

# CodeBrain: Towards Decoupled Interpretability and Multi-Scale Architecture for EEG Foundation Model

Jingying Ma<sup>1\*</sup>, Feng Wu<sup>1\*</sup>, Qika Lin<sup>1†</sup>, Yucheng Xing<sup>1</sup>, Chenyu Liu<sup>3</sup>,  
 Ziyu Jia<sup>4,5†</sup>, Mengling Feng<sup>1,2</sup>

<sup>1</sup>Saw Swee Hock School of Public Health, National University of Singapore, Singapore

<sup>2</sup>Institute of Data Science, National University of Singapore, Singapore

<sup>3</sup>College of Computing and Data Science, Nanyang Technological University, Singapore

<sup>4</sup>Beijing Key Laboratory of Brainnetome and Brain-Computer Interface, Institute of Automation, Chinese Academy of Sciences, Beijing, China

<sup>5</sup>Brainnetome Center, Institute of Automation, Chinese Academy of Sciences, Beijing, China

## Abstract

Electroencephalography (EEG) provides real-time insights into brain activity and supports diverse applications in neuroscience. While EEG foundation models (EFMs) have emerged to address the scalability issues of task-specific models, current approaches still yield clinically uninterpretable and weakly discriminative representations, inefficiently capture global dependencies, and neglect important local neural events. We present *CodeBrain*, a two-stage EFM designed to fill this gap. In the first stage, we introduce the **TFDual-Tokenizer**, which decouples heterogeneous temporal and frequency EEG signals into discrete tokens, quadratically expanding the representation space to enhance discriminative power and offering domain-specific interpretability by suggesting potential links to neural events and spectral rhythms. In the second stage, we propose the multi-scale **EEGSSM** architecture, which combines structured global convolution with sliding window attention to efficiently capture both sparse long-range and local dependencies, reflecting the brain’s small-world topology. Pretrained on the largest public EEG corpus, CodeBrain achieves strong generalization across 8 downstream tasks and 10 datasets under distribution shifts, supported by comprehensive ablations, scaling-law analyses, and interpretability evaluations. Both code and pretraining weights will be released in the future version.

## 1 Introduction

Electroencephalography (EEG) captures brain activity via scalp electrodes [1] and provides high temporal-resolution signals for neuroscience and cognitive research [2]. To enable automated analysis, researchers have developed various task-specific models for applications such as sleep staging [3, 4], emotion recognition [5, 6], motor imagery [7, 8], and other applications [9, 10]. However, building separate models from scratch for each task is resource-intensive and limits scalability, as shared knowledge across tasks cannot be effectively leveraged. Moreover, variations in channel configurations and input lengths across EEG tasks further hinder knowledge transfer. To tackle these issues, EEG foundation models (EFMs) are developed to learn universal representations for diverse downstream tasks [11].

Inspired by masked self-supervised pretraining in natural language processing [12, 13], current EFMs commonly adopt patch-wise representation learning: EEG signals are divided into patches, encoded into latent representations, and trained to reconstruct the masked portions. While this approach offers flexibility across varying channel configurations and input lengths by adjusting the patch number and arrangement, direct raw-signal reconstruction [14, 15] remains challenging due to the inherent noise and variability of EEG. To mitigate this, recent studies have introduced codebook-based tokenization [16, 17], which abstracts away low-level fluctuations and provides a more robust latent space. Despite these advances, existing EFMs still face critical limitations, calling for new architectures.

---

\*Equal Contributions.

†Corresponding Authors: qikalin@foxmail.com, jia.ziyu@outlook.com.

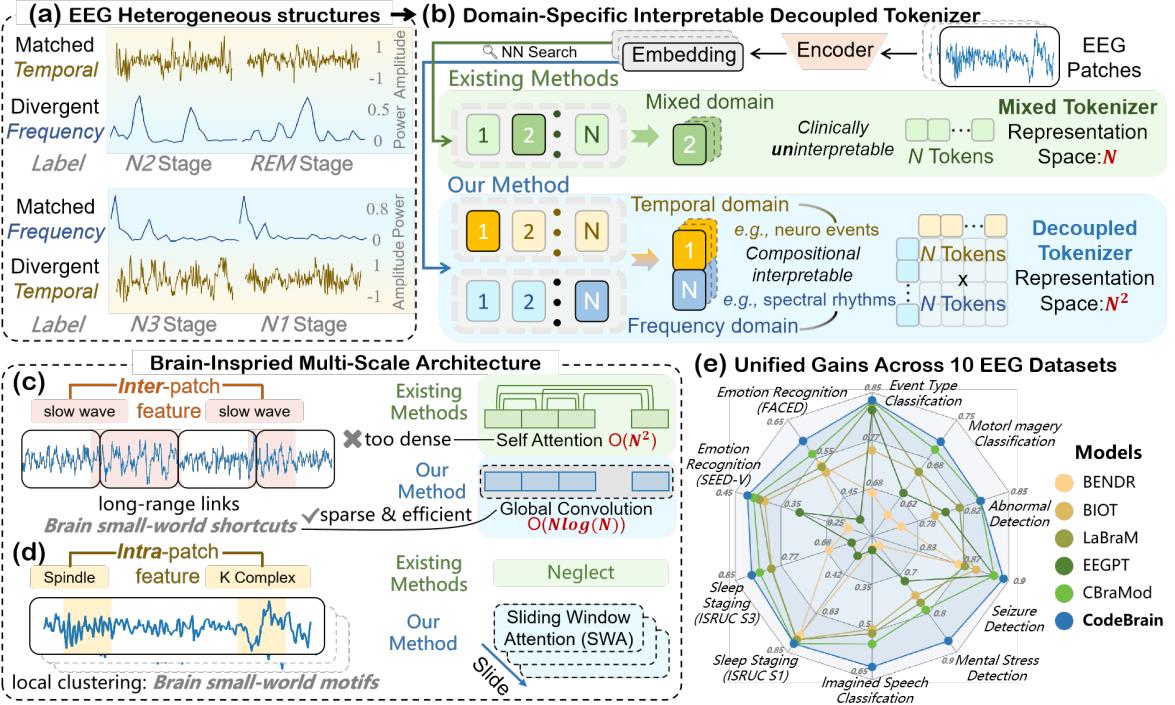


Figure 1: Rationale and overview of *CodeBrain* beyond existing EFM designs. (a) EEG signals are heterogeneous, as patches matched in one domain may diverge in the other. (b) We then propose a decoupled tokenizer for domain-specific interpretable representations while expanding the representation space. (c-d) Inspired by the brain’s small-world topology, a multi-scale architecture further captures inter-patch dependencies efficiently, while modeling overlooked intra-patch neural events. (e) These designs deliver performance gains across 10 EEG datasets.

**Failing to Decouple Heterogeneous EEG for Domain-Specific Interpretability.** Recent EFM designs adopt vector quantization for noise-robust representation [16, 17], following the VQ-VAE framework originally designed for images, where homogeneous visual features make a single tokenizer sufficient [12, 18]. In contrast, EEG exhibits heterogeneous structures: temporal and frequency components reflect distinct aspects of brain activity [19]. As illustrated in Fig. 1(a), signals matched in one domain may diverge in the other. Therefore, a mixed tokenizer can conflate domain-specific patterns [20], weakening representation capacity and producing tokens difficult to align with clinically interpretable neural events or spectral rhythms.

**Struggling with Efficiently Modeling Global Brain Dependencies.** EEG signals exhibit sparse global dependencies and strong local correlations, reflecting the brain’s small-world topology [21, 22, 23]. Efficient modeling of such structure requires capturing relationships in a scalable way. However, most EFM designs [24, 16, 14, 15] adopt Transformer architectures with fully connected self-attention [25]. This over-connected design is misaligned with the brain’s sparse structure and struggles to efficiently capture global dependencies due to its quadratic complexity with sequence length [26, 27, 28].

**Neglecting Local Dependencies within EEG Patches.** EEG signals exhibit rich local waveform structures over short temporal windows, reflecting crucial transient neural events (e.g., sleep waveforms in Fig. 1(d)) [29]. However, most existing EFM designs represent each EEG patch as a single token and apply attention mechanisms only at the patch level [15, 16], thereby ignoring important local dependencies within patches.

To address the above challenges, we propose *CodeBrain*, a novel EEG foundation model that integrates a decoupled tokenizer for domain-specific interpretability with a brain-inspired multi-scale architecture. *CodeBrain* is trained in two stages. In the first stage, we introduce the **TFDual-Tokenizer** (Fig. 1(b)), which decouples temporal and frequency EEG components into discrete tokens to support domain-specific interpretability. In the second stage, we develop **EEGSSM**, a masked self-supervised framework inspired by the brain’s small-world topology. **EEGSSM** adopts a structured

global convolution backbone, conceptually related to recent state-space sequence models [30, 31, 32], for sparse and efficient global modeling with a sliding window attention (SWA) mechanism for capturing local neural events overlooked by prior studies (Fig. 1(c–d)). Our detailed contributions are summarized as follows:

- **Decoupled Tokenizer for Domain-Specific Interpretability.** We propose the *TFDual-Tokenizer*, which decouples temporal and frequency EEG components into discrete tokens. This design quadratically expands the representation space, and qualitative analyses suggest that the resulting tokens align with neural events and spectral rhythms. A contrastive objective is further applied to the temporal branch to stabilize training. To the best of our knowledge, this is the first tokenizer in EFMs to provide domain-specific interpretability.
- **Brain Small-World Topology Inspired Multi-Scale Architecture.** We design *EEGSSM*, a patch-wise self-supervised framework for EEG. Guided by the brain’s small-world topology, it employs structured global convolution to capture sparse long-range temporal dependencies and sliding window attention to model local neural events. In addition, dynamic positional embeddings are used to flexibly learn spatial channel correlations.
- **Strong Generalization and Comprehensive Validation.** Pretrained on the largest publicly available EEG corpus, TUEG [33], CodeBrain achieves strong performance on 8 downstream tasks across 10 datasets (Fig. 1(e)) with distribution shifts in cohorts and channel configurations. This suggests the model design plays a central role in generalization. Comprehensive ablations, scaling-law analyses (3M–150M parameters), together with visualization and quantitative analyses, further confirm its robustness, scalability, and provide domain-specific interpretability.

## 2 Related Work

**EEG Foundation Models.** Inspired by foundation models in vision and language [34, 35], EEG research is shifting from task-specific models [36, 37, 38] to EFMs for learning expressive representations. Current EFMs can be divided into two categories. 1) Contrastive learning-based (CL) [39]: BENDER [40] first showed the ability of CL for EEG representations. Then, the Brant series [41, 42, 43] enables joint representation learning across physiological signals using CL. 2) Reconstruction-based: BIOT [24] pioneers cross-modal pretraining for biosignals, including EEG. Subsequent models focus specifically on EEG, learning representations by predicting masked discrete tokens [16, 44] or reconstructing raw signals [14, 45, 15]. Most existing EFMs adopt Transformer architecture, which is suboptimal for EEG due to poor handling of sparse dependencies, quadratic complexity, and ignoring local dependencies by treating each patch as a token.

**EEG Tokenization.** Tokenization has been key in NLP for generating generalizable and interpretable input representations [46, 47]. Inspired by this, early EFMs used patch-based continuous tokenization [24, 42] to handle EEG noise and variability, but without quantization, leading to unbounded and less interpretable representations. LaBraM [16] introduced vector quantization to learn discrete EEG tokens, following the VQ-VAE design from vision tasks [12]. However, this direct transfer overlooks EEG’s heterogeneous structure, limiting representation capacity. Moreover, the tokenizer of LaBraM is trained only on frequency-domain reconstruction due to convergence issues with raw signals reconstruction. Later efforts [17, 44] all relied on a single codebook to encode mixed domain features, limiting both the representation space and interpretability.

**State Space Model** Recent works have focused on enhancing classic State-Space Models (SSM) to more efficiently model sequential data using deep learning. For example, [48] used recurrent neural networks to learn parameters in SSM. However, the most significant progress came from [49] with the introduction of the structured state space model (S4), which reduces the computational complexity of SSM in modeling long sequences using a special state transition matrix [50]. These low-rank and normal matrices enable SSM to compute global convolution kernels efficiently through fast Fourier transform across the entire sequence. Subsequently, some works have further improved the shortcomings of S4 in areas such as model architecture [51] and convolution [52], and have started applying it to tasks such as natural language processing [53] and time series analysis [54]. Recently, some works have also

been applied to the EEG data, Tran et al. [55] leveraging SSMs to detect dementia. They extract temporal information from EEG signals through the Mamba architecture and combine it with frequency domain features to better manage the complexity of multivariate EEG. In the work of [56], SSM has also become the backbone network of the EEG foundation model. This further highlights the fast reasoning speed and efficient memory usage of the SSM model when processing EEG signals.

## 3 Methodology

### 3.1 Model Architecture

We introduce *CodeBrain*, a two-stage pretraining framework designed to learn interpretable and universal EEG representations (Fig. 2). The model is motivated by complementary goals: 1) **domain-specific interpretability** via decoupled tokenization of heterogeneous temporal and frequency information, achieved by the proposed **TFDual-Tokenizer**, and 2) **multi-scale modeling** of EEG sequences inspired by the brain’s small-world topology, addressed by the **EEGSSM** framework. This design lets Stage 1 learn a tokenizer of patch-level codes, while Stage 2 leverages it for EEG representations. We next provide a formal definition of the two stages to clarify their respective roles.

**Stage 1: Decoupled Tokenization.** Given a normalized EEG patch  $\mathbf{x} \in \mathbb{R}^L$ , where  $L$  is the patch length, our goal is to discretize  $\mathbf{x}$  into temporal and frequency tokens, enabling domain-specific representation learning. Specifically, let  $V_t \in \mathbb{R}^{K \times D}$  and  $V_f \in \mathbb{R}^{K \times D}$  denote the temporal and frequency codebooks, where  $K$  is the vocabulary size and  $D$  is the embedding dimension of each token. The tokenizer function is defined as:  $vt, vf = f_{\text{tokenizer}}(\mathbf{x}), vt, vf \in \mathbb{R}^D$ .

**Stage 2: EEG Representation Learning.** Given unlabeled EEG sequences  $\mathcal{X} = \{X_m\}_{m=1}^N$ , where each  $X_m \in \mathbb{R}^{C \times f \times T}$  consists of  $C$  channels, sampling rate  $f$ , and  $T$  seconds. We divide each sequence into  $n$  non-overlapping patches of  $t$  seconds, so each patch length is  $L = f \cdot t$ . The goal is to train an encoder  $f_{\text{enc}} : \mathbb{R}^{C \times n \times L} \rightarrow \mathbb{R}^{C \times n \times D}$  that produces latent representations  $Z_m = f_{\text{enc}}(X_m)$ .

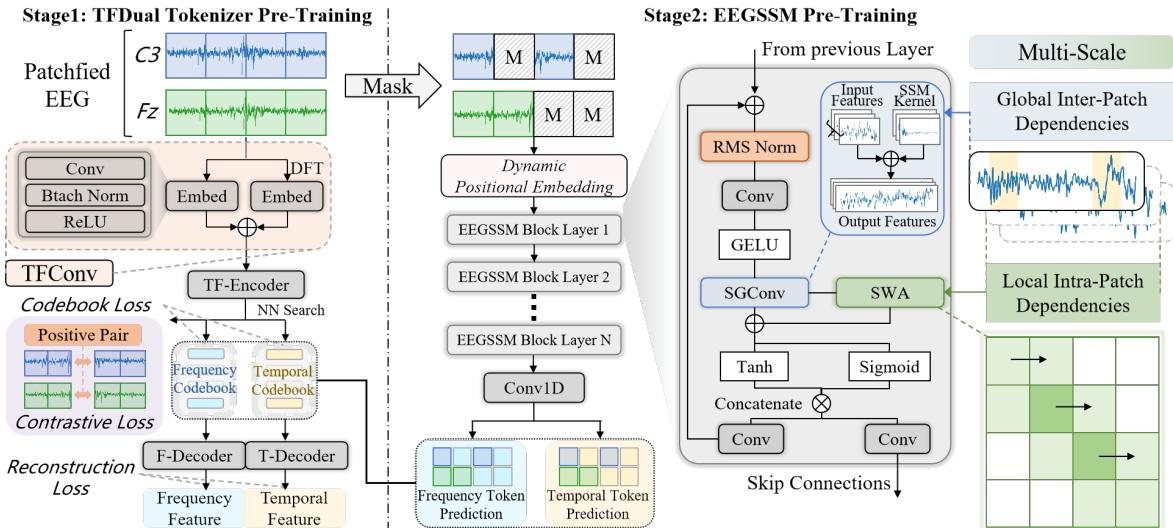


Figure 2: Overview of the *CodeBrain* framework. **Left:** *TFDual-Tokenizer* learns to discretize EEG signals into temporal and frequency tokens using two separate codebooks, by reconstructing both the temporal waveforms and the frequency-domain magnitude and phase. **Right:** *EEGSSM* learns representations by predicting the discrete tokens of masked patches generated by *TFDual-Tokenizer*.

### 3.2 TFDual-Tokenizer Pretraining

Our *TFDual-Tokenizer* includes a shared *neural encoder*, a *dual tokenizer* with separate codebooks, and two *decoders*. The *neural encoder* extracts joint time-frequency embeddings from EEG patches, which are then discretized into temporal and frequency tokens by the *dual tokenizer*. Each token stream is reconstructed by a decoder to supervise codebook learning in its respective domain.

**Neural Encoder** For each patch  $\mathbf{x}_i \in \mathbb{R}^L$ , we apply the Discrete Fourier Transform (DFT) [57] to obtain its frequency representation:

$$\mathbf{x}_i[k] = DFT(\mathbf{x}_i) \quad (1)$$

where  $\mathbf{x}_i[k]$  denotes the  $k$ -th frequency component.  $\mathbf{x}_i[k]$  and the  $\mathbf{x}_i$  are fed into the **TFCConv module**, where they are processed in parallel through stacks of convolutional, batch normalization, and ReLU layers. The temporal representation  $e_i^t = TFCConv(\mathbf{x}_i)$  and frequency representation  $e_i^f = TFCConv(\mathbf{x}_i[k])$  are concatenated to form a time-frequency embedding  $e_i^p = \text{Concat}\{e_i^t, e_i^f\}$ .

To get patch representation  $\tilde{e}_i$ , we then add a positional embedding  $e_{\text{pos}}$  and feed it into **TF-Encoder**:

$$\tilde{e}_i = \text{Encoder}(e_i^p + e_{\text{pos}}). \quad (2)$$

We choose a Transformer encoder here since this stage is *patch-to-token*, where its ability to model local contextual relations makes it well-suited for capturing patch-level patterns. To keep the tokenizer channel-agnostic,  $e_{\text{pos}}$  is shared temporal embeddings without channel-specific identities.

**Dual Tokenizer** We use two separate tokenizers with distinct codebooks for the temporal and frequency domains, denoted as  $vt_j, vf_j \in \mathbb{R}^D$ , where  $j = 1, \dots, K$ . Given the patch representation  $\tilde{e}_i$  from the neural encoder, each tokenizer independently selects the nearest code from its codebook:

$$pt_i = \arg \min_j \|\tilde{e}_i - vt_j\|^2, \quad pf_i = \arg \min_j \|\tilde{e}_i - vf_j\|^2, \quad (3)$$

where  $pt_i$  and  $pf$  denote the closest positions for the embeddings in the temporal and frequency domain codebook. The effectiveness of the Dual Tokenizer is based on the following proposition:

**Proposition 2.1** *Decoupling temporal and frequency codebooks yields representations that are no less effective than those from a joint codebook.*

*Proof.* See Appendix C.

For this Proposition, we provide empirical validation in Sections 4.4, 4.5 and analysis of Dual Tokenizer's interpretable structure in Appendix A.

**Frequency Codebook Training** To train the frequency codebook, we reconstruct amplitude and phase from the code embeddings. For each EEG patch, we apply the DFT to obtain the frequency representation:  $\mathbf{x}_i[k] = \text{Re}\{\mathbf{x}_i[k]\} + j \cdot \text{Im}\{\mathbf{x}_i[k]\}$  where  $\text{Re}\{\mathbf{x}_i[k]\}$  and  $\text{Im}\{\mathbf{x}_i[k]\}$  are the real part and imaginary part respectively, then the amplitude and phase can be calculated as:

$$A_i = \sqrt{\text{Re}(\mathbf{x}_i[k])^2 + \text{Im}(\mathbf{x}_i[k])^2}, \quad \phi_i = \arctan 2(\text{Im}(\mathbf{x}_i[k]), \text{Re}(\mathbf{x}_i[k])). \quad (4)$$

We use z-score normalization to ensure stable training. The code embedding  $vf_i$ , retrieved from the frequency codebook, is passed through the **F-Decoder**, which consists of a Transformer encoder followed by two linear layers:

$$y_i^A = \text{Encoder}(MLP(\tilde{e}_i)), \quad y_i^P = \text{Encoder}(MLP(\tilde{e}_i)). \quad (5)$$

where  $y_i^A$  and  $y_i^P$  are the predicted amplitude and phase, respectively. The frequency codebook's training objective is the mean squared error (MSE) loss:

$$\mathcal{L}_i^f = \|y_i^A - A_i\|_2^2 + \|y_i^P - \phi_i\|_2^2. \quad (6)$$

**Temporal Codebook Training** Direct reconstruction of temporal features might lead to non-convergence [16]. To address this, we combine contrastive loss with reconstruction loss to train a temporal codebook. Inspired by studies on physiological signals [58], we assume temporal dependencies exist between EEG segments, especially within the same channel. For an EEG segment  $X_m \in \mathbb{R}^{C \times n \times L}$ , we split it into two halves of length  $n/2$ . We use **TF-Encoder** in Eq. (2) to obtain the representations of these two parts  $X_{m1}, X_{m2} \in \mathbb{R}^{C \times \frac{n}{2} \times L}$  separately:

$$e_{mi}^h = \text{Encoder}(X_{mi}), i \in \{1, 2\}. \quad (7)$$

We encourage the latent representations of different parts within a single segment  $e_{m1}^h$  and  $e_{m2}^h$  to be similar, while making those of the same part across different segments  $e_{mi}^h$  and  $e_{sk}^h$  as distinct as possible, where  $X_s \in \mathbb{R}^{C \times n \times L}, s \neq m, k \neq i$ . SimCLR loss [59] is used for training:

$$\mathcal{L}_m^{CL} = -\log \frac{\exp(\text{sim}(e_{m1}^h, e_{m2}^h)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(e_{mi}^h, e_{sk}^h)/\tau)}, \quad (8)$$

where  $\text{sim}(e_{m1}^h, e_{m2}^h)$  is the cosine similarity,  $\tau > 0$  is the temperature parameter, and  $\mathbb{1}_{[k \neq i]}$  is an indicator function used to exclude itself. We introduce a **T-Decoder** to reconstruct raw signals from the temporal code embedding  $\tilde{e}_i$ , which consists of a Transformer encoder followed by a linear projection. Let  $y_t$  be the Transformer output from Eq. (2); the overall training objective is:

$$\mathcal{L}_i^t = \mathcal{L}_{CL} + \|y_i^t - \mathbf{x}_i\|_2^2, \quad y_i^t = \text{Encoder}(\text{MLP}(\tilde{e}_i)). \quad (9)$$

Finally, the training objective for the TFDual-Tokenizer can be defined as:

$$\begin{aligned} \mathcal{L}_{\text{tokenizer}} = \sum_{X_m \in \mathcal{X}} \mathcal{L}_m^{CL} + \sum_{X_m \in \mathcal{X}} \sum_{i=0}^n \mathcal{L}_i^f + \mathcal{L}_i^t + & \underbrace{\|\text{sg}(\tilde{e}_i) - vt_{pti}\|^2 + \|\text{sg}(\tilde{e}_i) - vf_{pf_i}\|^2}_{\text{codebook loss}} \\ & + \underbrace{\|\tilde{e}_i - \text{sg}(vt_{pti})\|^2 + \|\tilde{e}_i - \text{sg}(vf_{pf_i})\|^2}_{\text{commitment loss}}, \end{aligned} \quad (10)$$

where  $\text{sg}(\cdot)$  denotes stop-gradient to avoid updating encoder parameters.

### 3.3 EEGSSM Pretraining

In this stage, we introduce a novel convolutional structured state space model framework, called **EEGSSM**, composed of multiple *EEGSSM blocks*. To adapt to unseen channels, we first learn dynamic positional embeddings that capture inter-channel relationships through a lightweight convolutional module. The resulting features are processed by EEGSSM blocks. Another 1D convolutional layer then maps the output back to the token space for reconstructing the indices of masked tokens produced by the TFDual-Tokenizer.

**EEGSSM Block** Our EEGSSM block is composed of several blocks, which are integrated together through a residual connection mechanism. An EEGSSM block consists of a Layer Normalization, SGConv layer, SWA layer, and a gating component. Afterward, we feed the intermediate variables into the SGConv layer to obtain a global receptive field through convolution SSM.

*SGConv Layer.* SGConv is a structured SSM model (see Appendix B) using convolution architecture, and its convolution structure can be represented as a DFT formula:

$$y = F_N^{-1} D_k F_N u, D_k = \text{diag}(\bar{K} F_N), \quad (11)$$

where  $F_N$  denotes the DFT matrix of size  $N$ , and the convolution can be computed in  $O(N \log N)$  via FFT. As a type of convolutional SSM, SGConv improves the convolution kernel  $\bar{K}$  in Eq. (22) by introducing two features: sparse parameterization and kernel decay, making SGConv easier and more efficient to compute compared to the traditional S4 kernel. Let  $L$  be the length of the input sequence. The convolution kernel  $\bar{K}$  of SGConv is composed of several sub-kernels. Assuming the size of the first sub-kernel is  $d$ , with parameters  $w_i \in \mathbb{R}^{d \times d}$ , then the number of sub-kernels can be expressed as  $N = \log_2(\frac{L}{d}) + 1$ . The convolution kernel  $\bar{K}$  in Eq. (22) can thus be initialized as:

$$\bar{K} = \frac{1}{Z} [k_0, k_1, \dots, k_{N-1}], k_i = \alpha^i \text{Upsample}_{2^{\max[i-1, 0]} d}(w_i), \quad (12)$$

where  $\alpha$  denotes the decay coefficient, usually chosen to be 0.5, inducing the decaying structure, and  $\text{Upsample}(x)$  denotes upsample  $x$  to length  $l$ . We also introduce the normalization constant  $Z$  to ensure that the convolution operation does not change the scale of the input.

*Sliding Window Attention Layer.* We included a sliding window attention (SWA) layer to capture fine-grained local temporal dependencies. We apply a small fixed-length window, allowing the model to directly access the content within the context of the window through an attention mechanism by sliding it across the entire sequence, thereby addressing previous models' neglect of intra-patch temporal information. Furthermore, SWA maintains linear computational complexity to ensure training speed remains largely unaffected.

*Gate Mechanism.* we use a gating mechanism to control the output of the block. We employ a gated unit similar to Wavenet [33], which can suppress useless or irrelevant features and help stabilize

training in deep networks. We concatenate the global features output  $y_{sg}$  by the SGConv layer with the local features output  $y_{swa}$  by SWA and feed them into a gated unit:

$$z = \tanh(W_f \times \text{Concat}(y_{sg}, y_{swa})) \odot \sigma(W_g \times \text{Concat}(y_{sg}, y_{swa})), \quad (13)$$

$$y_1 = \text{Conv}(z), y_2 = \text{Conv}(z) \quad (14)$$

where  $\tanh(\cdot)$  and  $\sigma(\cdot)$  are tanh function and sigmoid function,  $\odot$  denotes an element-wise multiplication operator,  $W_f$  and  $W_g$  are learnable convolution filters.  $y_1$  becomes the input to the next block, while  $y_2$  will be aggregated to the output of SSM blocks through a skip connection.

**Pre-Training Objective** To help the EEGSSM model learn general EEG representations, we use a Masked Autoencoder (MAE) for self-supervised pre-training. For a patched sample  $X = \{x_i \mid i \in [1, 2, \dots, C]\}$ , we randomly generate a mask  $\mathcal{M} = \{m_i \mid i \in [1, 2, \dots, C]\}$  from a Bernoulli distribution of  $r$  proportion, where  $m_i \in \{0, 1\}$ . We reconstruct the token indices of masked EEG patches from the TFDual-Tokenizer by cross-entropy loss. Let  $y_i$  denotes the output of the EEGSSM block, the probability that the EEG signal matches the corresponding token  $v_i$  in the codebooks:

$$p(v_i|x_i) = \text{softmax}(\text{Conv1D}(y_i)). \quad (15)$$

Suppose the size of the pre-training set is  $N$ , the final cross-entropy loss is:

$$\mathcal{L}_p = - \sum_{j=0}^N \sum_{n \in \{m_i=1\}}^C p(v_{nj}|x_{nj}). \quad (16)$$

## 4 Experiments

### 4.1 Datasets

**Pre-Training.** We pretrain CodeBrain on the TUH EEG Corpus [33], the largest publicly available EEG dataset to date. Data processing follows a standardized pipeline: recordings shorter than 5 minutes are excluded, and the first and last minute of each segment is removed. We retain 19 commonly used EEG channels (C3, C4, Cz, F3, F4, Fp1, Fp2, F7, F8, Fz, O1, O2, P3, P4, Pz, T3, T4, T5, T6), selected based on the international 10–20 system for electrode placement [60]. We apply band-pass filtering (0.3–75 Hz), and notch filtering at 60 Hz to remove noise. The data is resampled to 200 Hz and divided into 30-second non-overlapping segments. Segments with absolute amplitudes over 100  $\mu$ V are filtered out. To normalize the signals, each value is divided by 100. Each segment is split into 1-second windows, resulting in 570 EEG patches per sample. After preprocessing, 1,109,545 samples (about 9,246 hours) are retained for pretraining.

**Downstream Tasks.** We evaluate CodeBrain on 8 downstream tasks across 10 public EEG datasets, which span diverse applications and exhibit distribution shifts from the pretraining dataset, to assess generalizability. Detailed dataset configurations are in Table 4.1. We perform cross-subject or cross-session splits with strict separation between training, validation, and test sets. For **FACED**, we use 80 subjects for training, 20 for validation, and the remaining 23 for testing. In **SEED-V**, each session consists of 15 trials, which are evenly divided into training, validation, and test sets. **ISRUUC\_S3** consists of 10 subjects, for which we apply an 8:1:1 cross-subject split. **MentalArithmetic** consists of 36 subjects, and we use a 7:1:1 cross-subject split. **BCIC2020-T3** follows the official competition protocol. In **CHB-MIT**, we use recordings from 19 subjects for training, and from 2 subjects each for validation and testing. Additional dataset details are provided in the Appendix F.

### 4.2 Experiment Settings

**Experiment Setup.** (1) Pretraining Setup. All experiments are conducted on NVIDIA 40GB A100 GPUs. The TFDual-Tokenizer is trained with temporal and frequency codebooks of 4096 codes (32 dimensions) for 20 epochs, and an 8-layer EEGSSM backbone with a masking ratio of 0.5 is trained for 10 epochs.

(2) Finetuning Strategy. We evaluate the quality of the pretrained representations under full finetuning. All downstream task datasets are resampled to 200 Hz to match the pretraining configuration.

Table 1: Summary of Downstream Tasks and Associated EEG Datasets.

Downstream Tasks	Datasets	#Channels	Length	#Samples	Class
Emotion Recognition	FACED [61]	32	10s	10,332	9-class
	SEED-V [62]	62	1s	117,744	5-class
Sleep Staging	ISRU C_S1 [63]	6	30s	86,320	5-class
	ISRU C_S3 [63]	6	30s	8,500	5-class
Imagined Speech Classification	BCIC 2020-T3 [64]	64	3s	6,000	5-class
	Mental Arithmetic [65]	20	5s	1,707	2-class
Seizure Detection	CHB-MIT [66]	16	10s	326,993	2-class
Motor Imagery Classification	SHU-MI [67]	32	4s	11,988	2-class
Event Type Classification	TUEV [33]	16	5s	112,491	6-class
Abnormal Detection	TUAB [33]	16	10s	409,455	2-class

A three-layer MLP is applied to aggregate channel information, compress the  $x$ -second sequence, and map the representation to the target class, with activation and dropout between layers.

**Baselines.** We compare our model with five publicly available EFM s that have released pretrained weights, representing different pretraining strategies. **BENDR** [40] adopts contrastive learning framework. **BIOT** [24] uses patch-based continuous tokenization. **LaBraM** introduces discrete neural tokens through vector quantization [16]. **EEGPT** [14] and **CBraMod** [15] rely on masked reconstruction of raw EEG signals.

**Evaluation Metric.** For multi-class classification, we report **Cohen’s Kappa**, **Weighted F1 score**, and **Balanced Accuracy**, with Kappa based on validation performance for testing. For binary classification, we use Area Under the ROC Curve (**AUROC**), Area Under the Precision-Recall Curve (**AUC-PR**), and **Balanced Accuracy**, with AUROC based on validation performance for testing. All experiments are repeated with five random seeds, and we report the mean and standard deviation.

More details about hyperparameter, baseline, and evaluation metric are provided in Appendix H, G.

Table 2: Comparison results of different methods on downstream tasks.

Methods	FACED (9-Class)			SEED-V (5-Class)		
	Cohen’s Kappa	Weighted F1	Balanced Acc	Cohen’s Kappa	Weighted F1	Balanced Acc
BENDR	$0.4716 \pm 0.0095$	$0.5340 \pm 0.0086$	$0.5320 \pm 0.0083$	$0.0335 \pm 0.0062$	$0.2026 \pm 0.0330$	$0.2231 \pm 0.0059$
BIOT	$0.4476 \pm 0.0254$	$0.5136 \pm 0.0112$	$0.5118 \pm 0.0118$	$0.2261 \pm 0.0262$	$0.3856 \pm 0.0203$	$0.3837 \pm 0.0187$
LaBraM	$0.4698 \pm 0.0102$	$0.5288 \pm 0.0188$	$0.5273 \pm 0.0107$	$0.2386 \pm 0.0209$	$0.3974 \pm 0.0111$	$0.3976 \pm 0.0138$
EEGPT	$0.4639 \pm 0.0023$	$0.3924 \pm 0.0017$	$0.4607 \pm 0.0014$	$0.1323 \pm 0.0062$	$0.3090 \pm 0.0052$	$0.3061 \pm 0.0044$
CBraMod	$0.5041 \pm 0.0122$	$0.5618 \pm 0.0093$	$0.5509 \pm 0.0089$	$0.2569 \pm 0.0143$	$0.4101 \pm 0.0108$	$0.4091 \pm 0.0097$
CodeBrain	<b><math>0.5406 \pm 0.0084</math></b>	<b><math>0.5953 \pm 0.0113</math></b>	<b><math>0.5941 \pm 0.0098</math></b>	<b><math>0.2735 \pm 0.0032</math></b>	<b><math>0.4235 \pm 0.0022</math></b>	<b><math>0.4137 \pm 0.0023</math></b>
Methods	ISRU C_S3 (5-Class)			BCIC 2020-T3 (5-Class)		
	Cohen’s Kappa	Weighted F1	Balanced Acc	Cohen’s Kappa	Weighted F1	Balanced Acc
BENDR	$0.5995 \pm 0.0151$	$0.6789 \pm 0.0142$	$0.6352 \pm 0.0095$	$0.0607 \pm 0.0093$	$0.2379 \pm 0.0165$	$0.2485 \pm 0.0075$
BIOT	$0.7168 \pm 0.0119$	$0.7834 \pm 0.0096$	$0.7598 \pm 0.0109$	$0.3650 \pm 0.0176$	$0.4917 \pm 0.0079$	$0.4920 \pm 0.0086$
LaBraM	$0.7194 \pm 0.0162$	$0.7843 \pm 0.0189$	$0.7617 \pm 0.0122$	$0.3800 \pm 0.0242$	$0.5054 \pm 0.0205$	$0.5060 \pm 0.0155$
EEGPT	$0.6160 \pm 0.0856$	$0.6375 \pm 0.0632$	$0.6650 \pm 0.0311$	$0.0567 \pm 0.0164$	$0.2441 \pm 0.0105$	$0.2453 \pm 0.0131$
CBraMod	$0.7407 \pm 0.0251$	$0.8056 \pm 0.0219$	$0.7844 \pm 0.0126$	$0.4216 \pm 0.0163$	$0.5383 \pm 0.0096$	$0.5373 \pm 0.0108$
CodeBrain	<b><math>0.7671 \pm 0.0091</math></b>	<b><math>0.8202 \pm 0.0071</math></b>	<b><math>0.7856 \pm 0.0031</math></b>	<b><math>0.5127 \pm 0.0065</math></b>	<b><math>0.6101 \pm 0.0053</math></b>	<b><math>0.6101 \pm 0.0052</math></b>
Methods	Mental Arithmetic (2-Class)			CHB MIT (2-Class)		
	AUROC	AUC-PR	Balanced Acc	AUROC	AUC-PR	Balanced Acc
BENDR	$0.6248 \pm 0.0765$	$0.3661 \pm 0.0672$	$0.5681 \pm 0.0448$	$0.8632 \pm 0.0526$	$0.3071 \pm 0.1240$	$0.5609 \pm 0.0432$
BIOT	$0.7536 \pm 0.0144$	$0.6004 \pm 0.0195$	$0.6875 \pm 0.0186$	$0.8761 \pm 0.0284$	$0.3277 \pm 0.0460$	$0.7068 \pm 0.0457$
LaBraM	$0.7721 \pm 0.0093$	$0.5999 \pm 0.0155$	$0.6909 \pm 0.0125$	$0.8679 \pm 0.0199$	$0.3287 \pm 0.0402$	$0.7075 \pm 0.0358$
EEGPT	$0.7162 \pm 0.0171$	$0.5081 \pm 0.0275$	$0.5597 \pm 0.0171$	$0.8892 \pm 0.0066$	$0.3073 \pm 0.0641$	$0.5481 \pm 0.0151$
CBraMod	$0.7905 \pm 0.0073$	$0.6267 \pm 0.0099$	$0.7256 \pm 0.0132$	$0.8892 \pm 0.0154$	$0.3689 \pm 0.0382$	<b><math>0.7398 \pm 0.0284</math></b>
CodeBrain	<b><math>0.8707 \pm 0.0209</math></b>	<b><math>0.7177 \pm 0.0421</math></b>	<b><math>0.7514 \pm 0.0203</math></b>	<b><math>0.8961 \pm 0.0174</math></b>	<b><math>0.4377 \pm 0.0288</math></b>	$0.7273 \pm 0.0240$

### 4.3 Comparison with Baselines

We ensure consistent data splits across all baselines. Results are reported on six representative downstream datasets, with additional results provided in Appendix I. As shown in Table 2, CodeBrain achieves consistent performance gains compared to baselines. For multi-class classification, it achieves the largest gain of +0.0911 in Cohen’s Kappa (21.6%), +0.0718 in Weighted F1 score (13.3%), +0.0728 in Balanced Acc (13.5%) on *BCIC 2020-T3* over the strongest baseline [15]. For binary classification, it achieves the largest gain of +0.0802 in AUROC (10.1%), +0.0910 in AUC-PR (14.5%) and +0.0258 in Balanced Acc (3.6%) on *Mental Arithmetic* over the strongest baseline. These results demonstrate the superior generalizability of CodeBrain.

### 4.4 Ablation Study

We conduct ablation studies on three datasets with the same five seeds as in the main experiments to evaluate the key components of CodeBrain (Table 4.3). Below is the detailed analysis:

(1) *Tokenizer configuration*: we compare the proposed *TFDual-Tokenizer* (**Dual**) with variants using a single domain codebook (**Temporal** or **Frequency**), or a shared codebook jointly reconstructing both domains (**Mixed**). Across all datasets, the Dual codebook consistently yields superior performance.

(2) *Contrastive learning* (CL): We evaluate the impact of contrastive learning in TFDual-Tokenizer pretraining. It leads to consistent gains, indicating a better capture of temporal patterns. Moreover, the CL facilitates convergence of the temporal codebook, with detailed analyses provided in Appendix E.

(3) *Components of EEGSSM*: Evaluating the **SWA**, **SGConv**, and **Gate** modules, which demonstrate improvements in EEG representation learning.

(4) *Scaling Laws*: Prior works [15, 16] explored scaling with 1-1000 hours of EEG data for pretraining. We extend to 1k-9k hours and 3M-150M models. As shown in Figure 3, Kappa consistently improves with more data and parameters. These results confirm that larger models yield consistent but diminishing returns. More detailed results are in Appendix M.

We further investigate several key design choices, including mask ratio, SWA window size, codebook size, and patch size in Appendix J. We report robustness experiments to demonstrate the adaptability of dynamic positional embeddings in Appendix L, and computational efficiency in Appendix K. We also provide low-resource comparisons (Appendix N) and pretraining curves (Appendix D), showing stable VQ convergence and efficient EEGSSM training.

Table 3: The results of ablation studies for tokenizer configurations and module components.

Dataset	Codebook	CL	SWA	SGConv	Gate	Cohen’s Kappa	Weighted F1	Balanced Accuracy
FACED 9-Class	Dual	✓	✓	✓	✓	<b>0.5406 ± 0.0084</b>	<b>0.5953 ± 0.0113</b>	<b>0.5941 ± 0.0098</b>
	Temporal	✓	✓	✓	✓	0.4618 ± 0.0072	0.5277 ± 0.0067	0.5217 ± 0.0056
	Frequency	✓	✓	✓	✓	0.5006 ± 0.0224	0.5607 ± 0.0201	0.5580 ± 0.0187
	Mixed	✓	✓	✓	✓	0.4676 ± 0.0061	0.5319 ± 0.0052	0.5281 ± 0.0049
	Dual	✗	✓	✓	✓	0.5222 ± 0.0082	0.5811 ± 0.0084	0.5765 ± 0.0074
	Dual	✓	✗	✓	✓	0.5192 ± 0.0092	0.5792 ± 0.0093	0.5736 ± 0.0075
	Dual	✓	✓	✗	✓	0.1936 ± 0.1637	0.2627 ± 0.1824	0.2858 ± 0.1467
	Dual	✓	✓	✓	✗	0.2578 ± 0.0340	0.3363 ± 0.0270	0.3431 ± 0.0316
SEED-V 5-Class	Dual	✓	✓	✓	✓	<b>0.2735 ± 0.0032</b>	<b>0.4235 ± 0.0022</b>	<b>0.4137 ± 0.0023</b>
	Temporal	✓	✓	✓	✓	0.2633 ± 0.0116	0.4152 ± 0.0092	0.4068 ± 0.0074
	Frequency	✓	✓	✓	✓	0.2665 ± 0.0208	0.4186 ± 0.0177	0.4098 ± 0.0147
	Mixed	✓	✓	✓	✓	0.2708 ± 0.0047	0.4214 ± 0.0044	0.4124 ± 0.0032
	Dual	✗	✓	✓	✓	0.2589 ± 0.0065	0.4129 ± 0.0056	0.4029 ± 0.0042
	Dual	✓	✗	✓	✓	0.2561 ± 0.0051	0.4106 ± 0.0042	0.4019 ± 0.0034
	Dual	✓	✓	✗	✓	0.2062 ± 0.0099	0.3707 ± 0.0075	0.3620 ± 0.0083
	Dual	✓	✓	✓	✗	0.2212 ± 0.0076	0.3826 ± 0.0057	0.3757 ± 0.0052
ISRUC_S3 5-Class	Dual	✓	✓	✓	✓	<b>0.7671 ± 0.0091</b>	<b>0.8202 ± 0.0071</b>	<b>0.7856 ± 0.0031</b>
	Temporal	✓	✓	✓	✓	0.7314 ± 0.0210	0.7916 ± 0.0181	0.7565 ± 0.0244
	Frequency	✓	✓	✓	✓	0.7390 ± 0.0601	0.7986 ± 0.0514	0.7728 ± 0.0361
	Mixed	✓	✓	✓	✓	0.7400 ± 0.0217	0.7999 ± 0.0171	0.7673 ± 0.0157
	Dual	✗	✓	✓	✓	0.7558 ± 0.0333	0.8130 ± 0.0264	0.7801 ± 0.0132
	Dual	✓	✗	✓	✓	0.7359 ± 0.0324	0.7950 ± 0.0259	0.7621 ± 0.0311
	Dual	✓	✓	✗	✓	0.6218 ± 0.0427	0.6956 ± 0.0279	0.6664 ± 0.0316
	Dual	✓	✓	✓	✗	0.5478 ± 0.0302	0.6429 ± 0.0307	0.6258 ± 0.0448

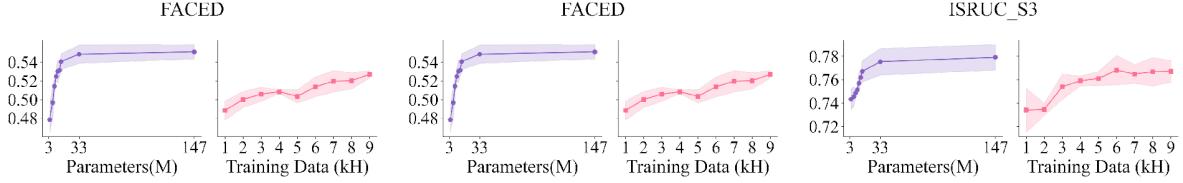


Figure 3: Model and training data scaling laws of CodeBrain across three datasets on Cohen’s Kappa.

#### 4.5 Vector Visualization

To demonstrate how the **TFDual-Tokenizer** models heterogeneous EEG, we visualize its learned temporal and frequency codes on the ISRUC\_S3 dataset. As shown in Figure 4(a)(b), each domain-specific codebook captures meaningful structures: temporal codes correspond to neural events (e.g., slow waves), while frequency codes highlight spectral rhythms such as dominant delta activity, either of which can support sleep stage discrimination. Yet in many cases a single domain is insufficient, and discrimination arises only from their composition, as in Figure 4(c)(d), where the same temporal code can pair with different frequency codes, and vice versa, to yield complementary representations. This decoupled design expands the representation space and enhances interpretability, with additional potential links to neurophysiological features and quantitative analyses provided in Appendix A.

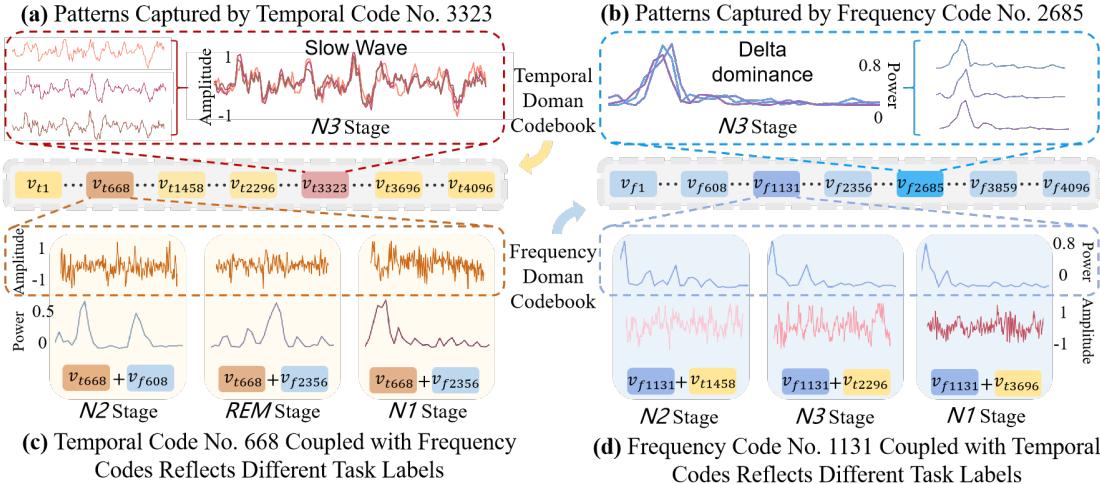


Figure 4: Decoupled time-frequency codes visualization on ISRUC\_S3 dataset.

## 5 Conclusion

In this paper, we present *CodeBrain*, a EEG foundation model that unifies interpretable tokenization with a brain-inspired multi-scale architecture. The TFDual-Tokenizer decouples heterogeneous EEG signals, expanding the representation space while suggesting domain-specific interpretability, and the EEGSSM architecture integrates structured global convolution with sliding-window attention to efficiently capture both long-range and local dependencies. Pretrained on the large-scale TUEG corpus, CodeBrain demonstrates strong generalization across 8 tasks and 10 datasets under distribution shifts, with comprehensive ablations, scaling-law analyses confirming robustness and scalability. These results establish CodeBrain as a strong foundation for neural time-series representation learning.

## References

- [1] Ernst Niedermeyer and FH Lopes da Silva. *Electroencephalography: basic principles, clinical applications, and related fields*. Lippincott Williams & Wilkins, 2005.
- [2] Fernando Lopes da Silva. Eeg and meg: relevance to neuroscience. *Neuron*, 80(5):1112–1128, 2013.
- [3] Hyojin Lee, You Rim Choi, Hyun Kyung Lee, Jaemin Jeong, Joopyo Hong, Hyun-Woo Shin, and Hyung-Sin Kim. Explainable vision transformer for automatic visual sleep staging on multimodal psg signals. *npj Digital Medicine*, 8(1):55, 2025.
- [4] Jingying Ma, Qika Lin, Ziyu Jia, and Mengling Feng. St-usleepnet: A spatial-temporal coupling prominence network for multi-channel sleep staging. In *Proceedings of the Thirty-Fourth International Joint Conference on Artificial Intelligence*, pages 4182–4190, 2025.
- [5] Ziyu Jia, Youfang Lin, Xiyang Cai, Haobin Chen, Haijun Gou, and Jing Wang. Sst-emotionnet: Spatial-spectral-temporal based attention 3d dense network for eeg emotion recognition. In *Proceedings of the 28th ACM international conference on multimedia*, pages 2909–2917, 2020.
- [6] Chenyu Liu, Xinliang Zhou, Zhengri Zhu, Liming Zhai, Ziyu Jia, and Yang Liu. Vbh-gnn: Variational bayesian heterogeneous graph neural networks for cross-subject emotion recognition. In *The Twelfth International Conference on Learning Representations*, 2024.
- [7] Zhenqi Li, Jing Wang, Ziyu Jia, and Youfang Lin. Learning space-time-frequency representation with two-stream attention based 3d network for motor imagery classification. In *2020 IEEE International Conference on Data Mining (ICDM)*, pages 1124–1129. IEEE, 2020.
- [8] Ziyu Jia, Youfang Lin, Jing Wang, Kaixin Yang, Tianhang Liu, and Xinwang Zhang. Mmcnn: A multi-branch multi-scale convolutional neural network for motor imagery classification. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 736–751. Springer, 2020.
- [9] Michele Guerra, Roberto Milanese, Michele Deodato, Madalina G Ciobanu, and Fausto Fasano. Exploring the diagnostic potential of llms in schizophrenia detection through eeg analysis. In *2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 6812–6819. IEEE, 2024.
- [10] Yongquan Hu, Shuning Zhang, Ting Dang, Hong Jia, Flora D Salim, Wen Hu, and Aaron J Quigley. Exploring large-scale language models to evaluate eeg-based multimodal data for mental health. In *Companion of the 2024 on ACM International Joint Conference on Pervasive and Ubiquitous Computing*, pages 412–417, 2024.
- [11] Xinliang Zhou, Chenyu Liu, Zhisheng Chen, Kun Wang, Yi Ding, Ziyu Jia, and Qingsong Wen. Brain foundation models: A survey on advancements in neural signal processing and brain discovery. *arXiv preprint arXiv:2503.00580*, 2025.
- [12] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *Advances in neural information processing systems*, volume 30, 2017.
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.
- [14] Guangyu Wang, Wenchao Liu, Yuhong He, Cong Xu, Lin Ma, and Haifeng Li. Eegpt: Pretrained transformer for universal and reliable representation of eeg signals. In *Advances in Neural Information Processing Systems*, volume 37, pages 39249–39280, 2024.
- [15] 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 Third International Conference on Learning Representations*, 2025.

- [16] 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.
- [17] Jathurshan Pradeepkumar, Xihao Piao, Zheng Chen, and Jimeng Sun. Single-channel eeg tokenization through time-frequency modeling. *arXiv preprint arXiv:2502.16060*, 2025.
- [18] Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. Finite scalar quantization: Vq-vae made simple. In *The Twelfth International Conference on Learning Representations*, 2024.
- [19] Fumikazu Miwakeichi, Eduardo Martínez-Montes, Pedro A Valdés-Sosa, Nobuaki Nishiyama, Hiroaki Mizuhara, and Yoko Yamaguchi. Decomposing eeg data into space–time–frequency components using parallel factor analysis. *NeuroImage*, 22(3):1035–1045, 2004.
- [20] Qijiong Liu, Xiaoyu Dong, Jiaren Xiao, Nuo Chen, Hengchang Hu, Jieming Zhu, Chenxu Zhu, Tetsuya Sakai, and Xiao-Ming Wu. Vector quantization for recommender systems: a review and outlook. *arXiv preprint arXiv:2405.03110*, 2024.
- [21] Ed Bullmore and Olaf Sporns. Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature reviews neuroscience*, 10(3):186–198, 2009.
- [22] Danielle Smith Bassett and ED Bullmore. Small-world brain networks. *The neuroscientist*, 12(6):512–523, 2006.
- [23] Yong He, Jinhui Wang, Liang Wang, Zhang J Chen, Chaogan Yan, Hong Yang, Hehan Tang, Chaozhe Zhu, Qiyong Gong, Yufeng Zang, et al. Uncovering intrinsic modular organization of spontaneous brain activity in humans. *PloS one*, 4(4):e5226, 2009.
- [24] Chaoqi Yang, M Westover, and Jimeng Sun. Biot: Biosignal transformer for cross-data learning in the wild. In *Advances in Neural Information Processing Systems*, volume 36, pages 78240–78260, 2023.
- [25] 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, 2017.
- [26] Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena : A benchmark for efficient transformers. In *The Ninth International Conference on Learning Representations*, 2021.
- [27] Jiazhen Hong, Geoffrey Mackellar, and Soheila Ghane. Eegm2: An efficient mamba-2-based self-supervised framework for long-sequence eeg modeling. *arXiv preprint arXiv:2502.17873*, 2025.
- [28] Anna Tegon, Thorir Mar Ingolfsson, Xiaying Wang, Luca Benini, and Yawei Li. Femba: Efficient and scalable eeg analysis with a bidirectional mamba foundation model. *arXiv preprint arXiv:2502.06438*, 2025.
- [29] William O Tatum IV. *Handbook of EEG interpretation*. Springer Publishing Company, 2021.
- [30] Jimmy TH Smith, Andrew Warrington, and Scott W Linderman. Simplified state space layers for sequence modeling. In *ICLR*, 2023.
- [31] Yuhong Li, Tianle Cai, Yi Zhang, Deming Chen, and Debadeepa Dey. What makes convolutional models great on long sequence modeling? *arXiv preprint arXiv:2210.09298*, 2022.
- [32] Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization of diagonal state space models. *Advances in Neural Information Processing Systems*, 35:35971–35983, 2022.
- [33] Iyad Obeid and Joseph Picone. The temple university hospital eeg data corpus. *Frontiers in neuroscience*, 10:196, 2016.

- [34] Wenhui Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. Internimage: Exploring large-scale vision foundation models with deformable convolutions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14408–14419, 2023.
- [35] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [36] Ziyu Jia, Youfang Lin, Jing Wang, Ronghao Zhou, Xiaojun Ning, Yuanlai He, and Yaoshuai Zhao. Graphsleepnet: Adaptive spatial-temporal graph convolutional networks for sleep stage classification. In *Ijcai*, volume 2021, pages 1324–1330, 2020.
- [37] Jiquan Wang, Sha Zhao, Haiteng Jiang, Yangxuan Zhou, Zhenghe Yu, Tao Li, Shijian Li, and Gang Pan. Caresleepnet: a hybrid deep learning network for automatic sleep staging. *IEEE Journal of Biomedical and Health Informatics*, 2024.
- [38] Zheng Chen, Yasuko Matsubara, Yasushi Sakurai, and Jimeng Sun. Long-term eeg partitioning for seizure onset detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 14221–14229, 2025.
- [39] Chi-Sheng Chen, Ying-Jung Chen, and Aidan Hung-Wen Tsai. Large cognition model: Towards pretrained eeg foundation model. *arXiv preprint arXiv:2502.17464*, 2025.
- [40] Demetres Kostas, Stephane Aroca-Ouellette, and Frank Rudzicz. Bendlr: Using transformers and a contrastive self-supervised learning task to learn from massive amounts of eeg data. *Frontiers in Human Neuroscience*, 15:653659, 2021.
- [41] Daoze Zhang, Zhizhang Yuan, Yang Yang, Junru Chen, Jingjing Wang, and Yafeng Li. Brant: Foundation model for intracranial neural signal. *Advances in Neural Information Processing Systems*, 36:26304–26321, 2023.
- [42] Zhizhang Yuan, Daoze Zhang, Junru Chen, Gefei Gu, and Yang Yang. Brant-2: Foundation model for brain signals. *CoRR*, 2024.
- [43] Daoze Zhang, Zhizhang Yuan, Junru Chen, Kerui Chen, and Yang Yang. Brant-x: A unified physiological signal alignment framework. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 4155–4166, 2024.
- [44] Weibang Jiang, Yansen Wang, Bao-liang Lu, and Dongsheng Li. Neurolm: A universal multi-task foundation model for bridging the gap between language and eeg signals. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [45] Navid Mohammadi Foumani, Geoffrey Mackellar, Soheila Ghane, Saad Irtza, Nam Nguyen, and Mahsa Salehi. Eeg2rep: enhancing self-supervised eeg representation through informative masked inputs. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5544–5555, 2024.
- [46] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *54th Annual Meeting of the Association for Computational Linguistics*, pages 1715–1725. Association for Computational Linguistics (ACL), 2016.
- [47] Taku Kudo and John Richardson. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, 2018.
- [48] Syama Sundar Rangapuram, Matthias W Seeger, Jan Gasthaus, Lorenzo Stella, Yuyang Wang, and Tim Januschowski. Deep state space models for time series forecasting. *Advances in neural information processing systems*, 31, 2018.
- [49] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.

- [50] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. *Advances in neural information processing systems*, 33:1474–1487, 2020.
- [51] Jimmy TH Smith, Andrew Warrington, and Scott W Linderman. Simplified state space layers for sequence modeling. *arXiv preprint arXiv:2208.04933*, 2022.
- [52] Aniruddh Raghu, Payal Chandak, Ridwan Alam, John Guttag, and Collin M. Stultz. Sequential multi-dimensional self-supervised learning for clinical time series, 2023.
- [53] Tri Dao, Daniel Y Fu, Khaled K Saab, Armin W Thomas, Atri Rudra, and Christopher Ré. Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint arXiv:2212.14052*, 2022.
- [54] Linqi Zhou, Michael Poli, Winnie Xu, Stefano Massaroli, and Stefano Ermon. Deep latent state space models for time-series generation. In *International Conference on Machine Learning*, pages 42625–42643. PMLR, 2023.
- [55] Xuan-The Tran, Linh Le, Quoc Toan Nguyen, Thomas Do, and Chin-Teng Lin. Eeg-ssm: Leveraging state-space model for dementia detection. *arXiv preprint arXiv:2407.17801*, 2024.
- [56] Yiyu Gui, MingZhi Chen, Yuqi Su, Guibo Luo, and Yuchao Yang. Eegmamba: Bidirectional state space model with mixture of experts for eeg multi-task classification, 2024.
- [57] James W Cooley and John W Tukey. An algorithm for the machine calculation of complex fourier series. *Mathematics of computation*, 19(90):297–301, 1965.
- [58] Dani Kiyasseh, Tingting Zhu, and David A Clifton. Clocs: Contrastive learning of cardiac signals across space, time, and patients. In *International Conference on Machine Learning*, pages 5606–5615. PMLR, 2021.
- [59] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmLR, 2020.
- [60] JASPER HH. The ten-twenty electrode system of the international federation. *Electroenceph clin Neurophysiol*, 10:367–380, 1958.
- [61] Jingjing Chen, Xiaobin Wang, Chen Huang, Xin Hu, Xinkle Shen, and Dan Zhang. A large finer-grained affective computing eeg dataset. *Scientific Data*, 10(1):740, 2023.
- [62] Wei Liu, Jie-Lin Qiu, Wei-Long Zheng, and Bao-Liang Lu. Comparing recognition performance and robustness of multimodal deep learning models for multimodal emotion recognition. *IEEE Transactions on Cognitive and Developmental Systems*, 14(2):715–729, 2021.
- [63] 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.
- [64] Ji-Hoon Jeong, Jeong-Hyun Cho, Young-Eun Lee, Seo-Hyun Lee, Gi-Hwan Shin, Young-Seok Kweon, José del R Millán, Klaus-Robert Müller, and Seong-Whan Lee. 2020 international brain–computer interface competition: A review. *Frontiers in human neuroscience*, 16:898300, 2022.
- [65] Wajid Mumtaz. MDD Patients and Healthy Controls EEG Data (New). Figshare, November 2016.
- [66] Ali Hossam Shoeb. *Application of machine learning to epileptic seizure onset detection and treatment*. PhD thesis, Massachusetts Institute of Technology, 2009.
- [67] Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220, 2000.

- [68] Richard B Berry, Rita Brooks, Charlene E Gamaldo, Susan M Harding, Carole Marcus, Bradley V Vaughn, et al. The aasm manual for the scoring of sleep and associated events. *Rules, Terminology and Technical Specifications, Darien, Illinois, American Academy of Sleep Medicine*, 176(2012):7, 2012.
- [69] Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In *The Tenth International Conference on Learning Representations*, 2022.
- [70] Daniel Y Fu, Elliot L Epstein, Eric Nguyen, Armin W Thomas, Michael Zhang, Tri Dao, Atri Rudra, and Christopher Ré. Simple hardware-efficient long convolutions for sequence modeling. In *International Conference on Machine Learning*, pages 10373–10391. PMLR, 2023.
- [71] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- [72] 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, 2018.
- [73] Yonghao Song, Xueyu Jia, Lie Yang, and Longhan Xie. Transformer-based spatial-temporal feature learning for eeg decoding. *arXiv preprint arXiv:2106.11170*, 2021.

## A Interpretability Analyses of TFDUAL-TOKENIZER Representations

We conduct interpretability analyses of the decoupled TFDUAL-TOKENIZER. As case studies, we visualize selected tokens in relation to well-established physiological markers, showing that frequency codes tend to align with spectral rhythms, whereas temporal codes correspond to neural events. These qualitative examples highlight that the tokenizer does not discretize signals in an abstract manner, but captures domain-specific structures with clear clinical relevance. Finally, we complement these visualizations with a quantitative analysis of class-specific token usage on 4 datasets, demonstrating that the learned tokens are not only interpretable but also discriminative, contributing to improved downstream classification.

### A.1 Exploring Frequency Token Patterns in Relation to Spectral Rhythms

To illustrate the interpretability of the learned frequency tokens, we take sleep staging as an example downstream task and analyze token embeddings. We focus on N2 and N3 stages, since they are characterized by the most distinctive spectral rhythms in clinical sleep scoring: *spindle* in the sigma band (11-16 Hz) for N2, and *slow wave* in the delta band (0.5-4 Hz) for N3.

As shown in Figure 5, frequency code No. 2298 captures a sigma bump, with prominent peaks localized in the spindle range. This aligns with the definition of N2 [68], which identifies spindle as a key marker. Quantitatively, this code exhibited a clear sigma peak in approximately 15.4% of its assigned patches. Similarly, Figure 6 shows that frequency code No. 32 predominantly encodes delta dominance, a typical spectral feature of N3 sleep. This code displayed clear delta dominance in about 18.3% of its assigned segments. Taken together, these findings suggest that our frequency-branch tokenizer helps establish a frequency vocabulary for EEG, suggesting potential links between the learned codes and clinically interpretable oscillations beyond class frequency bias.

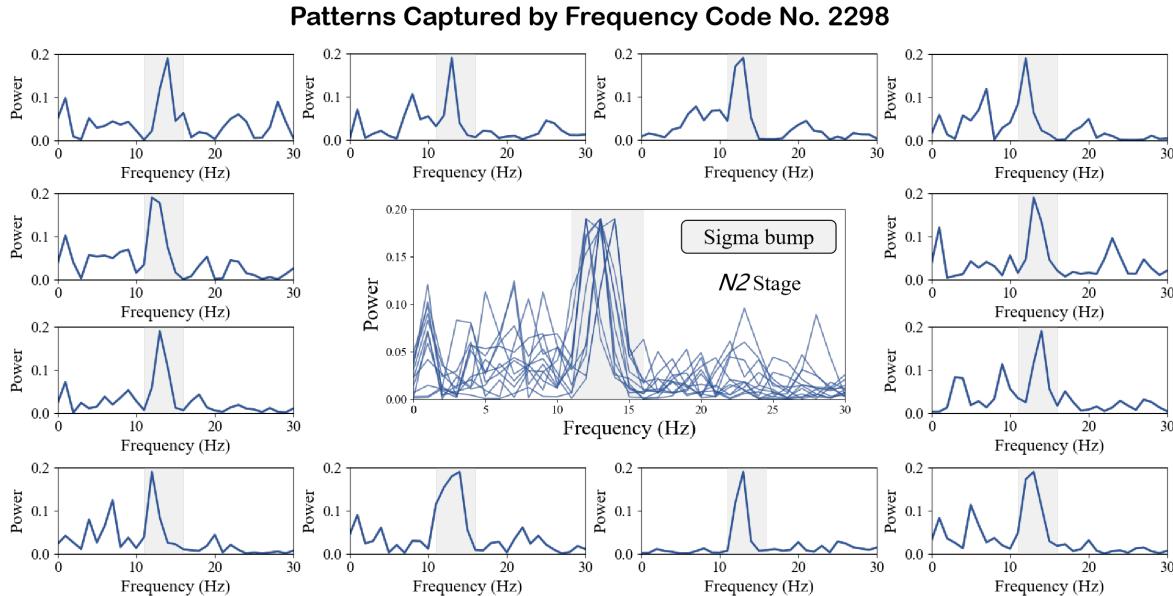


Figure 5: Frequency code capturing sigma bump, a typical spectral rhythm of N2 stage.

### A.2 Exploring Temporal Token Patterns in Relation to Neural Events

In addition to spectral rhythms, we further analyzed temporal codes on the sleep staging dataset to assess whether the tokenizer captures characteristic neural events in raw EEG. Figure 8 and Figure 7 show that certain codes are closely associated with K-complexes and sleep spindles, the hallmark waveforms of N2 sleep. Quantitatively,  $\approx 5.34\%$  of the assigned patches contained K-complexes and  $\approx 7.01\%$  contained spindles. These rates are consistent with clinical expectations (1–2 K-complexes

### Patterns Captured by Frequency Code No. 32

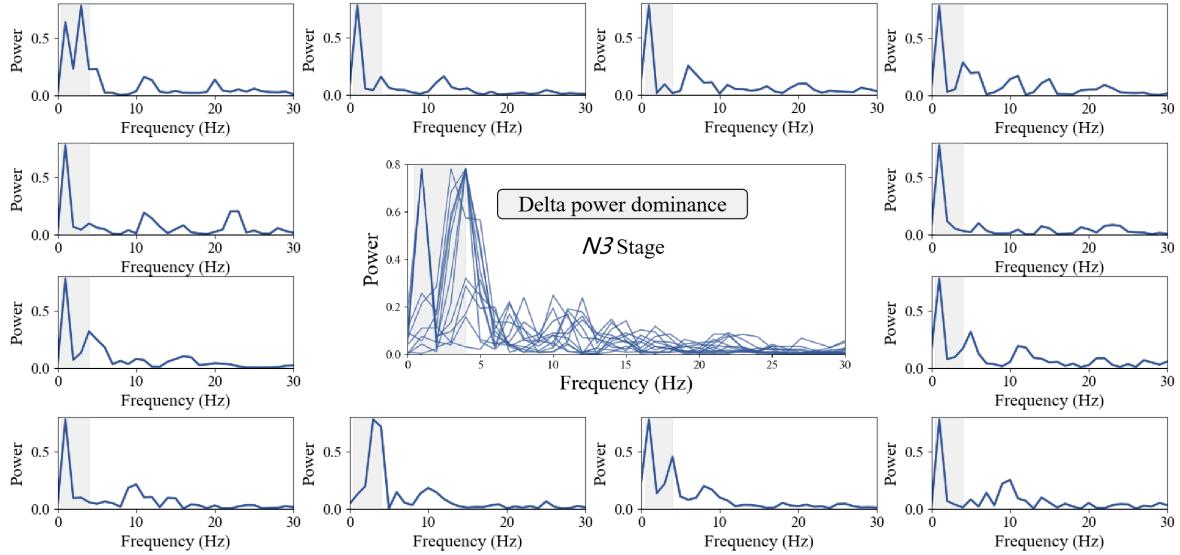


Figure 6: Frequency code capturing delta dominance, a typical spectral rhythm of N3 stage.

and 2–5 spindles per minute in N2 [68]) and, when considering the overall prevalence of N2 in five-class sleep staging, still substantially exceed what would be expected from random token assignment.

### Patterns Captured by Temporal Code No. 1459

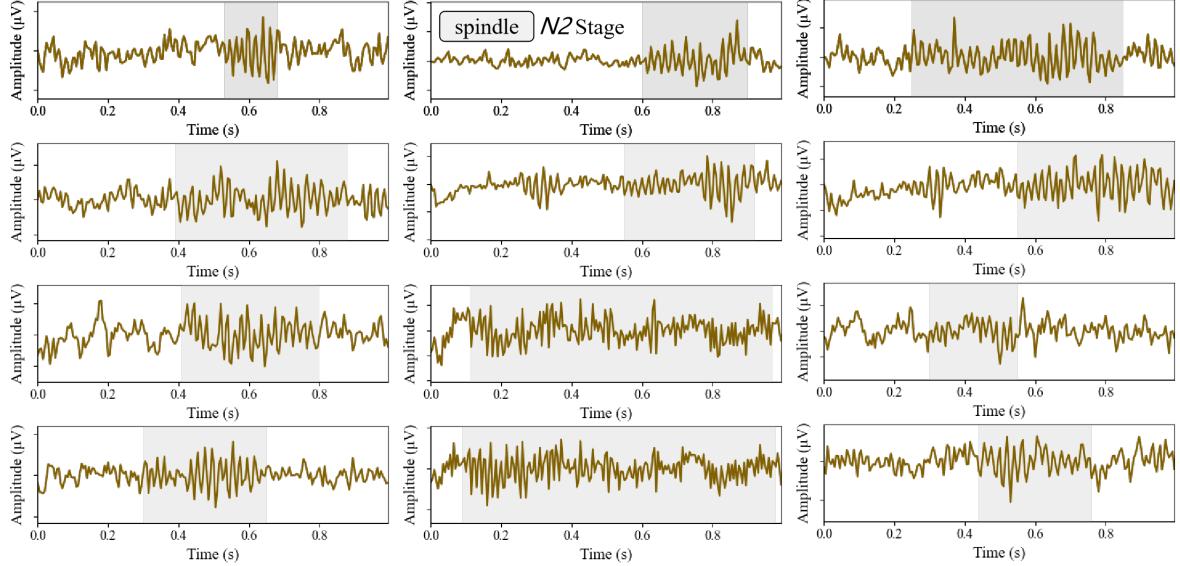


Figure 7: Temporal code capturing sleep spindle, a typical neural events of N2 stage.

Figure 9 shows that temporal code No. 1537 corresponds to slow waves, the defining feature of N3 sleep, occurring in  $\approx 40.4\%$  of its assigned patches. This exceeds the 20% per-epoch criterion for N3 scoring [68] and, given the overall prevalence of N3, indicates that this token provides a meaningful link to neural events.

Taken together, these observations suggest that the temporal branch of our tokenizer contributes to establishing a vocabulary of neurophysiological events, complementing the frequency-domain findings. While temporal waveforms are often noisier and harder to model than spectral rhythms [16], our decoupled design allows temporal tokens to emerge with meaningful associations to critical neural events. This indicates that the temporal branch offers useful links to clinically relevant waveforms and

### Patterns Captured by Temporal Code No. 134

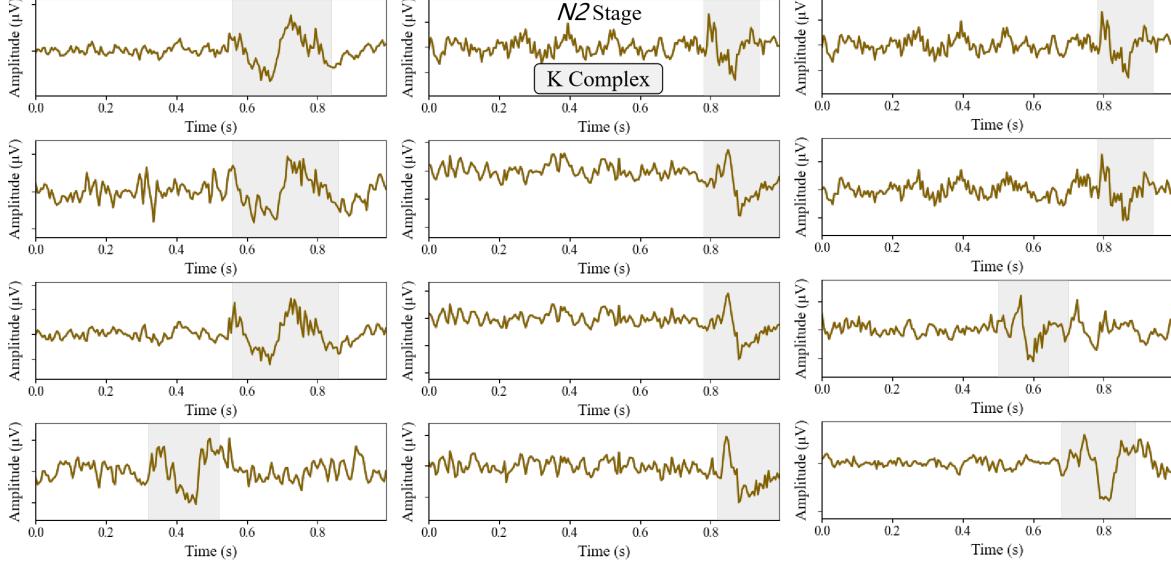


Figure 8: Temporal code capturing K complex, a typical neural events of N2 stage.

### Patterns Captured by Temporal Code No. 1537

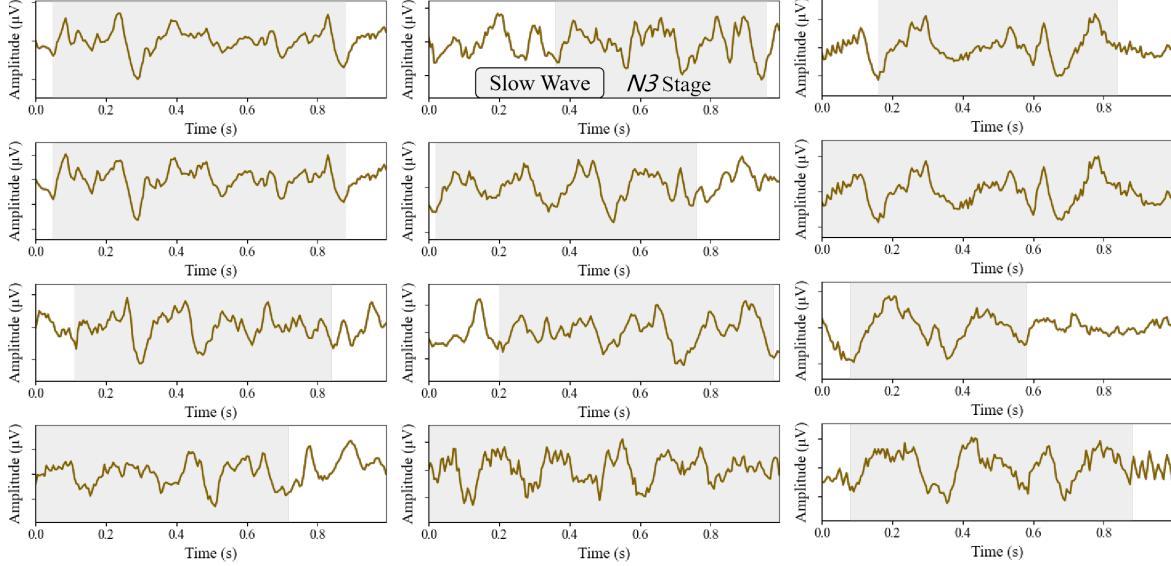


Figure 9: Temporal code capturing slow wave, a typical neural events of N3 stage.

demonstrates the effectiveness of our approach in capturing complementary structure.

### A.3 Class-Specific Token Ratio Analysis for TFDUAL-TOKENIZER

To better understand the contribution of the TFDUAL-TOKENIZER module, which decouples time and frequency domain representations, in learning discriminative EEG features, we conduct an analysis of *class-specific tokens*. A token (code) is considered class-specific if it predominantly appears in samples of a single class. Formally, for a given token  $c$ , let  $N_c^{(y)}$  denote the number of times  $c$  appears in class  $y$ .

The dominance ratio is defined as:

$$\text{Dominance}(c) = \frac{\max_y N_c^{(y)}}{\sum_y N_c^{(y)}} \quad (17)$$

If  $\text{Dominance}(c) > \tau$  (we use  $\tau = 1$ ), the token is deemed class-specific. The class-specific ratio for a codebook is then computed as:

$$\text{Class-Specific Token Ratio} = \frac{\# \text{ class-specific tokens}}{\text{Total tokens in codebook}} \quad (18)$$

Figure 10 illustrates the proportion of class-specific tokens derived from three configurations: using only the temporal codebook, only the frequency codebook, and a combination of both (TF-Decoupled). We employ two independent codebooks to capture complementary information via the proposed TFDUAL-TOKENIZER module.

Across all four datasets, including two for emotion recognition (FACED and SEED-V) and two for sleep staging (ISRUC\_S3 and ISRUC\_S1), the decoupled codebook consistently achieves the highest class-specific token ratios, reaching 54.7% on ISRUC\_S3 and 46.4% on FACED. These results confirm that decoupling temporal and frequency domain information significantly enhances the model’s ability to capture class-discriminative patterns.

Temporal and frequency domains exhibit inherently heterogeneous characteristics: temporal signals capture transient waveform dynamics, while frequency components reflect oscillatory structures over broader time scales. By decoupling these two domains and learning separate token representations, the model is able to preserve their distinct informational content. The TFDUAL-TOKENIZER then integrates temporally localized patterns—such as emotion-related bursts or sleep-stage-specific oscillations—that are highly class-discriminative. The TFDUAL-TOKENIZER enables richer and more effective EEG representations.

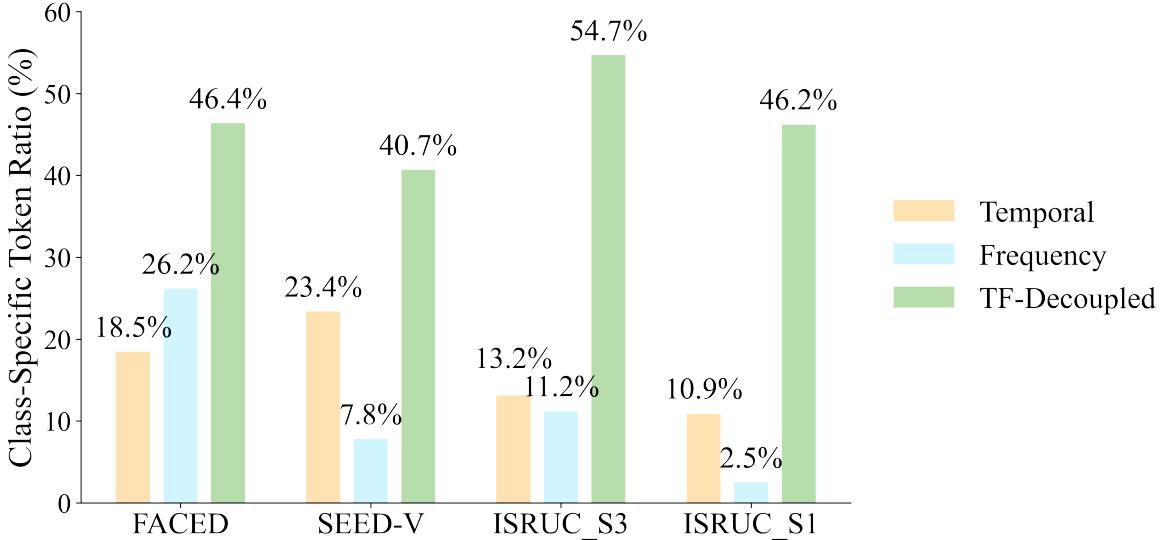


Figure 10: Class-Specific Code Ratio Across Different Codebooks.

## B Structure State Space Model

The state-space model is a classic model in control theory, and it represents the operational state of a system using first-order differential equations (ODE). A continuous state-space model can be defined in the following form:

$$x'(t) = Ax(t) + Bu(t), y(t) = Cx'(t) + Du(t), \quad (19)$$

where  $u(t)$  is a vector that represents the input of the system, while  $y(t)$  is a vector that represents the output of the system.  $x(t)$  and its derivative  $x'(t)$  represent the latent states of the system, typically in

the form of an N-D vector. And  $A, B, C, D$  here are the state, input, output, and feedforward matrices, defining the relationship between the input, output, and state vector. Following Gu et al. [69],  $\bar{D}$  is set equal to 0 since it can be replaced by the residual connection. Now, Equation 20 resembles an architecture similar to RNN, allowing us to recurrently compute  $x_k$ . Let the initial state be  $x_{k-1} = 0$ , and we can unroll Eq 20 as follows:

$$y_k = \overline{CA^k B u_0} + \overline{CA^{k-1} B u_1} + \dots + \overline{CA B u_{k-1}} + \overline{C B u_k} \quad (20)$$

$$y = \bar{K} u, \quad \bar{K} = (\overline{CB}, \overline{CA^1 B}, \dots, \overline{CA^k B}). \quad (21)$$

Therefore, SSM can be transformed from the form of a recurrent neural network to a convolutional neural network. During training, the  $\bar{K}$  can be considered as a 1-D globe convolution kernel, so the  $y$  can be calculated via the "long" convolution, allowing us to use the fast Fourier transform to efficiently compute the SSM convolutional kernel  $\bar{K}$ . However, directly computing the convolution in Equation 21 can be very expensive for long sequences. We can use the Fast Fourier Transform to accelerate it. The form of this convolution can be written as:

$$y = F_N^{-1} D_k F_N u, D_k = \text{diag}(\bar{K} F_N), \quad (22)$$

where  $F_N$  denotes the DFT matrix of size  $N$ . This FFT convolution has a computational complexity of  $O(n \log(n))$ . Following [31, 70], we hope to parameterize  $K$  directly rather than through  $\{A, B, C\}$  because we can eliminate complex parameterization and accelerate the entire convolution.

While Eq. 22 shows that the SSM can be computed using FFT-based convolution, it is important to formalize the complexity guarantee. We now state the following proposition.

**Proposition 1.** *Let  $(A, B, C, D)$  denote the discretized state-space matrices of an S4 layer, with input sequence  $u \in \mathbb{R}^N$ . The output  $y \in \mathbb{R}^N$  can be written as a linear convolution  $y = k * u$  with kernel*

$$k_0 = D, \quad k_n = CA^{n-1} B, \quad n \geq 1.$$

Using the convolution theorem and the FFT, this convolution can be computed in time  $\Theta(N \log N)$ .

*Proof.* This proof follows the proof in Lemma C.2 of (author?) [49]:

Expanding the recurrence

$$x_{n+1} = Ax_n + Bu_n, \quad y_n = Cx_n + Du_n, \quad (23)$$

yields

$$y_n = Du_n + \sum_{t=0}^{n-1} CA^{n-1-t} B u_t = \sum_{t=0}^n k_{n-t} u_t, \quad (24)$$

where  $k_0 = D$  and  $k_n = CA^{n-1} B$ . Thus  $y = k * u$ .

Let  $\mathcal{F}_N$  denote the  $N$ -point DFT. By the convolution theorem,

$$\mathcal{F}_N(y) = \mathcal{F}_N(k) \odot \mathcal{F}_N(u), \quad (25)$$

so that

$$y = \mathcal{F}_N^{-1}(\mathcal{F}_N(k) \odot \mathcal{F}_N(u)). \quad (26)$$

The cost of forward FFT, element-wise product, and inverse FFT is  $\Theta(N \log N)$ .

The S4 parameterization ensures  $A$  admits a diagonal-plus-low-rank (DPLR) structure, so the kernel takes the form

$$k_n = \sum_{j=1}^r \alpha_j \lambda_j^{n-1}, \quad r \ll N. \quad (27)$$

Each exponential sequence can be generated recursively in  $\mathcal{O}(N)$ , yielding all  $k_0, \dots, k_{N-1}$  in linear time.

Combining the above,

$$T(N) = \underbrace{\Theta(N \log N)}_{\text{FFT convolution}} + \underbrace{\mathcal{O}(N)}_{\text{kernel generation}} = \Theta(N \log N). \quad (28)$$

## C Theoretical Analysis of Decoupled Codebook Training

**Proposition 2.1.** Let  $X = (X_t, X_f)$  denote the temporal and frequency representations of an EEG segment, assumed approximately independent. Consider an additive reconstruction distortion  $d(x, \hat{x}) = d_t(x_t, \hat{x}_t) + d_f(x_f, \hat{x}_f)$ . Under a fixed total codebook size  $K = 2^R$ , a product codebook  $\mathcal{C}_t \times \mathcal{C}_f$  with  $|\mathcal{C}_t| = 2^{R_t}$  and  $|\mathcal{C}_f| = 2^{R_f}$ ,  $R_t + R_f = R$ , achieves a minimum expected distortion

$$D_{\text{prod}}^*(R) = \min_{R_t + R_f = R} (D_t^*(R_t) + D_f^*(R_f)), \quad (29)$$

which satisfies

$$D_{\text{mix}}^*(R) \geq D_{\text{prod}}^*(R), \quad (30)$$

where  $D_{\text{mix}}^*(R)$  is the minimum distortion of a single mixed codebook  $\mathcal{C}_{\text{mix}} \subset \mathbb{R}^{d_t+d_f}$  of size  $2^R$ .

*Proof.* The argument combines the separability of the rate–distortion (R–D) function for independent sources under additive distortion and the high-rate quantization approximation to the Shannon R–D limit.

Let  $X = (X_t, X_f)$  with  $X_t \perp X_f$ , and distortion

$$d(x, \hat{x}) = d_t(x_t, \hat{x}_t) + d_f(x_f, \hat{x}_f). \quad (31)$$

Then the Shannon R–D function satisfies

$$R(D) = \min_{D_t + D_f = D} (R_t(D_t) + R_f(D_f)), \quad (32)$$

where  $R_t(\cdot)$  and  $R_f(\cdot)$  are the marginal R–D functions for  $X_t$  and  $X_f$ . By convex duality, the optimal test channel factorizes:

$$p(\hat{x}_t, \hat{x}_f | x_t, x_f) = p(\hat{x}_t | x_t) p(\hat{x}_f | x_f), \quad (33)$$

and the optimal distortion allocation solves

$$\min_{D_t, D_f} R_t(D_t) + R_f(D_f) \quad \text{s.t. } D_t + D_f = D, \quad (34)$$

with KKT condition  $R'_t(D_t^*) = R'_f(D_f^*)$ .

At rate  $R_t$  bits for  $X_t$  and  $R_f$  bits for  $X_f$  (with  $R_t + R_f = R$ ), the minimal achievable distortions are  $D_t^*(R_t)$  and  $D_f^*(R_f)$ . In the high-rate regime, practical vector quantizers approach these Shannon limits, yielding

$$D_{\text{prod}}(R_t, R_f) \approx D_t^*(R_t) + D_f^*(R_f). \quad (35)$$

Optimizing over all feasible splits gives

$$D_{\text{prod}}^*(R) = \min_{R_t + R_f = R} (D_t^*(R_t) + D_f^*(R_f)). \quad (36)$$

Any mixed codebook  $\mathcal{C}_{\text{mix}} \subset \mathbb{R}^{d_t+d_f}$  with  $|\mathcal{C}_{\text{mix}}| = 2^R$  cannot beat the Shannon R–D limit, so

$$D_{\text{mix}}^*(R) \geq D^*(R). \quad (37)$$

But from (34)–(35),

$$D^*(R) = \min_{R_t + R_f = R} (D_t^*(R_t) + D_f^*(R_f)) = D_{\text{prod}}^*(R). \quad (38)$$

Therefore,

$$D_{\text{mix}}^*(R) \geq D_{\text{prod}}^*(R), \quad (39)$$

## D Pretraining Results

Our model follows a two-stage pretraining framework. In the first stage, we train the TFDUAL-TOKENIZER, which independently tokenizes EEG signals in both the temporal and frequency domains. This tokenizer is optimized to reconstruct the original raw EEG signals, amplitude, and phase components, thereby producing discrete code representations with structural interpretability. In the second stage, we pretrain the EEGSSM encoder using a masked modeling objective: given the original EEG signals as input, the model learns to predict the corresponding masked tokens generated from the TFDUAL-TOKENIZER.

This section reports the pretraining results of both stages, including loss convergence, reconstruction dynamics, and codebook utilization patterns.

## D.1 TFDUAL-TOKENIZER Pretraining Results

**Total Training Loss** The total pretraining loss curve of the TFDUAL-TOKENIZER is shown in Figure 11. The model demonstrates a rapid initial decrease in loss during the first few epochs, followed by a slower but consistent decline.

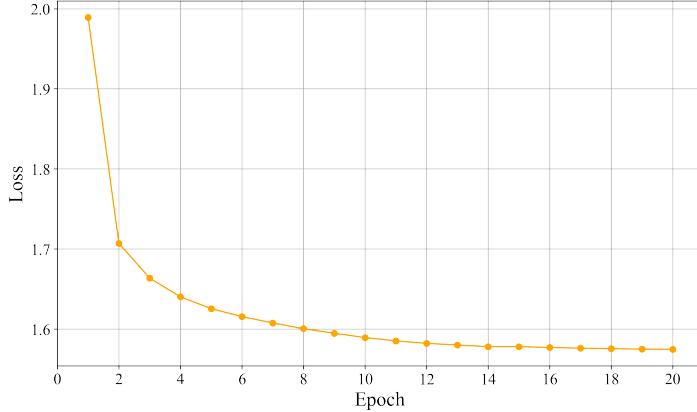


Figure 11: Pretraining Loss Curve of TFDUAL-TOKENIZER.

**Reconstruction Loss.** We report the pretraining reconstruction loss of the TFDUAL-TOKENIZER in Figure 12. The temporal codebook is trained to reconstruct raw EEG signals in the time domain, while the frequency codebook is trained to reconstruct the corresponding amplitude and phase components in the frequency domain. All three loss curves exhibit a sharp initial decline, followed by a gradual convergence, indicating stable optimization.

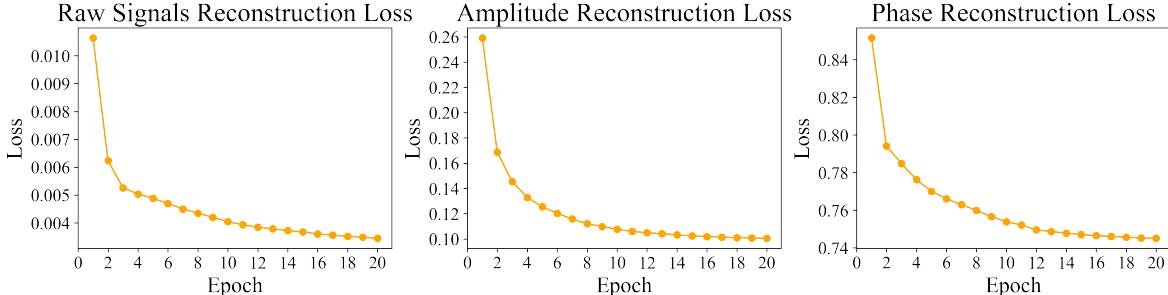


Figure 12: Pretraining Loss Curve of TFDUAL-TOKENIZER.

**Unused Codes Analysis.** During pretraining, we track the number of unused codes in both the temporal and frequency codebooks of the TFDUAL-TOKENIZER, each with a size of 4096. As shown in Figure 13, the frequency codebook demonstrates a rapid decrease in unused codes, while the temporal codebook shows a slower and more incremental reduction. A more detailed analysis of the temporal-frequency complementarity is provided in Section A.3.

## D.2 EEGSSM Pretraining Results

We plot the pretraining loss curve of EEGSSM in Figure 14. We select epoch 10 as the checkpoint for downstream fine-tuning. We observe that the pretraining loss of EEGSSM decreases rapidly from epoch 1 to 6 (9.04 → 6.39), then flattens gradually after epoch 10 (6.01 → 5.66). Using epoch 10 for fine-tuning is a balance between representation strength and generalization. Overtraining on EEG, prone to noise and inter-subject variability, can reduce transferability. Epoch 10 serves as a conservative yet effective checkpoint. This practice is consistent with trends in foundation models from NLP [13].

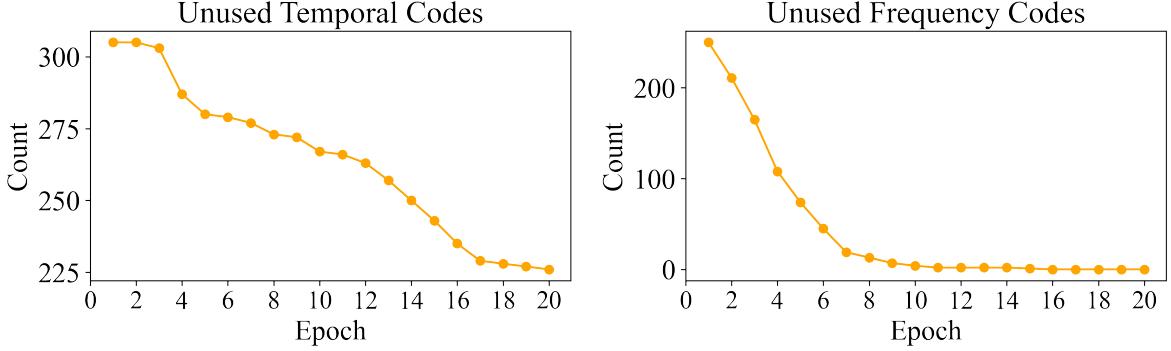


Figure 13: Unused code dynamics of the TFDUAL-TOKENIZER.

and vision [71], where mid-training checkpoints often lead to better downstream performance than final ones due to reduced overfitting to the pretext task.

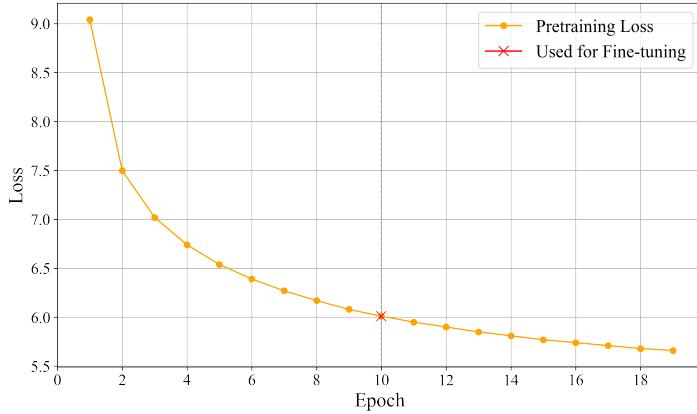


Figure 14: Pretraining Loss Curve of EEGSSM.

## E Improving Temporal Codebook Learning via Contrastive Loss

To improve the learning of the temporal codebook in our TFDUAL-TOKENIZER, we introduce a contrastive loss as one of the objectives during pretraining. This design is motivated by observations from prior work LaBraM [16], where the authors report that reconstructing raw EEG signals leads to unvergence, and thus omit the temporal reconstruction objective entirely.

To better understand this limitation, we first implemented a baseline reconstruction of raw signals within the LaBraM framework and observed that the training loss plateaued at a high value (between 0.128 and 0.131), showing no convergence over training epochs. To mitigate this, we introduce a *TFConv* module before the Transformer encoder, designed to extract temporal-frequency representations before tokenization. While this stabilizes training to some extent, we still observe significant issues with code utilization and loss convergence. Therefore, we incorporate a lightweight contrastive loss, applied over temporal representations before quantization, to encourage the model to organize similar input patterns closer in the latent space. As shown in Figure 15, this improves the optimization of the reconstruction loss and reduces the number of unused temporal codes during training.

These results demonstrate that contrastive regularization acts as an effective prior for stabilizing discrete token learning, particularly when reconstructing raw signals. It both improves convergence and mitigates codebook collapse in the temporal branch of the tokenizer.

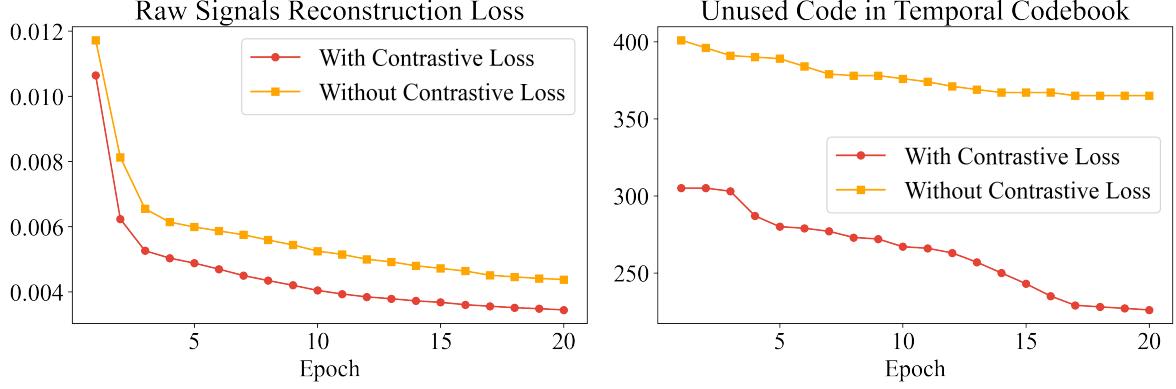


Figure 15: Effect of Contrastive Loss on Temporal Codebook Learning.

## F Dataset Description

We evaluate the CODEBRAIN across eight diverse downstream tasks covering ten publicly available EEG datasets. Notably, these datasets exhibit substantial variability in channel configurations (ranging from 6 to 64), sequence lengths (from 1s to 30s), and task complexities (2-class to 9-class), highlighting the versatility and robustness of CODEBRAIN across different EEG applications. The following sections describe each dataset in detail, including its task objective and data split strategy. A comprehensive analysis of each dataset is provided below.

**Emotion Recognition.** We conduct emotion recognition experiments on two widely used EEG datasets: FACED [61] and SEED-V [62].

The **FACED** dataset (Finer-Grained Affective Computing EEG Dataset) is a large-scale EEG dataset proposed by Chen et al. [61] for emotion recognition tasks. It consists of 32-channel EEG recordings sampled at 250 Hz from 123 participants, each exposed to 28 video clips designed to elicit nine distinct emotional states: *amusement, inspiration, joy, tenderness, anger, fear, disgust, sadness, and neutral*. These cover both positive and negative affective categories. Each EEG trial is 10 seconds long and is subsequently resampled to 200 Hz, resulting in a total of 10,332 clean EEG segments. For fair comparison, we adopt the same subject-wise split as in [15]: subjects 1–80 are used for training, 81–100 for validation, and 101–123 for testing, ensuring no subject overlap across splits and enabling evaluation of cross-subject generalization.

**SEED-V** [62] is an EEG dataset designed for emotion recognition, covering five emotional categories: *happy, sad, neutral, disgust, and fear*. It consists of 62-channel EEG recordings collected at 1000 Hz from 16 subjects, each participating in three sessions. Each session includes 15 trials, which are evenly divided into training, validation, and test sets (5 trials each). The EEG signals are segmented into 1-second windows, yielding a total of 117,744 samples, and resampled to 200 Hz for consistency. The dataset provides rich temporal structure and inter-subject variability, making it a strong benchmark for evaluating generalization in emotion-related EEG modeling.

**Sleep Staging.** We use two datasets, **ISRUC\_S1** and **ISRUC\_S3** [63], for sleep stage classification. Both datasets are annotated according to the American Academy of Sleep Medicine (AASM) standard [68], with five sleep stages: *Wake, NREM1 (N1), NREM (N2), NREM (N3), and REM*. Each EEG segment corresponds to a 30-second epoch.

**ISRUC\_S1** includes EEG recordings from 100 subjects using six channels at a sampling rate of 200 Hz. We adopt a subject-wise split, with 80 subjects for training, 10 for validation, and 10 for testing. As the transition rules between sleep stages carry important temporal patterns, we follow prior work [14, 15] and insert a Transformer layer on top of the projection head during fine-tuning to better capture sequence-level dependencies. We set the input sequence length to 20, and discard segments that cannot be evenly divided. In total, 86,320 labeled samples are retained.

**ISRUC\_S3** is a smaller dataset comprising recordings from 10 subjects, also sampled at 200 Hz with six channels, totaling 8,500 labeled segments. We follow an 8:1:1 subject-wise split for training,

validation, and testing.

**Imagined Speech Classification.** The **BCIC2020-T3** dataset [64] was released as part of the 2020 International Brain–Computer Interface Competition and focuses on imagined speech decoding. It contains EEG recordings from 15 participants who were instructed to silently imagine speaking five specific words or phrases, “hello”, “help me”, “stop”, “thank you”, and “yes”. EEG signals were collected using 64 scalp channels at a sampling rate of 256 Hz and were subsequently resampled to 200 Hz for preprocessing consistency. Each subject completed 80 trials per class, resulting in a total of 6,000 trials. The dataset provides predefined training, validation, and test splits, with 60, 10, and 10 trials per class, respectively, facilitating fair model evaluation like existing baselines [15].

**Mental Stress Detection.** The **Mental Arithmetic** dataset [65] supports the task of mental stress detection using EEG signals. It contains recordings from 36 subjects under two distinct cognitive conditions: *resting* and *active engagement* in mental arithmetic. EEG data labeled as “no stress” correspond to resting periods prior to the task, while “stress” labels are assigned to recordings during task performance. The signals were acquired using 20 electrodes placed according to the international 10–20 system, with an original sampling rate of 500 Hz. For consistency, the signals are resampled to 200 Hz and band-pass filtered between 0.5–45 Hz to suppress noise. Each recording is segmented into 5-second windows, yielding a total of 1,707 samples. We adopt a subject-wise split for fair evaluation with existing baselines [15]: subjects 1–28 for training, 29–32 for validation, and 33–36 for testing.

**Seizure Detection.** The **CHB-MIT** dataset [66] is a widely used benchmark for seizure detection from EEG signals. It contains long-term EEG recordings from 23 patients diagnosed with intractable epilepsy, collected at the Children’s Hospital Boston. The subjects underwent continuous monitoring over several days, during which seizures were recorded following the tapering of anti-epileptic medications. EEG signals were acquired using the international 10–20 system and originally sampled at 256 Hz. In our setting, we adopt 16 channels commonly used in prior work [24, 15], resample all signals to 200 Hz, and segment them into 10-second non-overlapping windows, yielding 326,993 labeled samples across seizure and non-seizure classes. We follow a subject-wise split: subjects 1–19 for training, 20–21 for validation, and 22–23 for testing. Notably, this dataset is highly imbalanced, with seizure events constituting only a small fraction of the total samples, posing significant challenges for model training and evaluation.

**Motor Imagery Classification.** The **SHU-MI** dataset [67] is designed for binary motor imagery classification, where participants are instructed to imagine movements of either the *left or right hand*. EEG signals were recorded from 25 subjects using a 32-channel setup at an original sampling rate of 250 Hz. To ensure consistency with the pre-training setting, all signals are resampled to 200 Hz and segmented into 4-second non-overlapping windows, resulting in 11,988 labeled samples. A subject-wise split is applied for fair model evaluation like existing baselines [15], with subjects 1–15 used for training, 16–20 for validation, and 21–25 for testing. This dataset supports the development of BCI systems that decode motor intentions from brain activity without actual movement.

**Event Type Classification.** The **TUEV** dataset [33] is a clinically annotated EEG corpus used for multi-class event type classification. It includes six event categories: *spike and sharp wave (SPSW)*, *generalized periodic epileptiform discharges (GPED)*, *periodic lateralized epileptiform discharges (PLED)*, *eye movements (EYEM)*, *artifacts (ARTF)*, and *background activity (BCKG)*. EEG signals were originally recorded at 256 Hz using 23 channels. In line with prior work [15], we preprocess the data by selecting 16 bipolar montage channels based on the international 10–20 system. The signals are band-pass filtered between 0.3–75 Hz to suppress low- and high-frequency noise, and a 60 Hz notch filter is applied to eliminate power line interference. All recordings are resampled to 200 Hz and segmented into 5-second windows, yielding 112,491 labeled samples. We follow the official train-test split and further divide the training subjects into training and validation sets in an 8:2 ratio, consistent with established benchmarks.

**Abnormal Detection.** The **TUAB** dataset [33] is employed for binary abnormal EEG detection, where each EEG recording is labeled as either *normal* or *abnormal* based on clinical interpretation. Originally recorded at 256 Hz using 23 channels, the dataset provides large-scale EEG recordings suitable for evaluating diagnostic models. To ensure fair comparison with prior work [15], we follow a similar preprocessing protocol. Specifically, we select 16 bipolar montage channels following the international 10–20 system, apply band-pass filtering between 0.3–75 Hz to eliminate low- and high-frequency artifacts, and remove 60 Hz power line interference using a notch filter. The EEG signals are then resampled to 200 Hz and segmented into 10-second windows, resulting in 409,455 labeled samples. We follow the official train-test split and further divide the training set into training and validation subsets using an 8:2 subject-wise ratio, consistent with existing benchmarks.

## G Baselines and Metrics Description

### G.1 Metrics

To comprehensively evaluate our model, we compare it with a set of strong baselines commonly used in EEG analysis. These baselines are evaluated using metrics tailored for class-imbalanced scenarios, which are prevalent in EEG datasets. The metrics include:

- **Balanced Accuracy**, which averages the recall across all classes and is particularly suitable for imbalanced multi-class classification tasks.
- **AUROC** and **AUC-PR**, which assess the performance of binary classifiers under different thresholds. While AUROC measures the trade-off between sensitivity and specificity, AUC-PR focuses on precision-recall trade-offs, especially informative under severe class imbalance.
- **Cohen’s Kappa**, which quantifies inter-class agreement beyond chance and is employed as the primary metric for multi-class classification.
- **Weighted F1 Score**, which combines precision and recall while adjusting for class support, ensuring fair performance measurement across imbalanced datasets.

For model selection and comparison, AUROC is used as the main evaluation metric for binary classification tasks, and Cohen’s Kappa is used for multi-class scenarios.

### G.2 Baselines

We compare our CODEBRAIN model against 5 publicly available EFM<sup>s</sup> that have released pre-trained weights, covering a diverse set of pretraining strategies.

We do not include traditional supervised baselines such as EEGNet [72] or ST-Transformer [73] in this comparison, as previous studies have consistently shown that EFMs significantly outperform these smaller supervised models. Our focus is therefore on evaluating the effectiveness of different foundation model designs and pretraining paradigms under comparable settings.

**BENDR** [40]: We adopted **BENDR (Bert-inspired Neural Data Representations)** as our baseline model, as introduced by Kostas et al. BENDR is a pioneering deep learning architecture for Electroencephalography (EEG) data, leveraging transformers and a contrastive self-supervised learning task. This approach enables the model to learn meaningful representations from vast amounts of unlabeled EEG data.

**BIOT** [24]: **BIOT (Biosignal Transformer for Cross-data Learning in the Wild)** is a transformer-based architecture designed to handle cross-dataset EEG signal classification under domain shifts. It leverages a domain-invariant attention mechanism and contrastive representation learning to enhance generalization across different recording conditions and subject populations.

**LaBraM**[16]: **LaBraM (Large Brain Model)** proposes a scalable transformer-based framework designed to learn generic EEG representations from large-scale brain signal datasets. By pretraining on a diverse corpus of EEG recordings, the model captures rich temporal and spatial features that transfer effectively to various downstream BCI tasks. The architecture incorporates efficient self-attention mechanisms and task-specific adapters to support flexible fine-tuning.

**EEGPT** [14]: **EEGPT** employs a dual self-supervised learning strategy that combines masked autoencoding with spatio-temporal representation alignment, enhancing feature quality by focusing

on high signal-to-noise ratio (SNR) representations rather than raw signals. The model’s hierarchical architecture decouples spatial and temporal processing, improving computational efficiency and adaptability to various brain-computer interface (BCI) applications.

**CBraMod** [15]: **CBraMod (Criss-Cross Brain Foundation Model)** is a transformer-based EEG foundation model that addresses the heterogeneous spatial and temporal dependencies inherent in EEG signals. It introduces a criss-cross transformer architecture comprising parallel spatial and temporal attention mechanisms, enabling separate yet simultaneous modeling of spatial and temporal relationships.

## H Hyperparameter Setting

We provide detailed hyperparameter configurations for the two-stage pretraining of our CODEBRAIN model and the fine-tuning settings across ten downstream tasks.

### H.1 Pretraining Settings

The pretraining process consists of two stages:

1. Training the TFDUAL-TOKENIZER
2. Training the EEGSSM

The hyperparameters used in each stage are summarized in Table 4 and Table 5, respectively.

Table 4: Hyperparameters for TFDUAL-TOKENIZER.

Hyperparameters	Values
<b>TFConv</b>	
Input channels	{1, 8, 4}
Output channels	{8, 4, 4}
Kernel size	{(1, 15), (1, 3), (1, 3)}
Stride	{(1, 8), (1, 1), (1, 1)}
Padding	{(0, 7), (0, 1), (0, 1)}
Transformer encoder layers	12
Transformer decoder layers	3
Hidden size	200
MLP size	800
Attention head number	8
Temporal Codebook size	$4096 \times 32$
Frequency Codebook size	$4096 \times 32$
Batch size	256
Peak learning rate	1e-4
Minimal learning rate	1e-5
Learning rate scheduler	Cosine
Optimizer	AdamW
Adam $\beta$	(0.9, 0.99)
Weight decay	1e-4
Total epochs	20
Data stride	200

### H.2 Parameters of EEGSSM

The model architecture parameters used in EEGSSM during pre-training are shown in Table 6.

Table 5: Hyperparameters of Pre-training.

Hyperparameters	Values
Epochs	10
Batch size	256
Dropout	0.1
Optimizer	Adam
Learning rate	1e-4
Adam $\beta$	(0.9, 0.999)
Adam $\epsilon$	1e-8
Weight decay	5e-3
Scheduler	CosineAnnealingLR
Minimal learning rate	1e-5
Clipping gradient norm	5

Table 6: Configuration of EEGSSM

Parameters	Values
Input size	200
Hidden dimension	200
Output size	200
Number of layers	8
Max sequence length	570
SGConv state	64
SGConv bidirectional	True
Layer normalization	True
Sliding window attention length	1s

### H.3 Fine-tuning Settings on Downstream Tasks

The CODEBRAIN model is fine-tuned on ten downstream EEG classification tasks, each with task-specific hyperparameters. Following the general strategy adopted by prior EFM, we adopt a lightweight three-layer MLP as the probe head for all downstream tasks and fine-tune the entire model end-to-end. Table 7 lists the fine-tuning configurations including learning rate, weight decay, dropout rate, and batch size for each task. For sleep staging tasks, due to their strong temporal structure, we follow prior work [14, 15] and insert an additional Transformer encoder on top of the projection head to jointly model the sequence of 20 consecutive EEG segments. This enables the model to capture inter-epoch transitions critical to sleep stage classification.

Table 7: Fine-tuning Hyperparameters for Downstream Tasks.

Dataset	Learning Rate	Weight Decay	Dropout	Batch Size
FACED	5e-5	5e-4	0.1	16
SEED-V	5e-5	1e-2	0.1	64
ISRUC_S1	1e-4	1e-1	0.2	48
ISRUC_S3	1e-4	1e-1	0.2	48
BCIC2020-T3	5e-5	5e-2	0.1	32
Mental Arithmetic	3e-5	1e-3	0.1	32
CHB-MIT	3e-5	1e-2	0.4	64
SHU-MI	5e-5	5e-3	0.3	64
TUEV	2e-5	5e-4	0.3	64
TUAB	1e-5	5e-5	0.4	512

## I Additional Evaluation on Other BCI Tasks

We report the performance of CODEBRAIN on four additional EEG datasets not included in the main text, covering diverse domains of sleep staging, motor imagery, event detection, and abnormality classification in Tables 8 to 11. These allow us to assess the cross-domain generalization ability of our pretrained model beyond the main text.

We note that both TUAB and TUEV originate from the TUH EEG corpus [33], which overlaps with our pretraining source (TUEG). To avoid overfitting to this distribution and promote generalization, we stop pretraining at epoch 10 as discussed in Section D.2. While this may limit gains on TUH datasets compared to previous EFM, such as CBraMod (trained for 40 epochs in the same pretraining dataset) [15], CODEBRAIN still achieves superior or competitive results.

**ISRUC\_S1** As shown in Table 8, CODEBRAIN achieves state-of-the-art performance on ISRUC\_S1 in terms of Cohen’s Kappa (0.7476) and Weighted F1 (0.8020), slightly surpassing CBraMod [15] by +0.34 and +0.09 points, respectively. Its Balanced Accuracy of 0.7835 is also competitive, trailing the best result by only -0.30. These results highlight the model’s ability to capture temporal dependencies and learn discriminative representations for 5-class sleep staging under a cross-subject setting.

Table 8: Performance Comparison on the ISURC\_S1 (5-Class) dataset.

Methods	Cohen’s Kappa	Weighted F1	Balanced Accuracy
BENDR[40]	$0.6956 \pm 0.0053$	$0.7569 \pm 0.0049$	$0.7401 \pm 0.0056$
BIOT[24]	$0.7192 \pm 0.0231$	$0.7790 \pm 0.0146$	$0.7527 \pm 0.0121$
LaBraM[16]	$0.7231 \pm 0.0182$	$0.7810 \pm 0.0133$	$0.7633 \pm 0.0102$
EEGPT[14]	$0.2223 \pm 0.0227$	$0.3111 \pm 0.0110$	$0.4012 \pm 0.0177$
CBraMod[15]	$0.7442 \pm 0.0152$	$0.8011 \pm 0.0099$	$0.7865 \pm 0.0110$
CodeBrain	<b><math>0.7476 \pm 0.0040</math></b>	<b><math>0.8020 \pm 0.0018</math></b>	$0.7835 \pm 0.0033$

**SHU-MI** As shown in Table 9, CODEBRAIN achieves the best overall performance on SHU-MI across all three metrics. It obtains an AUROC of 0.7124 and an AUC-PR of 0.7166, slightly improving over the previous best by +1.36 and +0.27 points, respectively. For Balanced Accuracy, it reaches 0.6431 (+0.61), with notably lower variance. These results underscore its strong generalization to motor imagery decoding under a cross-subject protocol.

Table 9: Performance Comparison on the SHU-MI (2-Class) dataset.

Methods	AUROC	AUC-PR	Balanced Accuracy
BENDR[40]	$0.5863 \pm 0.0280$	$0.5853 \pm 0.0268$	$0.5573 \pm 0.0227$
BIOT[24]	$0.6609 \pm 0.0127$	$0.6770 \pm 0.0119$	$0.6179 \pm 0.0183$
LaBraM[16]	$0.6604 \pm 0.0091$	$0.6761 \pm 0.0083$	$0.6166 \pm 0.0192$
EEGPT[14]	$0.6241 \pm 0.0071$	$0.6266 \pm 0.0133$	$0.5778 \pm 0.0162$
CBraMod[15]	$0.6988 \pm 0.0068$	$0.7139 \pm 0.0088$	$0.6370 \pm 0.0151$
CodeBrain	<b><math>0.7124 \pm 0.0050</math></b>	<b><math>0.7166 \pm 0.0106</math></b>	<b><math>0.6431 \pm 0.0066</math></b>

**TUEV** As shown in Table 10, CODEBRAIN achieves the highest Cohen’s Kappa (0.6912, an improvement of +0.0140 over the best baseline [15]) and Weighted F1 (0.8362) on TUEV. Although its Balanced Accuracy is lower than CBraMod [15], we attribute this to reduced sensitivity on the rare *SPSW* class. Since TUEV shares distributional overlap with our pretraining source (TUEG), we stop pretraining at epoch 10 to prevent overfitting, unlike CBraMod’s 40-epoch training as discussed in subsection D.2. We also report LaBraM’s results based on its original 23-channel setting [16], while

CODEBRAIN follows the 16-channel configuration used in CBraMod. Similarly, EEGPT [14] does not adopt a linear fine-tuning protocol but applies two convolutional layers before entering the foundation model, followed by an MLP head. While such architectural choices may enhance performance, we follow their respective fine-tuning settings. In addition, following the experimental setup of CBraMod, we also report the results of experiments conducted by removing all TUEV and TUAB samples from the TUEG dataset. It can be seen that although our model’s performance slightly declined on TUEV, it still surpassed other baselines. In such cases, CODEBRAIN still outperforming both LaBraM and EEGPT under their own fine-tuning settings clearly demonstrates the robustness of our approach.

Table 10: Performance Comparison on the TUEV (6-Class) dataset.

Methods	Cohen’s Kappa	Weighted F1	Balanced Accuracy
BENDR[40]	$0.4271 \pm 0.0238$	$0.6755 \pm 0.0216$	$0.4363 \pm 0.0245$
BIOT[24]	$0.5273 \pm 0.0249$	$0.7492 \pm 0.0082$	$0.5281 \pm 0.0225$
LaBraM[16]	$0.6637 \pm 0.0093$	$0.8312 \pm 0.0052$	$0.6409 \pm 0.0065$
EEGPT[14]	$0.6351 \pm 0.0134$	$0.8187 \pm 0.0063$	$0.6232 \pm 0.0114$
CBraMod[15]	$0.6772 \pm 0.0096$	$0.8342 \pm 0.0064$	<b><math>0.6671 \pm 0.0107</math></b>
CodeBrain (Excluding)	$0.6838 \pm 0.0291$	$0.8293 \pm 0.0163$	$0.6375 \pm 0.0182$
CodeBrain	<b><math>0.6912 \pm 0.0101</math></b>	<b><math>0.8362 \pm 0.0048</math></b>	$0.6428 \pm 0.0062$

**TUAB** As shown in Table 11, CODEBRAIN achieves the highest Balanced Accuracy (0.8294) on TUAB, slightly outperforming CBraMod [15]. Similar to TUEV, TUAB is part of the TUH EEG corpus family [33] and thus closely aligned with our pretraining source (TUEG). As discussed in subsection D.2, we adopt an early stopping strategy at epoch 10 to mitigate overfitting to this distribution, which may partly account for the slightly lower AUROC and AUC-PR compared to CBraMod, trained for 40 epochs on the same dataset. While LaBraM leverages a 23-channel montage [16] and EEGPT [14] employs two convolutional layers before the foundation model, we retain their respective fine-tuning protocols for comparison. In addition, the impact of duplicate data in the pre-training set on our model is smaller on the TUAB dataset, and in some experiments it can even surpass situations without leaking. This may be due to the larger size of the TUAB dataset. Despite these potentially stronger configurations, CODEBRAIN still exceeds both models under their own settings, highlighting its strong and consistent generalization.

Table 11: Performance Comparison on the TUAB (2-Class) dataset.

Methods	Balanced Accuracy	AUC-PR	AUROC
BENDR[40]	$0.7714 \pm 0.0248$	$0.8412 \pm 0.0215$	$0.8426 \pm 0.0237$
BIOT[24]	$0.7959 \pm 0.0057$	$0.8792 \pm 0.0023$	$0.8815 \pm 0.0043$
LaBraM[16]	$0.8140 \pm 0.0019$	$0.8965 \pm 0.0016$	$0.9022 \pm 0.0009$
EEGPT[14]	$0.8038 \pm 0.0040$	$0.8891 \pm 0.0018$	$0.8811 \pm 0.0015$
CBraMod[15]	$0.8289 \pm 0.0022$	<b><math>0.9258 \pm 0.0008</math></b>	<b><math>0.9227 \pm 0.0011</math></b>
CodeBrain (Excluding)	$0.8288 \pm 0.0064$	$0.9061 \pm 0.0039$	$0.9012 \pm 0.0020$
CodeBrain	<b><math>0.8294 \pm 0.0013</math></b>	$0.9100 \pm 0.0006$	$0.9030 \pm 0.0009$

## J Ablation on Design Choices

### J.1 Ablation on Mask Ratio

We conduct an ablation study to investigate the effect of the mask ratio in the EEGSSM pretraining framework. As shown in Tables 12–14 and Figure 16, downstream performance consistently exhibits a U-shaped trend with respect to the masking ratio across all three datasets: FACED, SEED-V, and ISRUC\_S3. Moderate masking (e.g., ratios around 0.4–0.6) leads to optimal performance, whereas excessively low (e.g., 0.1) or high (e.g., 0.9) ratios degrade generalization.

Table 12: Performance of CODEBRAIN on FACED Dataset under Different Mask Ratios.

Mask Ratio	Cohen’s Kappa	Weighted F1	Balanced Accuracy
0.1	0.5184 ± 0.0039	0.5746 ± 0.0158	0.5692 ± 0.0064
0.2	0.5239 ± 0.0036	0.5834 ± 0.0071	0.5821 ± 0.0040
0.3	0.5327 ± 0.0100	0.5859 ± 0.0036	0.5836 ± 0.0076
0.4	0.5391 ± 0.0045	0.5938 ± 0.0057	0.5904 ± 0.0058
0.5	<b>0.5406 ± 0.0084</b>	<b>0.5953 ± 0.0113</b>	<b>0.5941 ± 0.0098</b>
0.6	0.5295 ± 0.0075	0.5822 ± 0.0090	0.5793 ± 0.0112
0.7	0.5242 ± 0.0077	0.5800 ± 0.0065	0.5744 ± 0.0096
0.8	0.5157 ± 0.0065	0.5564 ± 0.0067	0.5528 ± 0.0040
0.9	0.5034 ± 0.0078	0.5457 ± 0.0084	0.5451 ± 0.0114

Table 13: Performance of CODEBRAIN on SEED-V Dataset under Different Mask Ratios.

Mask Ratio	Cohen’s Kappa	Weighted F1	Balanced Accuracy
0.1	0.2523 ± 0.0051	0.4081 ± 0.0030	0.3968 ± 0.0031
0.2	0.2633 ± 0.0033	0.4137 ± 0.0027	0.4071 ± 0.0035
0.3	0.2703 ± 0.0043	0.4200 ± 0.0048	0.4121 ± 0.0041
0.4	0.2734 ± 0.0042	<b>0.4244 ± 0.0029</b>	<b>0.4142 ± 0.0055</b>
0.5	<b>0.2735 ± 0.0032</b>	0.4235 ± 0.0022	0.4137 ± 0.0023
0.6	0.2699 ± 0.0041	0.4207 ± 0.0059	0.4106 ± 0.0048
0.7	0.2642 ± 0.0045	0.4158 ± 0.0033	0.4091 ± 0.0040
0.8	0.2603 ± 0.0031	0.4116 ± 0.0041	0.4009 ± 0.0025
0.9	0.2534 ± 0.0036	0.4085 ± 0.0024	0.3993 ± 0.0055

To further illustrate this pattern, we visualize the training loss curves across different mask ratios in Figure 17. Interestingly, higher mask ratios result in slower convergence and higher final training loss, which is expected due to the increased difficulty of the reconstruction task. In contrast, lower mask ratios lead to faster and smoother loss reduction, but do not necessarily yield better downstream performance. This observation suggests a possible *optimization-vs-generalization trade-off*: easier pretext tasks (low mask ratio) are more optimizable but may encourage the model to learn shortcut solutions with limited generalizability, while overly difficult tasks (high mask ratio) may hinder effective representation learning due to insufficient learning signal. Moderate masking strikes a balance by being sufficiently challenging to promote abstraction, while still being learnable, thereby facilitating better generalization across downstream tasks.

Table 14: Performance of CODEBRAIN on ISRUC\_S3 Dataset under Different Mask Ratios.

Mask Ratio	Cohen's Kappa	Weighted F1	Balanced Accuracy
0.1	$0.7252 \pm 0.0069$	$0.7865 \pm 0.0039$	$0.7686 \pm 0.0072$
0.2	$0.7403 \pm 0.0023$	$0.7982 \pm 0.0026$	$0.7752 \pm 0.0032$
0.3	$0.7501 \pm 0.0057$	$0.8013 \pm 0.0062$	$0.7766 \pm 0.0020$
0.4	$0.7608 \pm 0.0060$	$0.8168 \pm 0.0078$	$0.7846 \pm 0.0062$
0.5	<b><math>0.7671 \pm 0.0091</math></b>	$0.8202 \pm 0.0071$	<b><math>0.7856 \pm 0.0031</math></b>
0.6	$0.7661 \pm 0.0068$	<b><math>0.8219 \pm 0.0052</math></b>	$0.7826 \pm 0.0076$
0.7	$0.7577 \pm 0.0068$	$0.8161 \pm 0.0040$	$0.7782 \pm 0.0061$
0.8	$0.7471 \pm 0.0031$	$0.8090 \pm 0.0041$	$0.7658 \pm 0.0025$
0.9	$0.7504 \pm 0.0078$	$0.8125 \pm 0.0063$	$0.7668 \pm 0.0065$

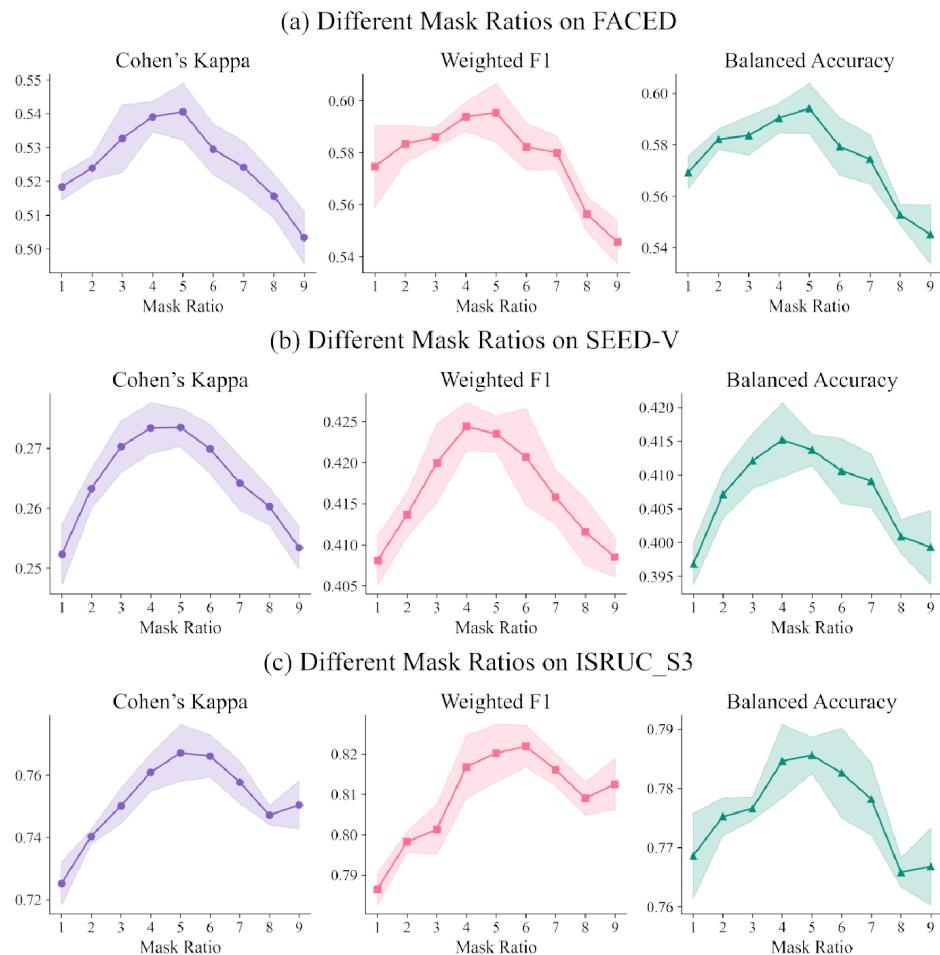


Figure 16: Performance Across Different Mask Ratios on FACED, SEED-V, and ISRUC\_S3.

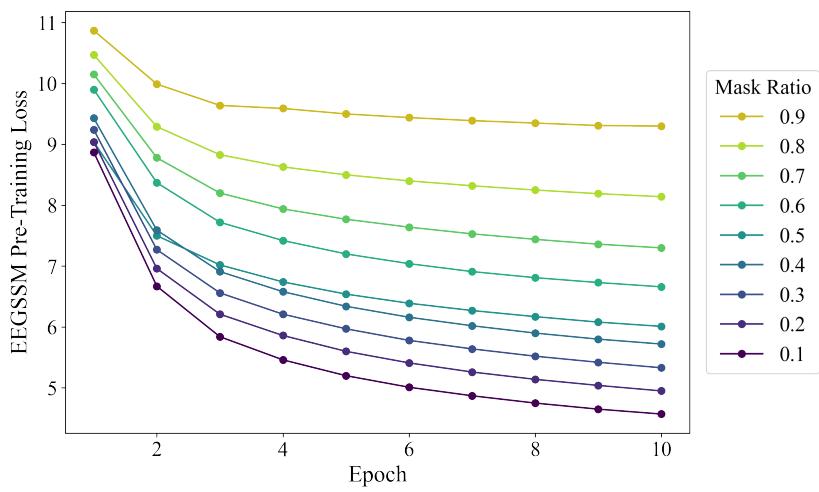


Figure 17: EEGSSM Pre-Training Loss Curve for Different Mask Ratios.

## J.2 Ablation on SWA Window Size

We conduct an ablation study to investigate the effect of the SWA window size in the EEGSSM framework. The window size of SWA means the length of the segment observed by the attention mechanism in the SWA mechanism. Usually, the window size of SWA is an odd number because the model generally focuses not only on a single segment but also needs to observe features in adjacent segments. Therefore, its size is usually  $2n + 1$ , where  $n$  represents the length of SWA. As shown in Figure 18 and Figure 19, SWA window size = 1 achieved the best performance on both datasets among the 5 SWA window sizes. With the increase of the SWA window size, the SWA is close to self-attention. Therefore, as the SWA window size increases, the model's performance actually improves because more data will be involved in the attention calculation. When the SWA window size is 1, it is a special case where SWA only calculates for one second. This is equivalent to the model performing attention mechanism calculations within one second. Overall, SWA can help the model achieve a certain performance improvement.

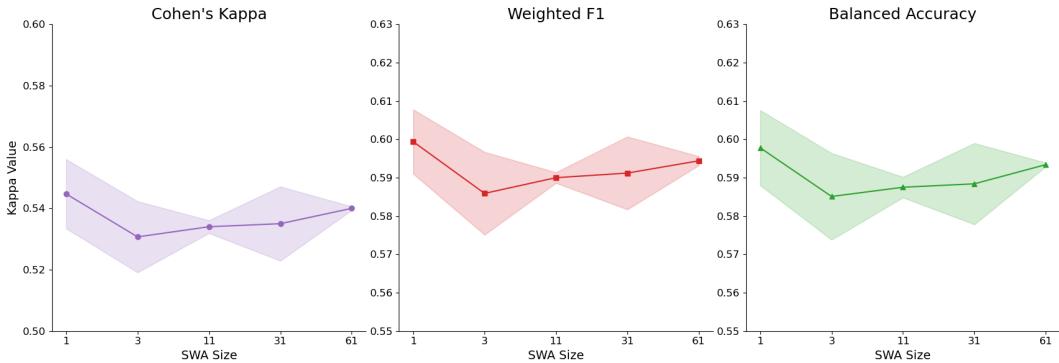


Figure 18: Performance of different SWA window sizes on the FACED dataset.

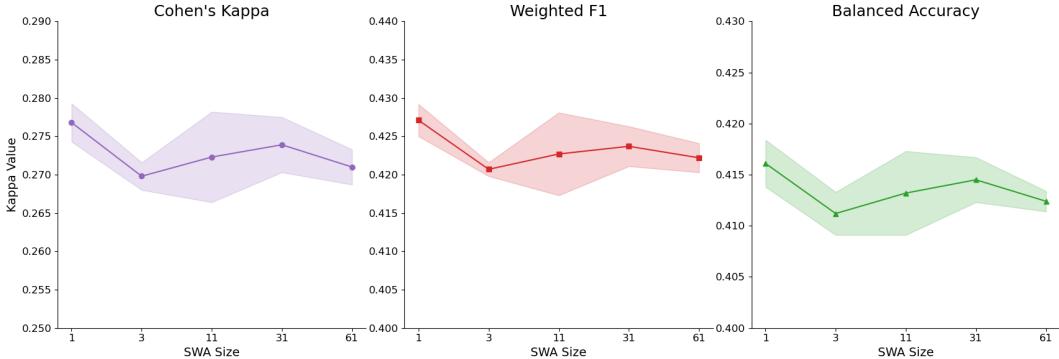


Figure 19: Performance of different SWA window sizes on the SEED-V dataset.

## J.3 Ablation on Codebook Size

The size of the codebook is an important parameter; a codebook that is too large may lead to unstable training, while a codebook size that is too small may result in mixed information. Ideally, the codebook should maintain a small amount of unused codes but not be 0. We tested several combinations of different time-domain and frequency-domain codebook sizes to observe their unused codes during Tokenizer training.

Table 15 shows the unused temporal codes and unused frequency codes under different codebook size combinations. When both the Temporal codebook size and Frequency codebook size are set to 2048, the unused frequency code is 0, indicating that there are duplicate frequency codes in the current codebook. When choosing a codebook size of 8192, the unused temporal codes reached 3401 and the

unused frequency codes reached 1260, indicating that a large number of codes in the completed training codebook were not used. We ultimately selected 4096 for both. Although the frequency domain often yields richer representations, enlarging its codebook may increase reliance on it, so we kept temporal and frequency codebooks equal. Considering capacity and utilization, 4096–4096 offers the best trade-off.

Table 15: Temporal and frequency codebook statistics

Temporal Codebook Size	Frequency Codebook Size	Unused Temporal Codes	Unused Frequency Codes
2048	2048	12	0
2048	4096	116	0
<b>4096</b>	<b>4096</b>	225	165
4096	8192	321	931
8192	8192	3401	1260

#### J.4 Ablation on Patch Size

The size of the patch window is also an important adjustable parameter, which affects temporal resolutions and masking strategies. To explore the impact of patch window size on the model, we used window sizes ranging from 0.5s to 5s for complete two-stage pre-training and full-parameter fine-tuning. Note that for some datasets where the patch size is larger than the channel length, such as SEED-V, we pad the portion exceeding the available data to match the patch size in this experiment.

Table 16 and Table 17 show the performance of our method with different patch sizes on the SEED-V dataset and ISRUC\_3 dataset. Notably, in SEED-V, patch lengths longer than 1s require heavy padding, causing large performance drops; in ISRUC\_S3, the shortest 0.5s patches achieve the worst performance, likely because they fragment key waveforms in sleep staging (e.g., Spindle  $\geq 0.5$ s). From these results, the 1s setting is supported by two key considerations:

**Broad compatibility with downstream task.** 1s is a divisor of most downstream sequence lengths (1–30s), minimizing padding and ensuring transferability. For example, on the SEED-V dataset, if the patch size chosen by the model is greater than 1s, some methods (such as padding) need to be adopted to enable model training. These methods may usually impair the model’s performance because they introduce additional noise or increase computational load. EEG datasets with durations less than 1 second are relatively rare, as most datasets have at least 1 second of data. If the data does not exactly match the whole seconds, the cost of processing such a dataset is also relatively small. Prior EEG foundation models (e.g., LaBram [16], CBraMod [15]) also adopt 1s patches for this reason.

**Semantic integrity.** Choosing a 1s patch length preserves the natural structure of EEG waveforms and prevents semantic fragmentation. Many physiologically meaningful EEG events have characteristic durations: for example, K-complexes are typically around 1s, spindles last 0.5–2s, and event-related potentials such as P300 occur in the range of 0.3–0.6s. If patches are shorter than these characteristic scales, the temporal branch of the TFDual-Tokenizer may only capture partial fragments of these waveforms, leading to loss of semantic context.

Table 16: Performance on SEED-V with different patch sizes

Patch Length	Cohen’s Kappa	Weighted F1	Balanced Accuracy
0.5s	$0.2640 \pm 0.0035$	$0.4035 \pm 0.0025$	$0.3841 \pm 0.0025$
<b>1s</b>	<b><math>0.2735 \pm 0.0032</math></b>	<b><math>0.4235 \pm 0.0022</math></b>	<b><math>0.4137 \pm 0.0023</math></b>
2s	$0.1545 \pm 0.0046$	$0.3313 \pm 0.0036$	$0.3279 \pm 0.0049$
5s	$0.1557 \pm 0.0043$	$0.3340 \pm 0.0029$	$0.3271 \pm 0.0045$

#### K Backbone Efficiency Comparison

To evaluate the computational efficiency of our proposed SGConv module, we conduct an ablation study by replacing it with three common sequence modeling modules: CNN, LSTM, and Transformer in EEGSSM block. We compare their model sizes, floating-point operations (FLOPs), and iteration times, as shown in Figure 20. Specifically, the CNN variant uses a 3-layer depthwise separable convolution

Table 17: Performance on ISRUC\_3 with different patch sizes

Patch Length	Cohen’s Kappa	Weighted F1	Balanced Accuracy
0.5s	0.7405±0.0102	0.7950±0.0097	0.7420±0.0081
1s	<b>0.7671±0.0091</b>	<b>0.8202±0.0071</b>	<b>0.7856±0.0031</b>
2s	0.7592±0.0079	0.8113±0.0081	0.7753±0.0087
5s	0.7601±0.0075	0.8096±0.0074	0.7791±0.0096

block, while the LSTM and Transformer variants use a single layer of standard LSTM and Transformer Encoder (implemented by Pytorch), respectively. In terms of parameter count, SGConv contains 15.17M parameters, fewer than CNN (16.22M), LSTM (17.35M), and Transformer (21.2M). For FLOPs, SGConv also achieves the lowest computational cost at 8.74G, compared to Transformer’s 27.79G. Regarding iteration time, SGConv is slightly slower than Transformer and CNN models in terms of training speed, but it outperforms the LSTM model. In summary, SGConv effectively reduces the number of parameters while maintaining computational complexity, which helps the model to be trained and inferred on smaller GPUs.

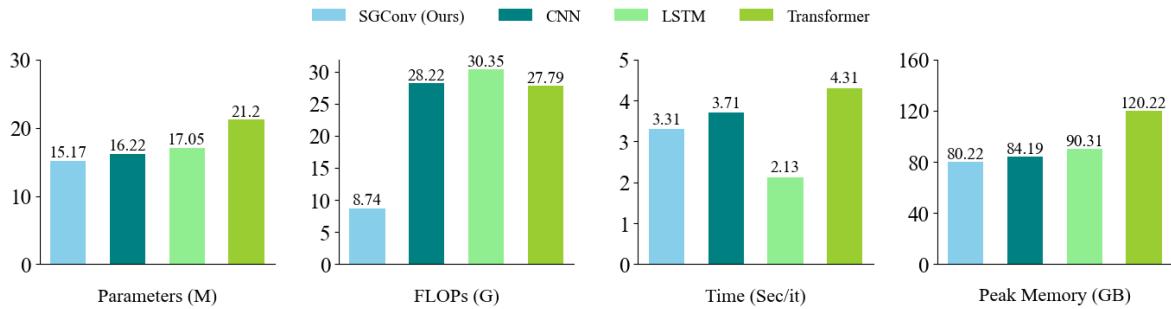


Figure 20: Computational Overhead of Using Different Backbones in the EEGSSM Module.

## L Channel Robustness

In real-world scenarios, the collection of EEG often encounters situations where channels are missing, especially when using machines from different manufacturers. To test the model’s performance on datasets with missing channel data, we randomly mask some channels in the training data for full parameter fine-tuning. We selected the FACED and SEED-V datasets for experiments because they represent short-sequence and long-sequence cases respectively, and their numbers of channels are relatively complete.

We evaluate the performance of the CodeBrain model and CBraMod model in scenarios with missing channels. We conducted three different experiments, randomly masking 12.5%, 25%, and 50% of the channels in each experiment, respectively. Figure 21 shows results of the experiment. It can be seen that our CodeBrain model outperforms the CBraMod model in all channel masking scenarios. On the FACED dataset, our model’s performance after masking 25% of the channels is still close to that of CBraMod without masking. The performance decline of the CBraMod model is also faster than that of our model. This trend is even more pronounced on the SEED-V dataset. Through this experiment, we can demonstrate that the channel robustness of CodeBrain is stronger than that of the CBraMod model, retaining most of its performance even in cases of channel failure.

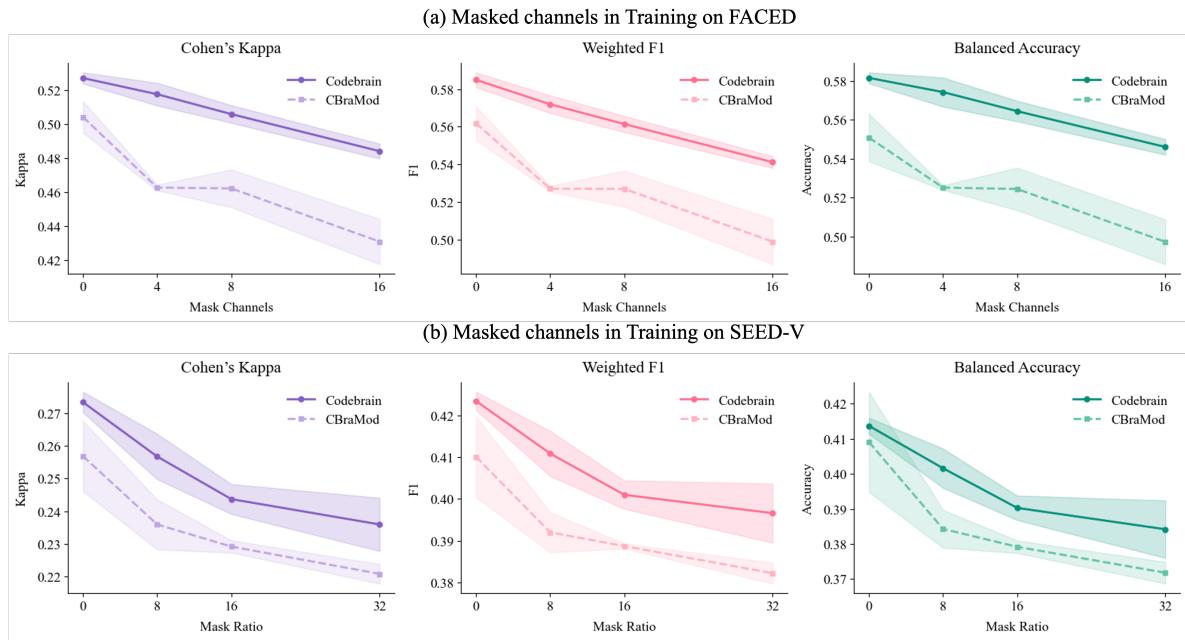


Figure 21: Performance after randomly masking different numbers of channels during the full parameter fine-tuning stage.

## M Detailed Results on Scaling Laws

We provide the detailed scaling law results for both data and model size across three representative EEG datasets (FACED, SEED-V, and ISRUC\_S3), covering three evaluation metrics. For brevity, only Cohen’s kappa scores are included in the main text, while full results are provided in this section. Prior work [15, 16] has explored the effect of scaling EEG foundation models using 1 to 1000 hours of pretraining data. We extend this analysis in two key dimensions:

1. Scaling the pretraining data volume from 1k up to 9k hours.
2. Investigating model scaling by varying the depth of the EEGSSM encoder from 3 layers (3.86M parameters) to 24 layers (146.75M parameters) and the hidden size from 128 to 384.

### M.1 Scaling Laws with Respect to Training Data Volume

We examine how the volume of pretraining data influences the downstream performance of CODEBRAIN. Specifically, we scale the pretraining duration from 1k to 9k hours and evaluate the resulting models on three downstream datasets: FACED, SEED-V, and ISRUC\_S3. Detailed quantitative results across three evaluation metrics (Cohen’s kappa, weighted F1 score, and balanced accuracy) are presented in Tables 18, 19, and 20, respectively. As shown in Figure 22, increasing the amount of pretraining data generally leads to consistent improvements across all datasets and metrics. On FACED and ISRUC\_S3, performance gains are steady throughout the entire range up to 9k hours, while on SEED-V, the trend is more modest and plateaus after 5k hours. These results highlight the importance of large-scale data for representation learning in EEG and suggest that further scaling may continue to yield performance benefits.

Table 18: Training Data Scaling Laws of CODEBRAIN on FACED Dataset.

Training Data	Cohen’s Kappa	Weighted F1	Balanced Accuracy
1000 Hours	$0.5014 \pm 0.0107$	$0.5462 \pm 0.0146$	$0.5452 \pm 0.0163$
2000 Hours	$0.5133 \pm 0.0068$	$0.5540 \pm 0.0120$	$0.5521 \pm 0.0133$
3000 Hours	$0.5189 \pm 0.0086$	$0.5688 \pm 0.0104$	$0.5687 \pm 0.0083$
4000 Hours	$0.5208 \pm 0.0032$	$0.5741 \pm 0.0041$	$0.5713 \pm 0.0024$
5000 Hours	$0.5171 \pm 0.0040$	$0.5692 \pm 0.0131$	$0.5661 \pm 0.0105$
6000 Hours	$0.5273 \pm 0.0085$	$0.5803 \pm 0.0064$	$0.5764 \pm 0.0069$
7000 Hours	$0.5328 \pm 0.0116$	$0.5854 \pm 0.0082$	$0.5809 \pm 0.0097$
8000 Hours	$0.5336 \pm 0.0082$	$0.5875 \pm 0.0071$	$0.5835 \pm 0.0072$
9000 Hours	$\mathbf{0.5406 \pm 0.0084}$	$\mathbf{0.5953 \pm 0.0113}$	$\mathbf{0.5941 \pm 0.0098}$

Table 19: Training Data Scaling Laws of CODEBRAIN on SEED-V Dataset.

Training Data	Cohen’s Kappa	Weighted F1	Balanced Accuracy
1000 Hours	$0.2584 \pm 0.0044$	$0.3946 \pm 0.0063$	$0.3799 \pm 0.0096$
2000 Hours	$0.2648 \pm 0.0062$	$0.4042 \pm 0.0091$	$0.3961 \pm 0.0084$
3000 Hours	$0.2689 \pm 0.0055$	$0.4117 \pm 0.0075$	$0.4028 \pm 0.0082$
4000 Hours	$0.2672 \pm 0.0043$	$0.4121 \pm 0.0054$	$0.4022 \pm 0.0049$
5000 Hours	$0.2678 \pm 0.0069$	$0.4120 \pm 0.0084$	$0.4026 \pm 0.0067$
6000 Hours	$0.2669 \pm 0.0031$	$0.4113 \pm 0.0067$	$0.4030 \pm 0.0102$
7000 Hours	$0.2686 \pm 0.0043$	$0.4129 \pm 0.0052$	$0.4027 \pm 0.0048$
8000 Hours	$0.2703 \pm 0.0049$	$0.4165 \pm 0.0034$	$0.4094 \pm 0.0054$
9000 Hours	$\mathbf{0.2735 \pm 0.0032}$	$\mathbf{0.4235 \pm 0.0022}$	$\mathbf{0.4137 \pm 0.0023}$

Table 20: Training Data Scaling Laws of CODEBRAIN on the ISRUC\_S3 Dataset.

Training Data	Cohen's Kappa	Weighted F1	Balanced Accuracy
1000 Hours	$0.7340 \pm 0.0187$	$0.7826 \pm 0.0154$	$0.7505 \pm 0.0192$
2000 Hours	$0.7347 \pm 0.0051$	$0.7869 \pm 0.0075$	$0.7524 \pm 0.0032$
3000 Hours	$0.7540 \pm 0.0106$	$0.8012 \pm 0.0093$	$0.7694 \pm 0.0089$
4000 Hours	$0.7590 \pm 0.0042$	$0.8108 \pm 0.0074$	$0.7753 \pm 0.0065$
5000 Hours	$0.7610 \pm 0.0058$	$0.8124 \pm 0.0102$	$0.7794 \pm 0.0086$
6000 Hours	$0.7681 \pm 0.0125$	$0.8190 \pm 0.0100$	$0.7856 \pm 0.0089$
7000 Hours	$0.7648 \pm 0.0079$	$0.8170 \pm 0.0076$	$0.7845 \pm 0.0076$
8000 Hours	$0.7668 \pm 0.0121$	$0.8182 \pm 0.0109$	$0.7851 \pm 0.0116$
9000 Hours	<b><math>0.7671 \pm 0.0091</math></b>	<b><math>0.8202 \pm 0.0071</math></b>	<b><math>0.7856 \pm 0.0031</math></b>

In addition, we visualize the pretraining optimization behavior across different data scales in Figure 23. As expected, larger pretraining data consistently lead to lower training loss, indicating more effective representation learning. Notably, the convergence curves become progressively smoother and more stable as training data volume increases, suggesting improved optimization stability in large-scale regimes. While smaller training data volumes (e.g., 1k–3k hours) show relatively high starting loss and slower convergence, larger training data volumes (6k–9k hours) reach lower final losses and exhibit diminishing returns, aligning with trends observed in downstream performance. These findings provide further empirical support for the scalability of EEG foundation models and reinforce the role of large data in enhancing both optimization and generalization.

## M.2 Scaling Laws with Respect to Model Size

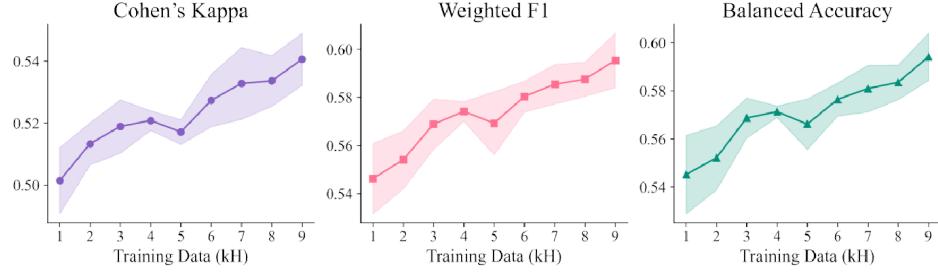
We further investigate how model parameters affect downstream performance by scaling the number of layers in the EEGSSM encoder from 3 to 8, resulting in parameter counts ranging from 6.82M to 15.17M. Detailed results across the FACED, SEED-V, and ISRUC\_S3 datasets are provided in Tables 21, 22, and 23, respectively. To visualize the trend more clearly, Figure 24 presents the performance curves as model size increases. Across all three datasets and evaluation metrics, we observe a consistent performance gain as the model size increases. The improvements are particularly pronounced on the FACED and ISRUC\_S3 datasets, where all three metrics show steady growth up to the largest model. In contrast, performance on SEED-V improves more modestly and begins to plateau beyond 13.5M parameters. These results suggest that increasing model capacity can enhance generalization ability, especially for datasets with richer structure or more complex temporal dynamics, while also indicating that optimal scaling may be task-dependent.

Table 21: Model Size Scaling Laws of CODEBRAIN on the FACED Dataset.

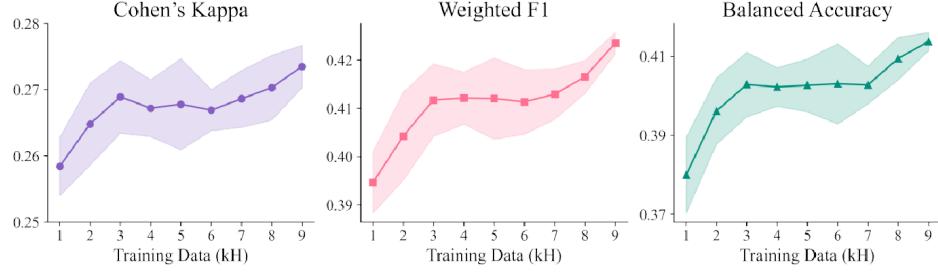
Layer	Hidden size	Params	Cohen's Kappa	Weighted F1	Balanced Accuracy
3	128	3.86M	$0.4786 \pm 0.0131$	$0.5231 \pm 0.0164$	$0.5287 \pm 0.0207$
3	200	6.82M	$0.4818 \pm 0.0165$	$0.5362 \pm 0.0113$	$0.5317 \pm 0.0182$
4	200	8.49M	$0.4988 \pm 0.0082$	$0.5506 \pm 0.0091$	$0.5497 \pm 0.0084$
5	200	10.16M	$0.5096 \pm 0.0049$	$0.5686 \pm 0.0104$	$0.5642 \pm 0.0067$
6	200	11.83M	$0.5244 \pm 0.0113$	$0.5705 \pm 0.0085$	$0.5778 \pm 0.0095$
7	200	13.50M	$0.5314 \pm 0.0069$	$0.5872 \pm 0.0080$	$0.5846 \pm 0.0067$
8	200	15.17M	$0.5406 \pm 0.0084$	$0.5953 \pm 0.0113$	$0.5941 \pm 0.0098$
12	256	33.15M	$0.5478 \pm 0.0013$	$0.5912 \pm 0.0031$	$0.5901 \pm 0.0128$
24	384	146.75M	<b><math>0.5503 \pm 0.0120</math></b>	<b><math>0.5964 \pm 0.0178</math></b>	<b><math>0.5985 \pm 0.0233</math></b>

To better understand the optimization behavior during pretraining, we plot the training loss curves

(a) Training Data Scaling Laws on FACED



(b) Training Data Scaling Laws on SEED-V



(c) Training Data Scaling Laws on ISRUC\_S3

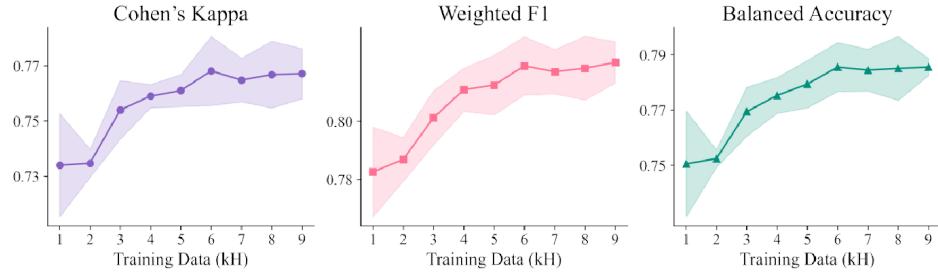


Figure 22: Training Data Scaling Laws on FACED, SEED-V, and ISRUC\_S3.

Table 22: Model Size Scaling Laws of CODEBRAIN on the SEED-V Dataset.

Layer	Hidden size	Params	Cohen's Kappa	Weighted F1	Balanced Accuracy
3	128	3.86M	$0.2576 \pm 0.0047$	$0.3969 \pm 0.0042$	$0.3896 \pm 0.0027$
3	200	6.82M	$0.2609 \pm 0.0078$	$0.4004 \pm 0.0112$	$0.3956 \pm 0.0098$
4	200	8.49M	$0.2638 \pm 0.0080$	$0.4108 \pm 0.0091$	$0.4030 \pm 0.0121$
5	200	10.16M	$0.2645 \pm 0.0102$	$0.4127 \pm 0.0158$	$0.4013 \pm 0.0084$
6	200	11.83M	$0.2663 \pm 0.0056$	$0.4202 \pm 0.0068$	$0.4079 \pm 0.0055$
7	200	13.50M	$0.2724 \pm 0.0057$	$0.4211 \pm 0.0060$	$0.4120 \pm 0.0086$
8	200	15.17M	$0.2735 \pm 0.0032$	$0.4235 \pm 0.0022$	$0.4137 \pm 0.0023$
12	256	33.15M	$0.2807 \pm 0.0029$	$0.4317 \pm 0.0036$	$0.4182 \pm 0.0028$
24	384	146.75M	$0.2831 \pm 0.0033$	$0.4342 \pm 0.0030$	$0.4216 \pm 0.0031$

for different model sizes in Figure 25. As expected, larger models consistently achieve lower final training loss, indicating stronger capacity to fit the pretraining objective. The loss reduction is particularly evident when increasing from 3 to 6 layers, while the gain starts to saturate beyond 7 layers. Even when increased to 24 layers, the reduction in pre-training loss brought by more than 100M parameters is not significant. This trend mirrors the downstream performance in Figure 24 and Tables 21–23, suggesting that both optimization efficiency and generalization benefit from increased model size—though with diminishing returns as parameter count grows. These results reinforce the scalability of EEGSSM and

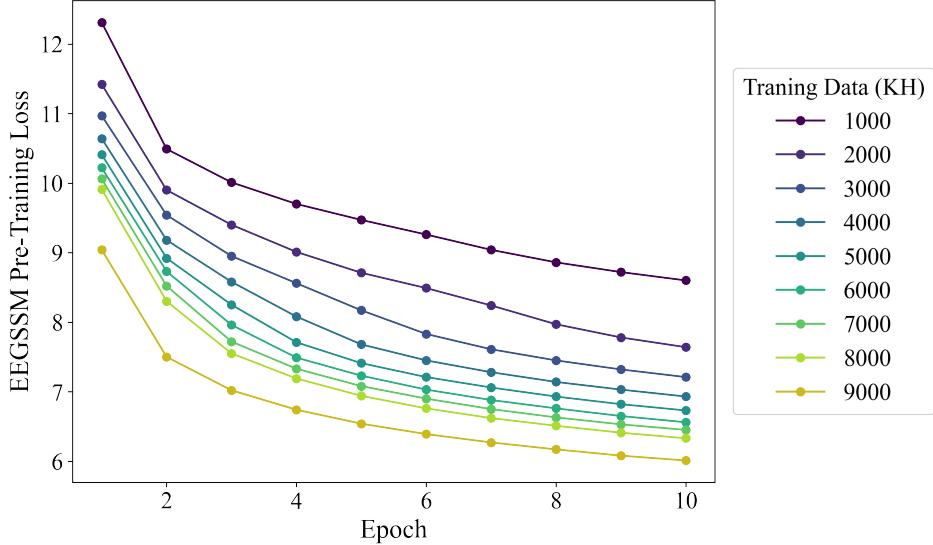


Figure 23: EEGSSM Pre-Training Loss Curve for Different Training Data Volume.

Table 23: Model Size Scaling Laws of CODEBRAIN on the ISRUC\_S3 Dataset.

Layer	Hidden size	Params(M)	Cohen's Kappa	Weighted F1	Balanced Accuracy
3	128	3.86M	0.7434±0.0087	0.7813±0.0107	0.7493±0.0058
3	200	6.82M	0.7456±0.0083	0.7862±0.0104	0.7514±0.0084
4	200	8.49M	0.7486±0.0014	0.7942±0.0040	0.7604±0.0035
5	200	10.16M	0.7516±0.0025	0.7985±0.0056	0.7639±0.0079
6	200	11.83M	0.7570±0.0082	0.8064±0.0068	0.7734±0.0061
7	200	13.50M	0.7620±0.0045	0.8153±0.0091	0.7824±0.0100
8	200	15.17M	0.7671±0.0091	0.8202±0.0071	0.7856±0.0031
12	256	33.15M	0.7753±0.0113	0.8316±0.0074	0.7940±0.0061
24	384	146.75M	<b>0.7791±0.0108</b>	<b>0.8352±0.0072</b>	<b>0.8008±0.0040</b>

underscore the importance of balancing capacity with task-specific requirements.

Our results demonstrate consistent improvements in downstream performance as both data volume and model capacity increase, suggesting that EEG foundation models may continue to benefit from further scaling, similar to trends observed in vision and language domains.

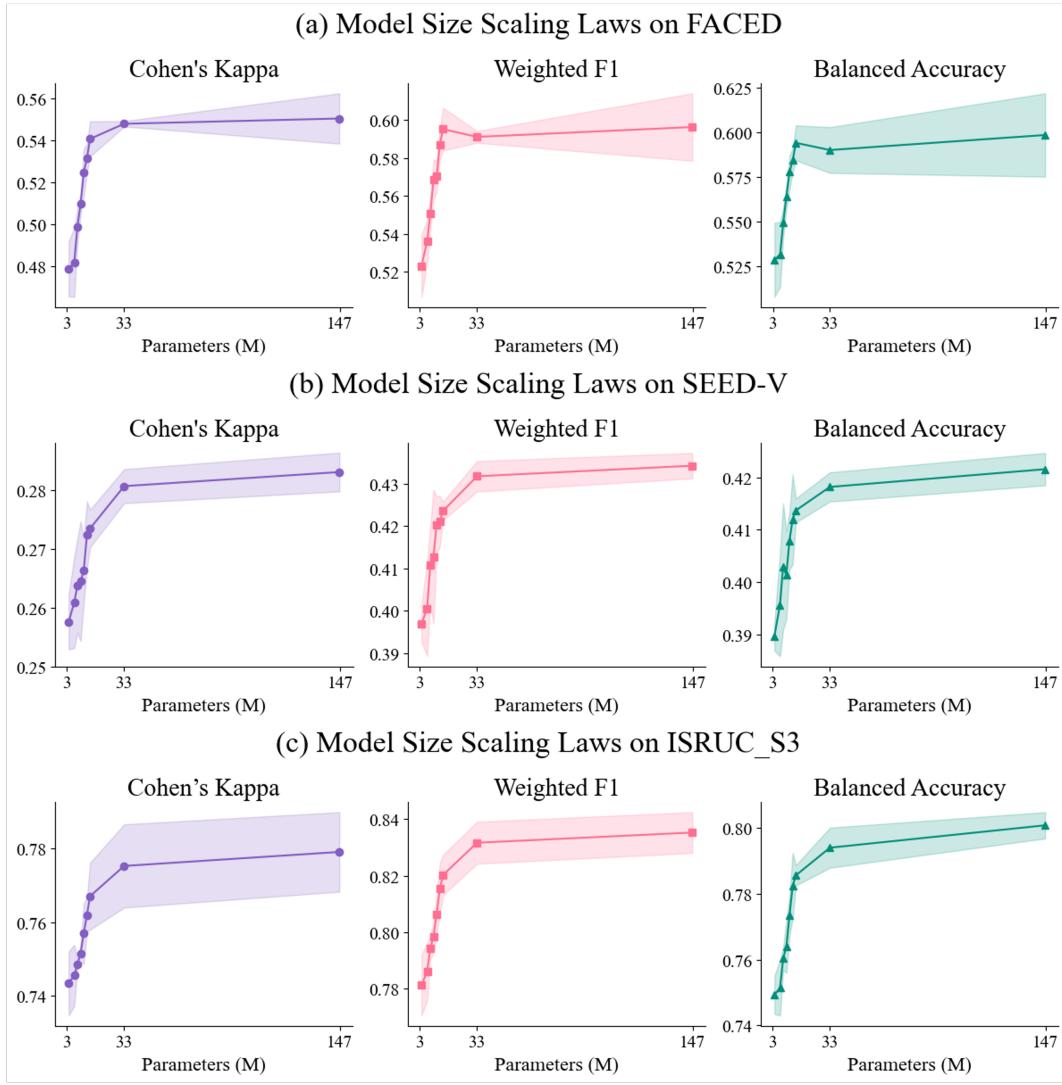


Figure 24: Model Size Scaling Laws on FACED, SEED-V, and ISRUC\_S3.

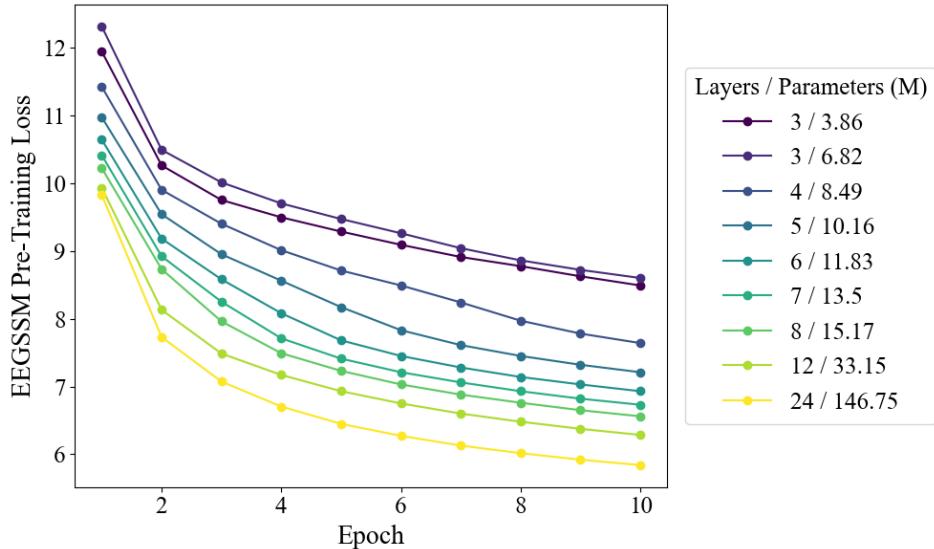


Figure 25: EEGSSM Pre-Training Loss Curve for Different Model Size.

## N Low-Resource Comparison With Existing Methods

To evaluate the effectiveness of our model under data-scarce conditions, we conduct experiments where only 30% of the training data is available. Tables 24 and 25 present the results on the FACED (9-class) and SEED-V (5-class) datasets, respectively, comparing our CODEBRAIN model with several strong baselines. We observe that CODEBRAIN consistently achieves the best performance among all methods under the 30% setting across all metrics. These improvements are highlighted with underlines in the tables. For reference, we also include the full-data results of CODEBRAIN. These results highlight CODEBRAIN’s strong data efficiency and suggest its potential for real-world EEG applications where annotated data is often scarce.

Table 24: Comparison under Partial Data Setting (30%) on the FACED Dataset (9-class).

Methods	Cohen’s Kappa	Weighted F1	Balanced Accuracy
BIOT (30%)	$0.2573 \pm 0.0346$	$0.3501 \pm 0.0341$	$0.3428 \pm 0.0329$
LaBraM (30%)	$0.2672 \pm 0.0371$	$0.3548 \pm 0.0325$	$0.3513 \pm 0.0315$
CBraMod (30%)	$0.3239 \pm 0.0265$	$0.4056 \pm 0.0256$	$0.4035 \pm 0.0233$
CodeBrain (30%)	<u><math>0.3356 \pm 0.0253</math></u>	<u><math>0.4114 \pm 0.0225</math></u>	<u><math>0.4104 \pm 0.0281</math></u>
CodeBrain (full)	<b><math>0.5406 \pm 0.0084</math></b>	<b><math>0.5953 \pm 0.0113</math></b>	<b><math>0.5941 \pm 0.0098</math></b>

Table 25: Comparison under Partial Data Setting (30%) on the SEED-V Dataset (5-class).

Methods	Cohen’s Kappa	Weighted F1	Balanced Accuracy
BIOT (30%)	$0.1775 \pm 0.0425$	$0.3492 \pm 0.0416$	$0.3505 \pm 0.0375$
LaBraM (30%)	$0.2044 \pm 0.0384$	$0.3700 \pm 0.0321$	$0.3686 \pm 0.0305$
CBraMod (30%)	$0.2291 \pm 0.0246$	$0.3886 \pm 0.0255$	$0.3877 \pm 0.0236$
CodeBrain (30%)	<u><math>0.2376 \pm 0.0284</math></u>	<u><math>0.3943 \pm 0.0259</math></u>	<u><math>0.3902 \pm 0.0271</math></u>
CodeBrain (full)	<b><math>0.2735 \pm 0.0032</math></b>	<b><math>0.4235 \pm 0.0022</math></b>	<b><math>0.4137 \pm 0.0023</math></b>