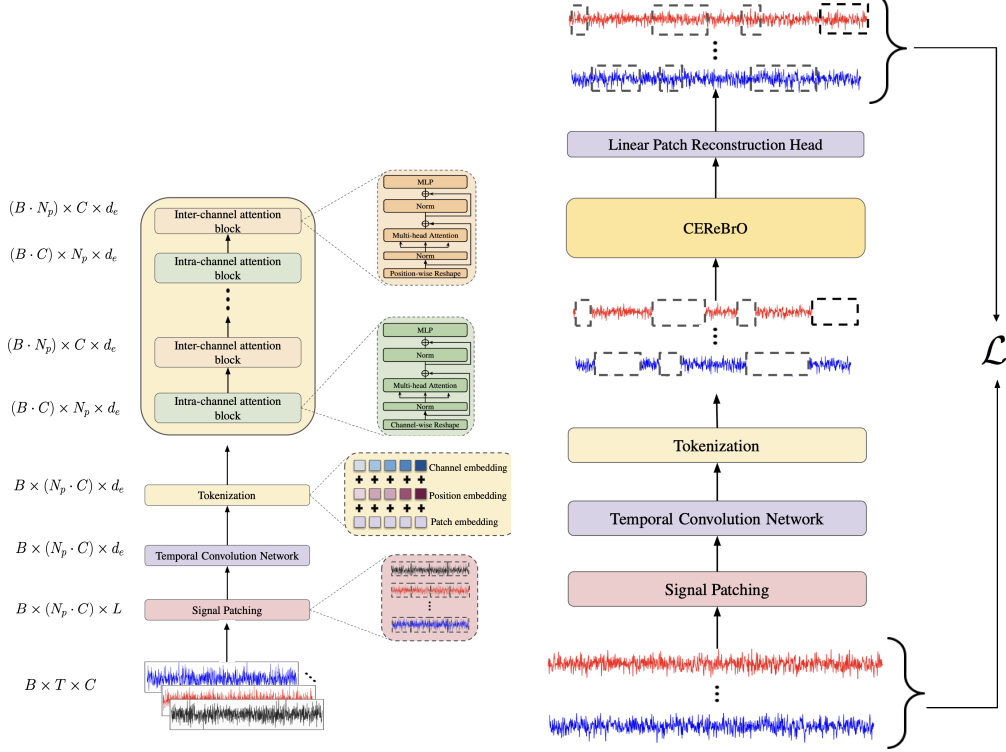
2

Figure 1: (a) Overview of the CEReBrO architecture    (b) Overview of our pre-training framework.

channels. LaBraM [17] introduces a learned neural tokenizer that maps patches of EEG waveforms to discrete codebook embeddings, which are then processed by a Transformer encoder within a symmetric MAE framework. The primary limitation of LaBraM is its pre-training which requires 2 stages: one stage to train the neural tokenizer and another stage to pre-train the foundation model. Similarly, EEGFormer [15] incorporates a vector quantizer within a Transformer-based model trained via autoencoding, adding complexity to the training process. CBraMod [19] leverages an innovative asymmetric conditional positional encoding (ACPE) that uses kernel operations to learn electrode configurations under the assumption of a fixed channel ordering, paired with a criss-cross attention mechanism that encodes spatial and temporal features in parallel. However, by pre-training solely on 19 standard channels from the Temple University EEG Corpus (TUEG) and assuming a fixed channel order, CBraMod's applicability to low-channel setups or alternative montages not seen during pre-training is limited.

## 3 Method

**CEReBrO converts raw EEG into position-aware tokens, processes them with an alternating-attention Transformer, and learns through a masked-autoencoding objective.** Fig. 1 summarizes the pipeline. We detail each component below.

### 3.1 Patching

**Localizing the signal in time improves frequency resolution and keeps sequence length manageable.** Following [22], we slice EEG waveforms into equally-sized non-overlapping patches to (i) enhance the locality (ii) reduce computation and memory, and (iii) attend to longer temporal dependencies.

Given a multi-channel EEG waveform $\mathbf{X} \in \mathbb{R}^{T \times C}$, where $T$ is the number of timesteps and $C$ is the number of channels, we segment $\mathbf{X}$ into non-overlapping patches of length $L$ with stride $S$. This results in a set of patches $\mathbf{P} \in \mathbb{R}^{N_p \times C \times L}$, where $N_p = \left\lfloor \frac{T-L}{S} + 1 \right\rfloor$ is the number of patches per channel.

## 3.2 Temporal Feature Extraction

Each patch $\mathbf{P}_{c,i} \in \mathbb{R}^L$, corresponding to channel $c$ and patch index $i$, is first processed by a single temporal encoder that performs hierarchical temporal filtering. Specifically, three successive 1D convolutions, each followed by GroupNorm and GELU, extract features at increasing levels of temporal granularity. The first layer captures coarse, broad-scale temporal patterns and downsamples the sequence, while the next two layers refine fine-scale oscillatory structure. The final feature map is flattened into

$$\mathbf{H}_{c,i} \in \mathbb{R}^{d_h},$$

which an MLP then projects to the patch embedding

$$\mathbf{E}_{c,i} \in \mathbb{R}^{N_p \times d_e}.$$

This design yields compact, hierarchical temporal representations and stable convergence dynamics with relatively few trainable parameters. It provides an efficient way to extract features from EEG signals, which are known to have high temporal resolution [17].

## 3.3 Electrode Position-Aware Channel Embeddings

**To generalize across montages we encode physical 3D positions rather than index IDs.** Each channel $c$ consists of two components: one representing an individual electrode and one acting either as an electrode (e.g. in a bipolar montage) or as a meta-electrode (e.g. in a mean referenced unipolar montage). Let $\mathbf{p}_1, \mathbf{p}_2 \in \mathbb{R}^{N_p \times 3}$ denote the normalized 3D Cartesian coordinates of the two electrodes that make up each channel. We used a map of polar and azimuthal angles for the 118 electrodes defined in the 10/10 and 10/20 systems [23] from which we can compute normalized 3D Cartesian coordinates for real electrodes. In the case of a meta-electrode, such as the mean of all active electrodes, $\mathbf{p}_2$ is the centroid of all available electrodes in a given training example. In the case of linked ear referencing, $\mathbf{p}_2$ is the mean of the M1 and M2 electrode coordinates. Similar approaches are used for other referencing styles.

We pass $\mathbf{p}_1$ and $\mathbf{p}_2$ through an MLP to obtain electrode vectors $\mathbf{e}_1$ and $\mathbf{e}_2 \in \mathbb{R}^{C \times \frac{d_e}{2}}$. These are concatenated to obtain channel embeddings $\mathbf{W}_{\text{chan},c} = \mathbf{e}_1 \oplus \mathbf{e}_2 \in \mathbb{R}^{C \times d_e}$. During pre-training, the MLP learns position-aware channel embeddings that are invariant to channel count, order, reference scheme, and montage. In contrast, ACPE (CBraMod) and absolute index encodings (LaBraM, BIOT) assume a fixed channel ordering and cannot discern between different reference schemes.

## 3.4 Temporal Embeddings

**Learnable temporal embeddings allow the model distinguish patches that are from the same channel but occur at different timesteps.** Temporal order between patches is captured using learnable positional embeddings $\mathbf{W}_{\text{pos}} \in \mathbb{R}^{N_p \times d_e}$, where $\mathbf{W}_{\text{pos},i} \in \mathbb{R}^{d_e}$ is the embedding corresponding to patch index $i$ in each channel. The final input representation to the encoder is given by:

$$\mathbf{E}_{c,i}^{\text{in}} = \mathbf{E}_{c,i} + \mathbf{W}_{\text{pos},i} + \mathbf{W}_{\text{chan},c}.$$

## 3.5 Alternating Attention

We employ a $N$-layer Transformer encoder [24] with $H$ attention heads and embedding dimension $d_e$. Fine-grained tokenization yields thousands of tokens per training sample, increasing the computation complexity of the self-attention mechanism that scales quadratically with sequence length. To mitigate this, we introduce *alternating attention* (Alg. 1), which switches between intra-channel (temporal) and inter-channel (spatial) heads on successive layers.

In intra-channel attention layers, we focus on temporal dependencies within each channel by computing attention over the sequence of patches in that channel. In inter-channel attention layers, we capture spatial correlations by computing attention across channels at each time step. We use an even number of encoder layers so that both intra-channel and inter-channel attention are equally represented. The inputs to the first layer are the previously introduced input embeddings $\mathbf{E}_{c,i}^{\text{in}}$. Each layer alternates between intra-channel and inter-channel multi-head attention, applies a feed-forward network, and is followed by layer normalization. The outputs $\mathbf{E}_{c,i}^{(j)}$ of the $j$-th layer serve as the inputs to the $(j+1)$-th layer.

**Algorithm 1** ALTERNATING ATTENTION MECHANISM

---

**Require:** Input tensor $\mathbf{T}$ of shape $[B, C \times N_p, d_e]$
    **Parameters:**
      $B$: batch size
      $C$: number of channels
      $N_p$: number of patches per channel
      $d_e$: embedding dimension
  1: **for** each encoder layer $i$ to $N$ **do**
  2:    **if** $i$ is odd **then**                                         $\triangleright$ Inter-channel attention
  3:        Reshape $\mathbf{T}$ to $[B \times N_p, C, d_e]$
  4:        Compute QKV projection
  5:        Multi-head attention over $C$ (channels)
  6:    **else**                                               $\triangleright$ Intra-channel attention
  7:        Reshape $\mathbf{T}$ to $[B \times C, N_p, d_e]$
  8:        Compute QKV projection
  9:        Multi-head attention over $N_p$ (patches)
 10:    **end if**
 11:    Reshape output back to $[B, C \times N_p, d_e]$
 12: **end for**
**Ensure:** Output tensor of shape $[B, C \times N_p, d_e]$

---

Table 1: Theoretical memory and time complexities of each attention type.

| Attention Type | Memory Complexity | Time Complexity |
|---|---|---|
| Intra-channel | $\mathcal{O}\left(CN_p^2\right)$ | $\mathcal{O}\left(CN_p^2 d_e\right)$ |
| Inter-channel | $\mathcal{O}\left(C^2 N_p\right)$ | $\mathcal{O}\left(C^2 N_p d_e\right)$ |
| Standard Self-Attention | $\mathcal{O}\left((CN_p)^2\right)$ | $\mathcal{O}\left((CN_p)^2 d_e\right)$ |

### 3.6 Complexity of Alternating Attention

For standard self-attention, the memory complexity is quadratic in sequence length because every token attends to every other. With alternating attention the cost depends on the layer type:

- **Intra-channel** layers attend within each channel over $N_p$ patches: $\mathcal{O}(CN_p^2)$.

- **Inter-channel** layers attend across $C$ channels at each time step: $\mathcal{O}(C^2 N_p)$

Table 1 summarizes the resulting memory and time complexity.

Empirical measurements are reported in Sect.4.2.

Our alternating attention mechanism is notable for two reasons: i) it allows us to jointly model intra-channel temporal dynamics and inter-channel spatial correlations within a single encoder, and ii) it enables efficient attention computation over long EEG token sequences. We illustrate our end-to-end pipeline in Fig. 1.

### 3.7 Lightweight Decoder

During both pre-training and fine-tuning, the output of the final encoder layer is passed to a single linear layer. In pre-training, this linear layer serves as a *patch-reconstruction head*; during fine-tuning, it becomes a classification head or regression head (depending on the fine-tuning task) that contains a few thousand trainable parameters. This is consistent with the foundation model literature [25, 13, 26, 17] In contrast, CBraMod uses a large Multi-Layer Perceptron (MLP) that contains 10M - 120M learnable parameters depending on the downstream task (which is orders of magnitude larger than the pre-trained encoder). Having such a large complex head makes fine-tuning memory and time-prohibitive, increases the chances of overfitting, may lead to inefficient parameter utilization or to the pre-trained encoder becoming a bottleneck [27].

## 3.8 Pre-training Objective

**We train CEReBrO with a masked auto-encoding loss that balances masked and visible tokens.**
We use MAE during pre-training to allow our model to learn useful representations from a large
unlabeled pre-training corpus. After patching of EEG waveforms, we randomly mask a fixed portion
of patches in each input sequence of patches $\mathbf{P}$. Masked positions are replaced with a single shared
learnable [MASK] token. The resulting sequence of tokens is passed to our Transformer encoder. The
encoder output is then passed to a linear layer that outputs a sequence $\hat{\mathbf{P}}$, which is a reconstruction of
$\mathbf{P}$. We define the following loss components:

$$\mathcal{L}_{\text{masked}} = \frac{1}{|\mathcal{M}|} \sum_{(c,i) \in \mathcal{M}} \|\mathbf{P}_{c,i} - \hat{\mathbf{P}}_{c,i}\|_2^2, \tag{1}$$

$$\mathcal{L}_{\text{visible}} = \frac{1}{|\overline{\mathcal{M}}|} \sum_{(c,i) \in \overline{\mathcal{M}}} \|\mathbf{P}_{c,i} - \hat{\mathbf{P}}_{c,i}\|_2^2, \tag{2}$$

where $\mathcal{M}$ and $\overline{\mathcal{M}}$ are the set of masked and visible token positions respectively. Our loss function
during pre-training is:

$$\mathcal{L} = \mathcal{L}_{\text{masked}} + \alpha \mathcal{L}_{\text{visible}} \tag{3}$$

Traditional MAE methods [28, 29, 26, 13, 25] compute the reconstruction loss only on masked patches
(i.e., $\alpha = 0$). We observed that doing so with our framework led to high-quality reconstructions of
masked patches but relatively poor and unpredictable reconstructions of visible patches. Inspired
by findings in[28], which showed that computing the loss uniformly on all tokens (i.e., $\alpha = 1$)
can reduce downstream performance, we introduce $\mathcal{L}_{\text{visible}}$ with a small weighting factor $\alpha = 0.1$.
Setting $\alpha = 0.1$ stabilizes visible-token reconstruction without collapsing to the identity solution,
echoing findings for lightweight vision transformers [30], demonstrating improved performance on
low-resolution image datasets.

## 4 Experiments

### 4.1 Pre-training

**Pre-training Corpus**     Our pre-training corpus includes the Temple University Hospital EEG Corpus
(TUEG) [31], with Temple University Abnormal Corpus (TUAB) subjects excluded to prevent
downstream contamination, complemented by the SEED series (SEED-IV, SEED-GER, SEED-FRA)
[32, 33, 34], BOAS [35], SleepEDFx[36], BCI-NER [37] and GWD [38, 39]. This results in over
25,000 hours of EEG recordings, which is $8\times$ the size of the corpus used by LaBraM and $2\times$ larger
than that of CBraMod. Notably, while CBraMod is limited to a fixed 19-channel configuration from
clinical TUEG data, our pre-training corpus includes data from various consumer-grade devices,
including headbands (BOAS), gel-based, wet and dry caps (GWD) with channel counts ranging from
2 to 64. In contrast to LaBraM, for which 14% of the pre-training corpus is private, CEReBrO's
pre-training corpus is 100% publicly accessible. More details are available in Section A.

### 4.2 Efficiency benchmarks

We benchmark memory footprint and forward-pass latency for alternating attention versus standard
self-attention. We sweep sequence length by varying the channel count $C \in \{1, \ldots, 64\}$ with
$N_p = 64$ patches per channel and test three model sizes ranging from 2.4M to 85M parameters. All
measurements are averaged over 100 forward passes on four NVIDIA RTX 2080 Ti GPUs.

Figure 2b shows that, at the largest sequence lengths, alternating attention is more than $2\times$ faster than
standard self-attention; Figure 2a confirms a $6\times$ reduction in peak GPU memory. These empirical
results align with the theoretical analysis in Table 1 and demonstrate that CEReBrO can train on high
channel count, long-duration EEG recordings that are computationally expensive for conventional
self-attention models such as LaBraM. This can be particularly useful for purposes like sleep staging,

monitoring of movement disorders, psychological disorders and cerebrovascular diseases which require long-term analysis of high channel count EEG signals [40, 41].
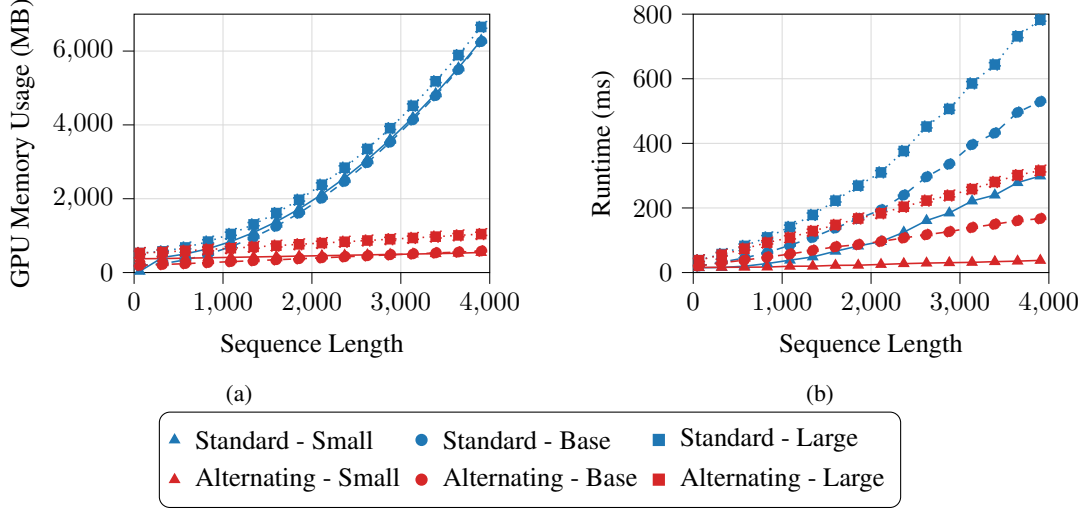


Figure 2: Forward pass GPU Memory Usage (a) and Runtime (b) vs. Sequence Length for alternating attention and standard self-attention in three CEReBrO model sizes. We use $N_p = 20$ and $C \in [1, 64]$ to simulate typical EEG configurations.

### 4.3 Downstream Experiments

We perform experiments on several downstream tasks to validate CEReBrO's effectiveness. Our main experiments are performed on 4 different downstream tasks, which are discussed in Section 4.4. We include several other downstream tasks in Section C.

**Baselines for Comparison**   We evaluate CEReBrO on all downstream tasks against both foundation and non-foundation-model baselines. Our non-foundation baselines include EEGNet [42], EEGConformer [43], SPaRCNet [44], ContraWR [45], CNN-Transformer [46], FFCL [47], and ST-Transformer [48], re-implemented using BIOT's public code [16] where necessary. Our foundation model baselines include BIOT [16], LaBraM [17] and CBraMod [19]. LaBraM has three available model sizes (Base, Large and Huge), but only LaBraM-Base's pre-trained weights are publicly available. We also provide the total number of trainable parameters during fine-tuning for each model. While CBraMod employs a 4M-parameter encoder, its approach involves flattening the transformer output and passing it to a classification or regression head comprising 10M-120M trainable parameters depending on the downstream task. In contrast, LaBraM, BIOT, and CEReBrO utilize lightweight classification heads with a few thousand trainable parameters.

**Performance Metrics**   For binary classification, we use Balanced Accuracy, AUC-PR, and AUROC. For multi-class classification, we evaluate using Balanced Accuracy, Cohen's Kappa, and Weighted F1. For regression, Pearson's Correlation, R2 Score, and RMSE are employed. All results are averaged over three random seeds and presented as mean ± standard deviation.

### 4.4 Results

**Anomaly Classification**   We evaluate our model on anomaly classification using the Temple University Hospital Abnormal EEG Corpus (TUAB) [31], consistent with existing models such as LaBraM and BIOT. TUAB contains EEG recordings labeled as either normal or abnormal, acquired with 23 channels at a 250 Hz sampling rate. For comparability across models, we adopt identical pre-processing steps and maintain the same train, validation, and test splits. Specifically, we use a bipolar montage containing 16 channels from the international 10-20 system. A band-pass filter (0.3–75 Hz) is applied to suppress low- and high-frequency noise, and a 60 Hz notch filter is used to eliminate power line interference. All EEG signals are subsequently resampled to 200 Hz and segmented into 409,455 samples of 10 seconds each. Results are shown in table Table 2.

Table 2: Performance of different methods in anomaly classification (TUAB) and sleep stage classification (ISRUC).

| Methods | Model Size | TUAB, 2-class | | | ISRUC, 5-class | | |
|---|---|---|---|---|---|---|---|
| | | Balanced Accuracy | AUPR | AUROC | Balanced Accuracy | Cohen's Kappa | Weighted F1 |
| EEGNet | 0.003M | $0.7642 \pm 0.0036$ | $0.8299 \pm 0.0043$ | $0.8412 \pm 0.0031$ | $0.7154 \pm 0.0121$ | $0.7040 \pm 0.0173$ | $0.7513 \pm 0.0124$ |
| EEGConformer | 0.55M | $0.7758 \pm 0.0049$ | $0.8427 \pm 0.0054$ | $0.8445 \pm 0.0038$ | $0.7400 \pm 0.0133$ | $0.7143 \pm 0.0162$ | $0.7634 \pm 0.0151$ |
| SPaRCNet | 0.79M | $0.7896 \pm 0.0018$ | $0.8414 \pm 0.0018$ | $0.8676 \pm 0.0012$ | $0.7487 \pm 0.0075$ | $0.7097 \pm 0.0132$ | $0.7624 \pm 0.0092$ |
| ContraWR | 1.6M | $0.7746 \pm 0.0041$ | $0.8421 \pm 0.0104$ | $0.8456 \pm 0.0074$ | $0.7402 \pm 0.0126$ | $0.7178 \pm 0.0156$ | $0.7610 \pm 0.0137$ |
| CNN-Transformer | 3.2M | $0.7777 \pm 0.0022$ | $0.8433 \pm 0.0039$ | $0.8461 \pm 0.0013$ | $0.7363 \pm 0.0087$ | $0.7129 \pm 0.0121$ | $0.7719 \pm 0.0105$ |
| FFCL | 2.4M | $0.7848 \pm 0.0038$ | $0.8448 \pm 0.0065$ | $0.8569 \pm 0.0051$ | $0.7277 \pm 0.0182$ | $0.7016 \pm 0.0291$ | $0.7614 \pm 0.0197$ |
| ST-Transformer | 3.5M | $0.7966 \pm 0.0023$ | $0.8521 \pm 0.0026$ | $0.8707 \pm 0.0019$ | $0.7381 \pm 0.0205$ | $0.7013 \pm 0.0352$ | $0.7681 \pm 0.0175$ |
| BIOT | 3.2M | $0.7959 \pm 0.0057$ | $0.8792 \pm 0.0023$ | $0.8815 \pm 0.0043$ | $0.7527 \pm 0.0121$ | $0.7192 \pm 0.0231$ | $0.7790 \pm 0.0146$ |
| LaBraM-Base | 5.8M | $0.8140 \pm 0.0019$ | $0.8965 \pm 0.0016$ | $0.9022 \pm 0.0009$ | $0.7633 \pm 0.0102$ | $0.7231 \pm 0.0182$ | $0.7810 \pm 0.0133$ |
| CBraMod | (68M, 4M) | $\mathbf{0.8249 \pm 0.0025}$ | $\mathbf{0.9221 \pm 0.0015}$ | $\mathbf{0.9156 \pm 0.0017}$ | $\mathbf{0.7865 \pm 0.0110}$ | $\mathbf{0.7442 \pm 0.0152}$ | $\mathbf{0.8011 \pm 0.0099}$ |
| CEReBrO | 2.4M | $0.8215 \pm 0.0027$ | $0.9096 \pm 0.0028$ | $0.9023 \pm 0.0020$ | $0.7312 \pm 0.0021$ | $0.6943 \pm 0.0114$ | $0.7525 \pm 0.0052$ |

Table 3: Performance of different methods in seizure detection (CHB-MIT) and motor imagery classification (SHU-MI).

| Methods | Model Size | CHB-MIT, 2-class | | | SHU-MI, 2-class | | |
|---|---|---|---|---|---|---|---|
| | | Balanced Accuracy | AUPR | AUROC | Balanced Accuracy | AUPR | AUROC |
| EEGNet | 0.003M | $0.5658 \pm 0.0106$ | $0.1914 \pm 0.0182$ | $0.8048 \pm 0.0136$ | $0.5889 \pm 0.0177$ | $0.6311 \pm 0.0142$ | $0.6283 \pm 0.0152$ |
| EEGConformer | 0.55M | $0.5976 \pm 0.0141$ | $0.2209 \pm 0.0215$ | $0.8226 \pm 0.0170$ | $0.5900 \pm 0.0107$ | $0.6370 \pm 0.0093$ | $0.6351 \pm 0.0101$ |
| SPaRCNet | 0.79M | $0.5876 \pm 0.0191$ | $0.1247 \pm 0.0119$ | $0.8143 \pm 0.0148$ | $0.5978 \pm 0.0097$ | $0.6510 \pm 0.0062$ | $0.6431 \pm 0.0082$ |
| ContraWR | 1.6M | $0.6344 \pm 0.0002$ | $0.2264 \pm 0.0174$ | $0.8097 \pm 0.0114$ | $0.5873 \pm 0.0128$ | $0.6315 \pm 0.0105$ | $0.6273 \pm 0.0113$ |
| CNN-Transformer | 3.2M | $0.6389 \pm 0.0067$ | $0.2479 \pm 0.0227$ | $0.8662 \pm 0.0082$ | $0.5975 \pm 0.0169$ | $0.6412 \pm 0.0076$ | $0.6323 \pm 0.0082$ |
| FFCL | 2.4M | $0.6262 \pm 0.0104$ | $0.2049 \pm 0.0346$ | $0.8271 \pm 0.0051$ | $0.5692 \pm 0.0252$ | $0.5943 \pm 0.0172$ | $0.6014 \pm 0.0168$ |
| ST-Transformer | 3.5M | $0.5915 \pm 0.0195$ | $0.1422 \pm 0.0094$ | $0.8237 \pm 0.0491$ | $0.5992 \pm 0.0206$ | $0.6394 \pm 0.0122$ | $0.6431 \pm 0.0111$ |
| BIOT | 3.2M | $0.7068 \pm 0.0457$ | $0.3277 \pm 0.0460$ | $0.8761 \pm 0.0284$ | $0.6179 \pm 0.0183$ | $0.6770 \pm 0.0119$ | $0.6609 \pm 0.0127$ |
| LaBraM-Base | 5.8M | $0.7075 \pm 0.0358$ | $0.3287 \pm 0.0402$ | $0.8679 \pm 0.0199$ | $0.6166 \pm 0.0192$ | $0.6761 \pm 0.0083$ | $0.6604 \pm 0.0091$ |
| CBraMod | (68M, 24M) | $\mathbf{0.7398 \pm 0.0284}$ | $\mathbf{0.3689 \pm 0.0382}$ | $\mathbf{0.8892 \pm 0.0154}$ | $\mathbf{0.6370 \pm 0.0151}$ | $\mathbf{0.7139 \pm 0.0088}$ | $\mathbf{0.6988 \pm 0.0068}$ |
| CEReBrO | 2.4M | $0.7221 \pm 0.0465$ | $0.3064 \pm 0.0161$ | $0.8704 \pm 0.0133$ | $0.6140 \pm 0.0138$ | $0.7043 \pm 0.0085$ | $0.6987 \pm 0.0149$ |

**Seizure Detection**    We evaluate our model on seizure detection using the CHB-MIT Scalp EEG Database, which contains recordings from 23 pediatric subjects with epileptic seizures. CHB-MIT contains EEG signals sampled at 256 Hz and labeled as seizure or non-seizure. Following BIOT's pre-processing pipeline, we selected 16 channels from a bipolar montage, resample all signals to 200 Hz, and segment the signals into 326,993 10-second samples. Subjects 1–19 are used for training, 20–21 for validation, and 22–23 for testing. Results are shown in table Table 3.

**Sleep Stage Classification**    We evaluate our model on sleep stage classification using the Subgroup-1 from the ISRUC-Sleep dataset [49], which consists of 100 all-night PSG recordings from 100 adult subjects. EEG signals are recorded from 6 channels at a 200 Hz sampling rate. All signals are segmented into 89,240 30-second epochs and annotated by sleep experts into five stages (Wake, N1, N2, N3, REM) following the AASM guidelines. Subjects 1–80 are used for training, 81–90 for validation, and 91–100 for testing. In line with standard practice, sleep staging is treated as a sequence-to-sequence classification task using sequences of 20 consecutive epochs to capture stage transitions. Following CBraMod, each model serves as a sample encoder, with a one-layer transformer applied as the sequence encoder. Results are shown in table Table 2.

**Motor Imagery Classification**    We evaluate our method on SHU-MI [50], which contains 32-channel EEG recordings from 25 subjects originally sampled at 250 Hz and labeled with left or right hand imagery. We segment the signals into 11,988 samples resampled at 200 Hz. Subjects 1–15, 16–20, and 21–25 are used for training, validation, and testing, respectively. Results are shown in table Table 3.

**The Role of Electrode Position-Aware Channel Embeddings**    We discuss the theoretical foundation for our electrode position-aware channel embeddings in Section 3. To validate the empirical validity of these channel embeddings, we compare them with 1) no channel embeddings and 2) learnable index-based channel embeddings (similar to those used in LaBraM and BIOT). Results are shown in Table 4. Generally, the worst performing model is the one with no channel embeddings. With no explicit channel embedding the model has to infer generalizable spatial correlations purely from signal correlations. This can be challenging especially for smaller downstream datasets. We observed that this model has a tendency to overfit and lead to poor performance on the test data. Index-based channel embeddings provide an improvement over no channel embeddings, but we do

not fully benefit from pre-trained representations because the channels seen in pre-training are not necessarily the same (or in the same order) as those seen downstream. Electrode position-aware embeddings offer the best performance, because they inject a notion of inter-channel distance, allowing CEReBrO to model physiologically meaningful correlations.

| Channel Embedding Type | TUAB, 2-class | | | ISRUC, 5-class | | |
|---|---|---|---|---|---|---|
| | Balanced Accuracy | AUPR | AUROC | Balanced Accuracy | Cohen's Kappa | Weighted F1 |
| Electrode Position-Aware | **0.8134 ± 0.0030** | **0.8960 ± 0.0054** | **0.8889 ± 0.0033** | **0.7312 ± 0.0021** | **0.6943 ± 0.0114** | **0.7525 ± 0.0052** |
| Index-based | 0.8053 ± 0.0010 | 0.8917 ± 0.0013 | 0.8837 ± 0.0020 | 0.6830 ± 0.0028 | 0.6437 ± 0.0274 | 0.7157 ± 0.0250 |
| No channel embeddings | 0.7966 ± 0.0026 | 0.8820 ± 0.0016 | 0.8756 ± 0.0027 | 0.7117 ± 0.0051 | 0.6449 ± 0.0126 | 0.7338 ± 0.0128 |

| Channel Embedding Type | CHB-MIT, 2-class | | | SHU-MI, 2-class | | |
|---|---|---|---|---|---|---|
| | Balanced Accuracy | AUPR | AUROC | Balanced Accuracy | AUC-PR | AUROC |
| Electrode Position-Aware | **0.7221 ± 0.0465** | **0.3064 ± 0.0161** | **0.8704 ± 0.0133** | **0.6342 ± 0.0135** | **0.6980 ± 0.0130** | **0.6936 ± 0.0137** |
| Index-based | 0.5707 ± 0.0106 | 0.2353 ± 0.0991 | 0.8008 ± 0.0273 | 0.5990 ± 0.0107 | 0.6285 ± 0.0075 | 0.6425 ± 0.0104 |
| No channel embeddings | 0.5262 ± 0.0045 | 0.2578 ± 0.0011 | 0.8087 ± 0.0048 | 0.5750 ± 0.0114 | 0.6180 ± 0.0105 | 0.6241 ± 0.0084 |

Table 4: Ablation on channel-embedding strategies in CEReBrO's downstream performance.

**The Role of Pre-training**   To assess the impact of pre-training on downstream performance, we compare CEReBrO's performance using pre-trained weights against its performance when initialized with random weights. Results are displayed in Table 5. We observe that generally, pre-training leads to higher performance across all metrics in all downstream tasks. We attribute this improvement to the fact that pre-training discovers stable spatio-temporal feature extractors, guiding subsequent fine-tuning into broad, flat regions of the loss landscape that confer strong generalization. In contrast, a model trained with random initialization must first allocate representational capacity to learn the basic spatio-temporal statistics of the data, which steers optimization toward narrower, sharper minima. We observed that when trained with random weight initialization, CEReBrO exhibits a greater propensity to overfit even with regularization methods applied.

| Training Style | TUAB, 2-class | | | ISRUC, 5-class | | |
|---|---|---|---|---|---|---|
| | Balanced Accuracy | AUPR | AUROC | Balanced Accuracy | Cohen's Kappa | Weighted F1 |
| With Pre-training | **0.8134 ± 0.0030** | **0.8960 ± 0.0054** | **0.8889 ± 0.0033** | **0.7312 ± 0.0021** | **0.6943 ± 0.0114** | 0.7525 ± 0.0052 |
| W/o Pre-training | 0.8038 ± 0.0023 | 0.8911 ± 0.008 | 0.8814 ± 0.0016 | 0.7272 ± 0.0051 | 0.6876 ± 0.0066 | 0.7556 ± 0.0043 |

| Training Style | CHB-MIT, 2-class | | | SHU-MI, 2-class | | |
|---|---|---|---|---|---|---|
| | Balanced Accuracy | AUPR | AUROC | Balanced Accuracy | AUC-PR | AUROC |
| With Pre-training | **0.7221 ± 0.0465** | **0.3064 ± 0.0161** | **0.8704 ± 0.0133** | **0.6342 ± 0.0135** | **0.6980 ± 0.0130** | **0.6936 ± 0.0137** |
| W/o Pre-training | 0.6733 ± 0.0110 | 0.2799 ± 0.0118 | 0.8156 ± 0.0030 | 0.5791 ± 0.0379 | 0.6500 ± 0.0104 | 0.6454 ± 0.0096 |

Table 5: Ablation on the role of pre-training in CEReBrO's downstream performance.

# 5   Conclusion

We present CEReBrO, a self-supervised EEG foundation model that leverages novel, learnable electrode position-aware channel embeddings and an alternating attention architecture. Our approach yields four main benefits. First, the combination of per-channel tokenization, learnable temporal embeddings and electrode position-aware channel embeddings enables robust handling of diverse montages, including both unipolar (including all reference schemes) and bipolar configurations. Second, our alternating spatial and temporal attention layers provide a $2\times$ speedup and $6\times$ reduction in memory footprint compared to standard self-attention, which is prevalent in existing EEG foundation models, while achieving comparable or superior performance to the current state of the art. Third, by pre-training on over 25,000 hours of data drawn from several datasets (including clinical and consumer data), our model flexibly generalizes from low channel count consumer-grade devices to high-density clinical setups.

# 6   Limitations

Although EEG foundation models generally yield higher performance, they are often larger and incur greater computational overhead than conventional EEG architectures. This scale-up raises the barrier to entry, making real-time inference on embedded or low-power hardware particularly

challenging. Consequently, deploying these models in resource-constrained environments often requires specialized accelerators or cloud-based inference. For future work, we aim to investigate on-device deployment of compact EEG foundation models.

# References

[1] Michal Teplan. Fundamental of eeg measurement. *MEASUREMENT SCIENCE REVIEW*, 2, 01 2002.

[2] Yannick Roy, Hubert Banville, Isabela Albuquerque, Alexandre Gramfort, Tiago H. Falk, and Jocelyn Faubert. Deep learning-based electroencephalography analysis: a systematic review. *Journal of Neural Engineering*, 16(5):051001, Aug 2019.

[3] Mamunur Rashid, Bifta Sama Bari, and Md. Golam Sadeque. Mathematical principals and modeling of EEG signal exploration. *IOSR Journal of VLSI and Signal Processing*, 6(6):05–09, Nov. - Dec. 2016.

[4] Peter J. Brockwell and Richard A. Davis. *Introduction to Time Series and Forecasting*. Springer Texts in Statistics. Springer, Cham, 3 edition, 2016.

[5] Megan A. Boudewyn, Molly A. Erickson, Kurt Winsler, John Daniel Ragland, Andrew Yonelinas, Michael Frank, Steven M. Silverstein, Jim Gold, Angus W. 3rd MacDonald, Cameron S. Carter, Deanna M. Barch, and Steven J. Luck. Managing EEG studies: How to prepare and what to do once data collection has begun. *Psychophysiology*, 60(11):e14365, Nov 2023.

[6] Marina Diachenko, Simon J. Houtman, Erika L. Juarez-Martinez, Jennifer R. Ramautar, Roos Weiler, Huibert D. Mansvelder, Hilgo Bruining, Peter Bloem, and Klaus Linkenkaer-Hansen. Improved manual annotation of EEG signals through convolutional neural network guidance. *eNeuro*, 9(5):ENEURO.0160–22.2022, Sep 2022.

[7] Bharat Gurnani and Kirandeep Kaur. Data annotators: The unacclaimed heroes of artificial intelligence revolution in ophthalmology. *Indian Journal of Ophthalmology*, 70(5):1847, May 2022.

[8] Damián Pascual, Amir Aminifar, and David Atienza. A self-learning methodology for epileptic seizure detection with minimally-supervised edge labeling. In *2019 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 764–769, 2019.

[9] William O. Tatum, Barbara A. Dworetzky, and Donald L. Schomer. Artifact and recording concepts in EEG. *Journal of Clinical Neurophysiology: Official Publication of the American Electroencephalographic Society*, 28(3):252–263, June 2011.

[10] Robin Tibor Schirrmeister, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin Glasstetter, Katharina Eggensperger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. Deep learning with convolutional neural networks for EEG decoding and visualization. *Human Brain Mapping*, 38(11):5391–5420, 2017.

[11] Sally V Mathias and Meriem Bensalem-Owen. Artifacts that can be misinterpreted as interictal discharges. *Journal of Clinical Neurophysiology*, 36(4):264–274, 2019.

[12] Thorir Mar Ingolfsson, Simone Benatti, Xiaying Wang, Adriano Bernini, Pauline Ducouret, Philippe Ryvlin, Sandor Beniczky, Luca Benini, and Andrea Cossettini. Minimizing artifact-induced false-alarms for seizure detection in wearable eeg devices with gradient-boosted tree classifiers. *Scientific Reports*, 14(1):2980, 2024.

[13] Christopher Wang, Vighnesh Subramaniam, Adam Uri Yaari, Gabriel Kreiman, Boris Katz, Ignacio Cases, and Andrei Barbu. BrainBERT: Self-supervised representation learning for intracranial recordings. In *The Eleventh International Conference on Learning Representations*, 2023.

[14] Wenhui Cui, Woojae Jeong, Philipp Thölke, Takfarinas Medani, Karim Jerbi, Anand A. Joshi, and Richard M. Leahy. Neuro-gpt: Towards a foundation model for eeg, 2024.

[15] Yuqi Chen, Kan Ren, Kaitao Song, Yansen Wang, Yifan Wang, Dongsheng Li, and Lili Qiu. Eegformer: Towards transferable and interpretable large-scale eeg foundation model, 2024.

[16] Chaoqi Yang, M. Brandon Westover, and Jimeng Sun. Biot: Cross-data biosignal learning in the wild, 2023.

[17] Weibang Jiang, Liming Zhao, and Bao liang Lu. Large brain model for learning generic representations with tremendous EEG data in BCI. In *The Twelfth International Conference on Learning Representations*, 2024.

[18] Daoze Zhang, Zhizhang Yuan, Yang Yang, Junru Chen, Jingjing Wang, and Yafeng Li. Brant: Foundation model for intracranial neural signal. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

[19] Jiquan Wang, Sha Zhao, Zhiling Luo, Yangxuan Zhou, Haiteng Jiang, Shijian Li, Tao Li, and Gang Pan. CBramod: A criss-cross brain foundation model for EEG decoding. In *The Thirteenth International Conference on Learning Representations*, 2025.

[20] Demetres Kostas, Stephane Aroca-Ouellette, and Frank Rudzicz. Bendr: using transformers and a contrastive self-supervised learning task to learn from massive amounts of eeg data, 2021.

[21] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. In *Thirty-fourth Conference on Neural Information Processing Systems*, 2020.

[22] Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *The Eleventh International Conference on Learning Representations*, 2023.

[23] BESA GmbH. Electrodes and surface locations. `https://wiki.besa.de/index.php?title=Electrodes_and_Surface_Locations`, 2025. Accessed: 2025-04-12.

[24] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[25] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.

[26] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling, 2022.

[27] Julius Pettersson and Petter Falkman. Comparison of lstm, transformers, and mlp-mixer neural networks for gaze based human intention prediction. *Frontiers in Neurorobotics*, 17:1157957, 2023.

[28] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners, 2021.

[29] Po-Yao Huang, Hu Xu, Juncheng B Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.

[30] Jen Hong Tan. Pre-training of lightweight vision transformers on small datasets with minimally scaled images, 2024.

[31] Iyad Obeid and Joseph Picone. The temple university hospital eeg data corpus. *Frontiers in Neuroscience*, 10:196, 2016.

[32] Wei-Long Zheng and Bao-Liang Lu. Investigating critical frequency bands and channels for eeg-based emotion recognition with deep neural networks. *IEEE Transactions on Autonomous Mental Development*, 7(3):162–175, 2015.

[33] Wei-Long Zheng, Wei Liu, Yifei Lu, Bao-Liang Lu, and Andrzej Cichocki. Emotionmeter: A multimodal framework for recognizing human emotions. *IEEE Transactions on Cybernetics*, pages 1–13, 2018.

[34] Wei Liu, Wei-Long Zheng, Ziyi Li, Si-Yuan Wu, Lu Gan, and Bao-Liang Lu. Identifying similarities and differences in emotion recognition with eeg and eye movements among chinese, german, and french people. *Journal of Neural Engineering*, 19(2):026012, 2022.

[35] Eduardo López-Larraz, María Sierra-Torralba, Sergio Clemente, Galit Fierro, David Oriol, Javier Minguez, Luis Montesano, and Jens G. Klinzing. "bitbrain open access sleep dataset". 2025.

[36] B Kemp, A H Zwinderman, B Tuk, H A C Kamphuisen, and J J L Oberyé. Analysis of a sleep-dependent neuronal feedback loop: the slow-wave microcontinuity of the eeg. *IEEE Transactions on Biomedical Engineering*, 47(9):1185–1194, 2000.

[37] M. Perrin, E. Maby, S. Daligault, O. Bertrand, and J. Mattout. Objective and subjective evaluation of online error correction during p300-based spelling. *Advances in Human-Computer Interaction*, 2012(4), 2012.

[38] Reinmar J. Kobler, Andreea I. Sburlea, and Gernot R. Müller-Putz. Tuning characteristics of low-frequency eeg to positions and velocities in visuomotor and oculomotor tracking tasks. *Scientific Reports*, 2018.

[39] Andreas Schwarz, Joana Pereira, Reinmar Kobler, and Gernot R. Müller-Putz. Unimanual and bimanual reach-and-grasp actions can be decoded from human eeg. *IEEE Transactions on Bio-Medical Engineering*, September 2019.

[40] William O. Tatum IV. Long-term eeg monitoring: a clinical approach to electrophysiology. *Journal of Clinical Neurophysiology*, 18(5):442–455, September 2001.

[41] William O. Tatum IV, Brin Freund, and W. David Freeman. Eeg monitoring in the intensive care unit. *MedLink Neurology*, 2023. UPDATED 07.24.2023, RELEASED 01.11.2018, EXPIRES FOR CME 07.24.2026, Editor: John M Stern MD.

[42] Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and Brent J Lance. Eegnet: a compact convolutional neural network for eeg-based brain–computer interfaces. *Journal of Neural Engineering*, 15(5):056013, July 2018.

[43] Yonghao Song, Qingqing Zheng, Bingchuan Liu, and Xiaorong Gao. Eeg conformer: Convolutional transformer for eeg decoding and visualization. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31:710–719, 2023.

[44] J. Jing, W. Ge, S. Hong, M. B. Fernandes, Z. Lin, C. Yang, S. An, A. F. Struck, A. Herlopian, I. Karakis, J. J. Halford, M. C. Ng, E. L. Johnson, B. L. Appavu, R. A. Sarkis, G. Osman, P. W. Kaplan, M. B. Dhakar, L. Arcot Jayagopal, Z. Sheikh, O. Taraschenko, S. Schmitt, H. A. Haider, J. A. Kim, C. B. Swisher, N. Gaspard, M. C. Cervenka, A. A. Rodriguez Ruiz, J. W. Lee, M. Tabaeizadeh, E. J. Gilmore, K. Nordstrom, J. Y. Yoo, M. G. Holmes, S. T. Herman, J. A. Williams, J. Pathmanathan, F. A. Nascimento, Z. Fan, S. Nasiri, M. M. Shafi, S. S. Cash, D. B. Hoch, A. J. Cole, E. S. Rosenthal, S. F. Zafar, J. Sun, and M. B. Westover. Development of expert-level classification of seizures and rhythmic and periodic patterns during eeg interpretation. *Neurology*, 100(17):e1750–e1762, April 2023. PMID: 36878708; PMCID: PMC10136013.

[45] Chaoqi Yang, Danica Xiao, M. Brandon Westover, and Jimeng Sun. Self-supervised eeg representation learning for automatic sleep staging, 2023.

[46] Wei Yan Peh, Yuanyuan Yao, and Justin Dauwels. Transformer convolutional neural networks for automated artifact detection in scalp eeg, 2022.

[47] Hongli Li, Man Ding, Ronghua Zhang, and Chunbo Xiu. Motor imagery eeg classification algorithm based on cnn-lstm feature fusion network. *Biomedical Signal Processing and Control*, 72:103342, February 2022.

[48] Yonghao Song, Xueyu Jia, Lie Yang, and Longhan Xie. Transformer-based spatial-temporal feature learning for eeg decoding, 2021.

[49] Sirvan Khalighi, Teresa Sousa, José Moutinho Santos, and Urbano Nunes. Isruc-sleep: A comprehensive public dataset for sleep researchers. *Computer Methods and Programs in Biomedicine*, 124:180–192, 2016.

[50] Jun Ma, Banghua Yang, Wenzheng Qiu, Yunzhe Li, Shouwei Gao, and Xinxing Xia. A large EEG dataset for studying cross-session variability in motor imagery brain-computer interface. *Scientific Data*, 9(1):531, 2022.

[51] Ervin Sejdić, Igor Djurović, and Jin Jiang. Time-frequency feature representation using energy concentration: An overview of recent advances. *Digital Signal Processing*, 19(1):153–183, 2009.

[52] Ali Hossam Shoeb. *Application of Machine Learning to Epileptic Seizure Onset Detection and Treatment*. PhD thesis, Massachusetts Institute of Technology, 2009.

[53] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.

# A Additional Experiment Details

This section describes the details of our pre-training.

## A.1 Data Pre-processing

In pre-training, all signals are downsampled to 200 Hz and segmented into 30-second signals, which is longer than most downstream tasks studied in literature. Pre-training on longer signals than seen downstream helps the model learn long-range time dependencies [20]. Table 6 details all the data used in our pre-training corpus. We use data recorded from clinical devices (e.g. TUEG) and from consumer-grade devices (e.g. GWD and BOAS).

| Dataset | # Samples | # Channels | Sample Duration (s) | Sampling Frequency (Hz) |
|---|---|---|---|---|
| TUEG (w/o TUAB subjects) | 2,619,175 | 17-27 | 30 | 256 |
| Sleep-EDFx | 464,925 | 2 | 30 | 100 |
| BOAS | 120,095 | 2 | 30 | 256 |
| SEED-GER | 6,237 | 62 | 30 | 1000 |
| SEED | 4,590 | 62 | 30 | 1000 |
| GWD | 4,417 | 6-58 | 30 | 256 |
| SEED-IV | 4,500 | 62 | 30 | 1000 |
| BCI-NER | 3,541 | 56 | 30 | 256 |
| SEED-FRA | 3,052 | 62 | 30 | 1000 |

Table 6: Pre-training corpus details.

## A.2 Pre-training Hyperparameters

Our pre-training hyperparameters are available in Table 7. We detail the pre-training recipes for CEReBrO-Tiny (2.4M parameters), CEReBrO-Small (5.8M parameters) and CEReBrO-Base (23M parameters).

| | CEReBrO Tiny | CEReBrO Small | CEReBrO Base |
|---|---|---|---|
| Batch Size | | 4096 | |
| Scheduler | | Cosine + Linear warmup | |
| Optimizer | | AdamW | |
| Betas | | [0.9, 0.98] | |
| Peak learning rate | | 1.25e-3 | |
| Minimum learning rate | | 2.5e-7 | |
| Maximum allowed epochs | | 50 | |
| Training stopped at epoch | | 30 | |
| Warmup epochs | | 3 | |
| Masking ratio | | 0.5 | |
| Encoder layers | 8 | 10 | 12 |
| Number of attention heads per layer | | 12 | |
| Embedding dimension | 180 | 200 | 400 |
| MLP Size | 720 | 800 | 1600 |
| Weight decay | | 0.05 | |
| Patch size | | 200 | |

Table 7: Pre-training hyperparameters. The same hyperparameters are used across all model sizes, with a few exceptions specified in this table.

## A.3 Fine-tuning Hyperparameters

Our fine-tuning hyperparameters are available in Table 8. We detail the fine-tuning recipes for CEReBrO-Tiny (2.4M parameters), CEReBrO-Small (5.8M parameters) and CEReBrO-Base (23M parameters).

| | CEReBrO Tiny | CEReBrO Small | CEReBrO Base |
|---|:---:|:---:|:---:|
| **Batch Size** | | 256 | |
| **Scheduler** | | Cosine + Linear warmup | |
| **Optimizer** | | AdamW | |
| **Betas** | | [0.9, 0.999] | |
| **Peak learning rate** | | 5e-4 | |
| **Minimum learning rate** | | 2.5e-6 | |
| **Total epochs** | | 50 | |
| **Warmup epochs** | | 5 | |
| **Layer-wise learning rate decay factor** | | 0.85 | |
| **Weight decay** | | 0.05 | |
| **Drop path** | | 0.2 | |
| **Label smoothing** | | 0.1 | |

Table 8: Fine-tuning hyperparameters for CEReBrO.

## A.4 Pre-training Reconstructions

In Figure 3 we depict signal reconstructions using data not seen in pre-training. The grey patches denote masked portions of the signal. The target signal appears in red and the predicted signal is in blue. As expected, CEReBrO produces semantically plausible reconstructions for a wide variety of signals. The high quality reconstructions on unseen data suggests that CEReBrO's learned embeddings are strong and generalizable.
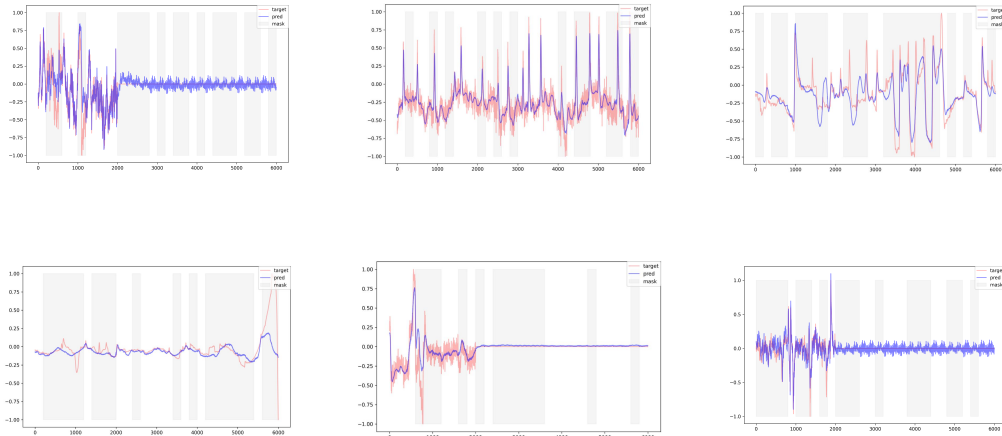


Figure 3: Example reconstructions from the pre-trained CEReBrO.

# B Efficiency Experiments

## B.1 Memory and Runtime Comparisons With Existing Foundation Models

In this experiment, we compare BIOT, LaBraM, CBraMod and CEReBrO's memory footprints and runtimes. We use the original implementations for all models.

In Table 9, CEReBrO leads in fine-tuning efficiency: at only 388 MFLOPs it requires approximately 55% the compute of LaBraM-Base, BIOT and CBraMod. As discussed in Section 3, although CBraMod uses a 4M parameter encoder, it requires a giant MLP in fine-tuning. In the case of CHB-MIT, this MLP contains 64M trainable parameters, which bring the total number of trainable parameters to 68M. In constrast, BIOT, LaBraM and CEReBrO only require lightweight classification/regression heads with a few thousand trainable parameters.

| Methods | Params | FLOPs |
|---|---|---|
| EEGNet | 0.003M | 8.9M |
| Conformer | 0.55M | 29.6M |
| SPaRCNet | 0.79M | 65.7M |
| ContraWR | 1.6M | 66.4M |
| CNN-Transformer | 3.2M | 79.1M |
| FFCL | 2.4M | 209.9M |
| ST-Transformer | 3.5M | 42.1M |
| BIOT | 3.2M | 703.92M |
| LaBraM-Base | 5.8M | 703.92M |
| CBraMod | 68.0M | 695.73M |
| CEReBrO | 2.4M | 388.02M |

Table 9: Model parameter counts and FLOPs during fine-tuning on CHB-MIT 16-channel, 10-second patches.



(a) Peak GPU Memory vs Input Sequence Length
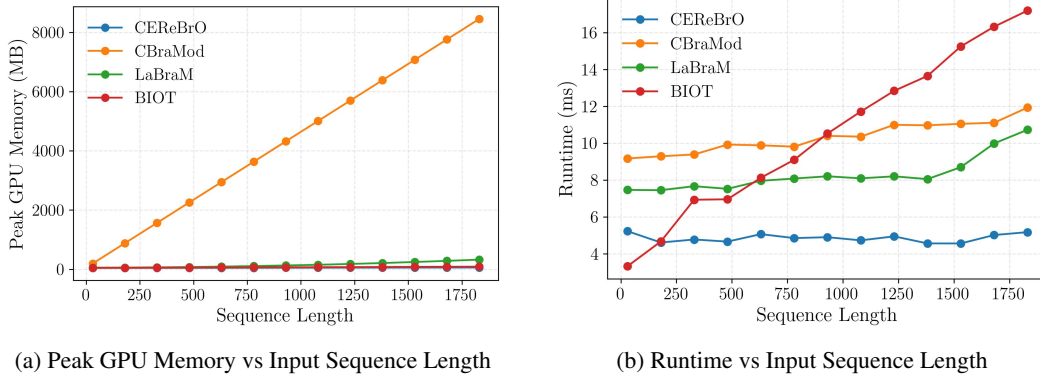
(b) Runtime vs Input Sequence Length

Figure 4: Performance comparison across EEG foundation models: (a) GPU memory usage and (b) runtime as a function of input sequence length. We average results over 10 forward passes for each model using a single NVIDIA GH200 GPU with 120GB RAM. Standard deviations in measurements are negligible at this scale.

In Figure 4a, we see that CEReBrO consumes less than 100MB even at 1920 tokens per forward pass (corresponding to 30 second, 64-channel patches). BIOT and LaBraM have a similar, but slightly larger memory footprint to CEReBrO, still on the order of approximately 100 MB. In contrast, CBraMod's memory utilization jumps to over 8 GB for the larger token sequences. This is driven by the giant MLP that CBraMod uses during fine-tuning, which can exceed 100M trainable parameters for certain downstream tasks.

In Figure 4b, we see that up to 1920 tokens per forward pass (corresponding to 30 second, 64-channel patches) , CEReBrO sustains a near-flat 4–5 ms runtime, outperforming all existing foundation models. CBraMod and LaBraM hover at 9–12 ms, while BIOT's latency climbs steeply (from 3.5 ms up to 17 ms). BIOT's slow runtime can be attributed to the computation of the Short Time Fourier Transform [51], which is not vectorized in the original implementation [16].

CBraMod's increased performance comes from its large MLP with 10M-100M parameters and multi-gigabyte memory footprint. CEReBrO remains lean at 2.4 million parameters, delivering competitive performance at a lower cost (2× faster runtime and up to 80× less memory) compared to CBraMod.

Compact EEG foundation models are critical for real-world impact. Small models fit on wearable and edge devices where memory is scarce and power is limited. Local inference preserves patient privacy and cuts cloud costs by avoiding continuous data uploads. Lower hardware requirements democratize research, enabling rapid iteration on moderate quality infrastructure. By balancing performance and efficiency, CEReBrO paves the way for deployable EEG foundation models.

## B.2 Attention Mechanism Comparisons

To validate the efficiency of the alternating attention mechanism, we compare it with self-attention[24] (used in LaBraM), linear attention (used in BIOT) and criss-cross attention (used in CBraMod). In this experiment, the patch embedding, temporal embedding, channel embedding and model size are the same across all configurations. Only the attention mechanism is modified. Results are shown in Figure 5.



(a) Peak GPU Memory vs Input Sequence Length          (b) Runtime vs Input Sequence Length
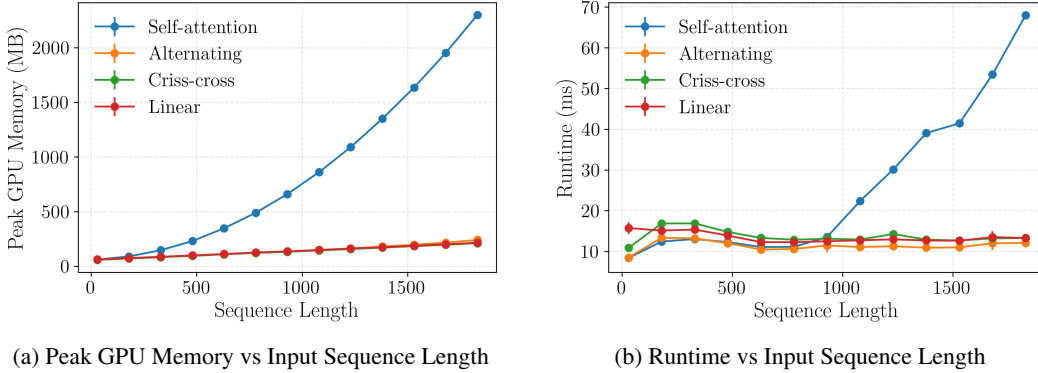
Figure 5: Efficiency comparison for different attention mechanisms: (a) GPU memory usage and (b) runtime as a function of input sequence length. We average results over 10 forward passes for each model using a single NVIDIA GH200 GPU with 120GB RAM. Standard deviations in measurements are negligible at this scale.

In Figure 5a and Figure 5b respectively, we see that for the same model size, alternating, criss-cross and linear attention yield similar runtime and peak memory usage for all sequence lengths. In contrast, we see that standard self-attention yields up to $10\times$ more memory and $6\times$ more latency than any other listed attention mechanism.

# C  Additional Details Regarding Downstream Tasks

**Downstream Tasks and Datasets**    To comprehensively evaluate the performance of our method, we select up to 10 downstream BCI tasks. All the downstream tasks with the corresponding datasets are presented in Table 10. For all the downstream datasets, we resampled the EEG signals to 200 Hz, consistent with the pre-training data. More details and experiments can be found in Section C.

Table 10: Overview of downstream tasks and datasets.

| Task | Dataset | Sampling Frequency (Hz) | # Channels | Duration (s) | # Samples | Label |
|---|---|---|---|---|---|---|
| Motor Imagery Classification | SHU-MI | 250 | 32 | 4 | 11,988 | 2-class |
| Sleep Staging | ISRUC | 200 | 6 | 30 | 89,240 | 5-class |
| Seizure Detection | CHB-MIT | 256 | 16 | 10 | 326,993 | 2-class |
| Abnormal Detection | TUAB | 256 | 16 | 10 | 409,455 | 2-class |

We provide more details regarding downstream tasks, including experiment settings, task descriptions and baseline descriptions.

## C.1  Downstream Task Descriptions

**Motor Imagery Classification** involves decoding motor imagery from EEG signals. We use SHU-MI [50]

**Sleep Staging** involves classifying sleep stages using EEG signals. We use ISRUC [49].

**Seizure Detection** involves identifying epileptic seizures from EEG signals. We use CHB-MIT [52].

**EEG Anomaly Detection** involves categorizing abnormal patterns or events in EEG signals. We use the TUAB [31] dataset.

## C.2 Baselines for Comparison

**EEGNet** [42] is a compact convolutional neural network built with depthwise and separable convolutions.

**EEGConformer** [43] leverages CNNs to learn low-level local features and self-attention to capture global correlations within local temporal features.

**SPaRCNet** [44] employs a 1D CNN with dense residual connections for feature extraction.

**ContraWR** [45] first transforms biosignals into multi-channel spectrograms and then utilizes a 2D-CNN based on ResNet [53] to extract features.

**CNN-Transformer** [46] uses CNNs for local feature extraction and transformers to model global dependencies.

**FFCL** [47] combines CNN and LSTM in parallel, with CNN extracting spatial features and LSTM capturing temporal dynamics.

**ST-Transformer** [48] relies on attention mechanisms within a transformer framework to learn both spatial and temporal characteristics of EEG signals.

**BIOT** [16] is an EEG foundation model that learns generic representations via a linear transformer and combined supervised-unsupervised pre-training; note that the pre-trained BIOT accepts only up to 18 channels, so multiple BIOT models are applied for EEG signals with more than 18 channels.

**LaBraM** [17] is an EEG foundation model that learns generic EEG representations by predicting neural quantized tokens for masked EEG patches using a full-attention transformer.

**CBraMod** [19] is an EEG foundation model that learns generic EEG representations by predicting neural tokens for masked EEG patches using a criss-cross attention.

## C.3 The Impact of CBraMod's MLP Head

CBraMod uses an encoder with 4M trainable parameters, accompanied by a large MLP head, which increases the total number of parameters to 68M or 24M depending on the task. The total model size varies by downstream tasks because the large MLP flattens the Transformer sequence output which is variable in size. For a fair comparison between CBraMod and all existing EEG foundation models, we provide CBraMod's downstream performance with a standard classification head, identical to the one used in CEReBrO. In Table 11 and Table 12, it is clear that CBraMod's performance across all downstream tasks decreases when we use a standard MLP classification head instead of a giant MLP head. Unlike LaBraM, BIOT and CEReBrO which compress representations into a single representative token (common practice in pre-trained models), CBraMod flattens the token output yielding a much more expressive classification head. The downside is this makes fine-tuning much more memory and time-intensive (see Section B). In addition it is unclear whether this method validates the usefulness of the pre-trained encoder's embeddings, or if the boosted performance is a result of the large number of added parameters.

Table 11: Performance of different methods in anomaly classification (TUAB).

| Methods | Model Size | Balanced Accuracy | AUPR | AUROC |
|---|---|---|---|---|
| EEGNet | 0.003M | 0.7642 ± 0.0036 | 0.8299 ± 0.0043 | 0.8412 ± 0.0031 |
| EEGConformer | 0.55M | 0.7758 ± 0.0049 | 0.8427 ± 0.0054 | 0.8445 ± 0.0038 |
| SPaRCNet | 0.79M | 0.7896 ± 0.0018 | 0.8414 ± 0.0018 | 0.8676 ± 0.0012 |
| ContraWR | 1.6M | 0.7746 ± 0.0041 | 0.8421 ± 0.0104 | 0.8456 ± 0.0074 |
| CNN-Transformer | 3.2M | 0.7777 ± 0.0022 | 0.8433 ± 0.0039 | 0.8461 ± 0.0013 |
| FFCL | 2.4M | 0.7848 ± 0.0038 | 0.8448 ± 0.0065 | 0.8569 ± 0.0051 |
| ST-Transformer | 3.5M | 0.7966 ± 0.0023 | 0.8521 ± 0.0026 | 0.8707 ± 0.0019 |
| BIOT | 3.2M | 0.7959 ± 0.0057 | 0.8792 ± 0.0023 | 0.8815 ± 0.0043 |
| CBraMod (with standard MLP head) | 4M | 0.8064 ± 0.0054 | 0.8933 ± 0.0032 | 0.8892 ± 0.0047 |
| LaBraM-Base | 5.8M | 0.8140 ± 0.0019 | 0.8965 ± 0.0016 | 0.9022 ± 0.0009 |
| CBraMod (with giant MLP head) | 68M | **0.8249 ± 0.0025** | **0.9221 ± 0.0015** | **0.9156 ± 0.0017** |
| CEReBrO (with reduced pre-training corpus) | 2.4M | 0.8146 ± 0.0060 | 0.9037 ± 0.0049 | 0.8946 ± 0.0050 |
| CEReBrO | 2.4M | 0.8215 ± 0.0027 | 0.9096 ± 0.0028 | 0.9023 ± 0.0020 |

Table 12: Performance of different methods in seizure detection (CHB-MIT) and motor imagery classification (SHU-MI).

| Methods | Model Size | CHB-MIT, 2-class | | | SHU-MI, 2-class | | |
|---|---|---|---|---|---|---|---|
| | | Balanced Accuracy | AUPR | AUROC | Balanced Accuracy | AUPR | AUROC |
| EEGNet | 0.003M | 0.5658 ± 0.0106 | 0.1914 ± 0.0182 | 0.8048 ± 0.0136 | 0.5889 ± 0.0177 | 0.6311 ± 0.0142 | 0.6283 ± 0.0152 |
| EEGConformer | 0.55M | 0.5976 ± 0.0141 | 0.2209 ± 0.0215 | 0.8226 ± 0.0170 | 0.5900 ± 0.0107 | 0.6370 ± 0.0093 | 0.6351 ± 0.0101 |
| SPaRCNet | 0.79M | 0.5876 ± 0.0191 | 0.1247 ± 0.0119 | 0.8143 ± 0.0148 | 0.5978 ± 0.0097 | 0.6510 ± 0.0062 | 0.6431 ± 0.0082 |
| ContraWR | 1.6M | 0.6344 ± 0.0002 | 0.2264 ± 0.0174 | 0.8097 ± 0.0114 | 0.5873 ± 0.0128 | 0.6315 ± 0.0105 | 0.6273 ± 0.0113 |
| CNN-Transformer | 3.2M | 0.6389 ± 0.0067 | 0.2479 ± 0.0227 | 0.8662 ± 0.0082 | 0.5975 ± 0.0169 | 0.6412 ± 0.0076 | 0.6323 ± 0.0082 |
| FFCL | 2.4M | 0.6262 ± 0.0104 | 0.2049 ± 0.0346 | 0.8271 ± 0.0051 | 0.5692 ± 0.0252 | 0.5943 ± 0.0172 | 0.6014 ± 0.0168 |
| ST-Transformer | 3.5M | 0.5915 ± 0.0195 | 0.1422 ± 0.0094 | 0.8237 ± 0.0491 | 0.5992 ± 0.0206 | 0.6394 ± 0.0122 | 0.6431 ± 0.0111 |
| BIOT | 3.2M | 0.7068 ± 0.0457 | 0.3277 ± 0.0460 | 0.8761 ± 0.0284 | 0.6179 ± 0.0183 | 0.6770 ± 0.0119 | 0.6609 ± 0.0127 |
| LaBraM-Base | 5.8M | 0.7075 ± 0.0358 | 0.3287 ± 0.0402 | 0.8679 ± 0.0199 | 0.6166 ± 0.0192 | 0.6761 ± 0.0083 | 0.6604 ± 0.0091 |
| CBraMod (with standard MLP head) | 4M | 0.6083 ± 0.0301 | 0.3764 ± 0.0395 | 0.8800 ± 0.0137 | 0.5012 ± 0.0028 | 0.4937 ± 0.0069 | 0.4950 ± 0.0003 |
| CBraMod (with giant MLP head) | 68M | **0.7398 ± 0.0284** | **0.3689 ± 0.0382** | **0.8892 ± 0.0154** | **0.6370 ± 0.0151** | **0.7139 ± 0.0088** | **0.6988 ± 0.0068** |
| CEReBrO | 2.4M | 0.7404 ± 0.0251 | 0.3378 ± 0.0297 | 0.8804 ± 0.0185 | 0.6140 ± 0.0138 | 0.7043 ± 0.0085 | 0.6987 ± 0.0149 |

## C.4 Scaling CEReBrO To Larger Sizes

In this section, we evaluate 3 model sizes of CEReBrO (Tiny, Small and Base) with 2.4M, 5.8M and 23M parameters accordingly.

Table 13: Performance of different methods in anomaly classification on the TUAB dataset (2-class).

| Methods | Model Size | Balanced Accuracy | AUPR | AUROC |
|---|---|---|---|---|
| EEGNet | 0.003M | 0.7642 ± 0.0036 | 0.8299 ± 0.0043 | 0.8412 ± 0.0031 |
| EEGConformer | 0.55M | 0.7758 ± 0.0049 | 0.8427 ± 0.0054 | 0.8445 ± 0.0038 |
| SPaRCNet | 0.79M | 0.7896 ± 0.0018 | 0.8414 ± 0.0018 | 0.8676 ± 0.0012 |
| ContraWR | 1.6M | 0.7746 ± 0.0041 | 0.8421 ± 0.0104 | 0.8456 ± 0.0074 |
| CNN-Transformer | 3.2M | 0.7777 ± 0.0022 | 0.8433 ± 0.0039 | 0.8461 ± 0.0013 |
| FFCL | 2.4M | 0.7848 ± 0.0038 | 0.8448 ± 0.0065 | 0.8569 ± 0.0051 |
| ST-Transformer | 3.5M | 0.7966 ± 0.0023 | 0.8521 ± 0.0026 | 0.8707 ± 0.0019 |
| BIOT | 3.2M | 0.7959 ± 0.0057 | 0.8792 ± 0.0023 | 0.8815 ± 0.0043 |
| LaBraM-Base | 5.8M | 0.8140 ± 0.0019 | 0.8965 ± 0.0016 | 0.9022 ± 0.0009 |
| CBraMod | 68M | **0.8249 ± 0.0025** | **0.9221 ± 0.0015** | **0.9156 ± 0.0017** |
| CEReBrO-Tiny | 2.4M | 0.8134 ± 0.0030 | 0.8960 ± 0.0054 | 0.8889 ± 0.0033 |
| CEReBrO-Small | 5.8M | 0.8144 ± 0.0070 | 0.8938 ± 0.0057 | 0.8827 ± 0.0024 |
| CEReBrO-Base | 23M | 0.8151 ± 0.0096 | 0.8966 ± 0.0073 | 0.8907 ± 0.0041 |

## C.5 The Role of Different Attention Mechanisms

| | TUAB, 2-class | | | SHU-MI, 2-class | | |
| --- | --- | --- | --- | --- | --- | --- |
| **Attention Mechanism** | **Balanced Accuracy** | **AUPR** | **AUROC** | **Balanced Accuracy** | **AUPR** | **AUROC** |
| Alternating | **0.8215 ± 0.0027** | **0.9096 ± 0.0028** | **0.9023 ± 0.0020** | 0.6140 ± 0.0138 | **0.7043 ± 0.0115** | **0.6987 ± 0.0149** |
| Axial | 0.8187 ± 0.0038 | 0.9057 ± 0.0013 | 0.8995 ± 0.0023 | 0.6116 ± 0.00122 | 0.6916 ± 0.0128 | 0.6848 ± 0.0174 |
| Linear | 0.8134 ± 0.0023 | 0.9004 ± 0.0025 | 0.8926 ± 0.0015 | 0.6146 ± 0.0131 | 0.6968 ± 0.0121 | 0.6795 ± 0.0174 |
| Full | 0.8193 ± 0.0021 | 0.9031 ± 0.0019 | 0.8947 ± 0.0018 | **0.6175 ± 0.00124** | 0.6806 ± 0.0102 | 0.6741 ± 0.00145 |

Table 14: Ablation on attention mechanisms in CEReBrO's downstream performance.