

---

# SAMBA: TOWARD A LONG-CONTEXT EEG FOUNDATION MODEL VIA SPATIAL EMBEDDING AND DIFFERENTIAL MAMBA

---

**Jiazhen Hong\***  
Emotiv Research  
Melbourne, Australia  
jiazhen@emotiv.com

**Geoffrey Mackellar**  
Emotiv Research  
Sydney, Australia  
geoff@emotiv.com

**Soheila Ghane**  
Emotiv Research  
Melbourne, Australia  
soheila@emotiv.com

## ABSTRACT

Long-sequence electroencephalogram (EEG) modeling is essential for developing generalizable EEG representation models. This need arises from the high sampling rate of EEG data and the long recording durations required to capture extended neurological patterns in brain activity. Transformer-based models have shown promise in modeling short sequences of a few seconds; however, their quadratic complexity limits scalability to longer contexts. Moreover, variability in electrode montage across available datasets, along with inter-subject differences in brain signals, pose significant challenges to developing a generalizable and robust foundation model. We propose *SAMBA*, a self-supervised learning framework with a Mamba-based U-shaped encoder-decoder architecture, which effectively captures long-range temporal dependencies and spatial variability in EEG data. Leveraging the inherent ability of Mamba in processing long context sizes, we introduce: (1) *Temporal Semantic Random Masking* for semantic-level sequence reconstruction, (2) a *Multi-Head Differential Mamba* module to suppress redundancy and emphasize salient temporal structures, and (3) a *Spatial-Adaptive Input Embedding* that learns unified embeddings in a three-dimensional Euclidean space, enabling robustness across devices. Experiments on thirteen EEG datasets across diverse tasks, electrode configurations, and sequence durations demonstrate that *SAMBA* consistently outperforms state-of-the-art methods while maintaining low memory consumption and inference time. We also show the learned spatial weight maps from our embedding module align closely with task-relevant neurophysiological regions, demonstrating the learnability and interpretability of *SAMBA*. These results highlight *SAMBA*'s scalability and practical potential as a foundation model for real-time brain-computer interface applications. The code is available at: <https://github.com/Jiazhen-Hong/SAMBA>

**Keywords** Electroencephalography (EEG), Long-Sequence Modeling, EEG Self-Supervised Learning, Differential Mamba, EEG Spatial Embedding

## 1 Introduction

Long-sequence electroencephalogram (EEG) modeling is essential for a wide range of real-world applications, including both high sampling rate scenarios, such as steady-state visual evoked potential detection at 30,000 Hz [1]; and long-duration monitoring scenarios, such as Alzheimer's disease detection using 12-second recordings, sleep stage classification from 20-second windows [2], emotion recognition using 2-minute EEG signals [3], and driver fatigue monitor lasting up to two hours [4]. Despite their importance, robust and generalizable EEG modeling remains highly challenging due to the following factors:

(1) *High memory usage in long-sequence modeling.* Transformer-based self-supervised learning (SSL) approaches such as MAEEG [5], BENDR [6], BIOT [7], EEG2Rep [8], LabraM [9], and EEGPT [10] have achieved success in short-term EEG scenarios by leveraging downsampling as a standard preprocessing technique. However, in long-term EEG monitoring tasks, even modest downsampling to 128 Hz yields extremely long input sequences (e.g., 12,800

---

\*Work done during internship at Emotiv Research.

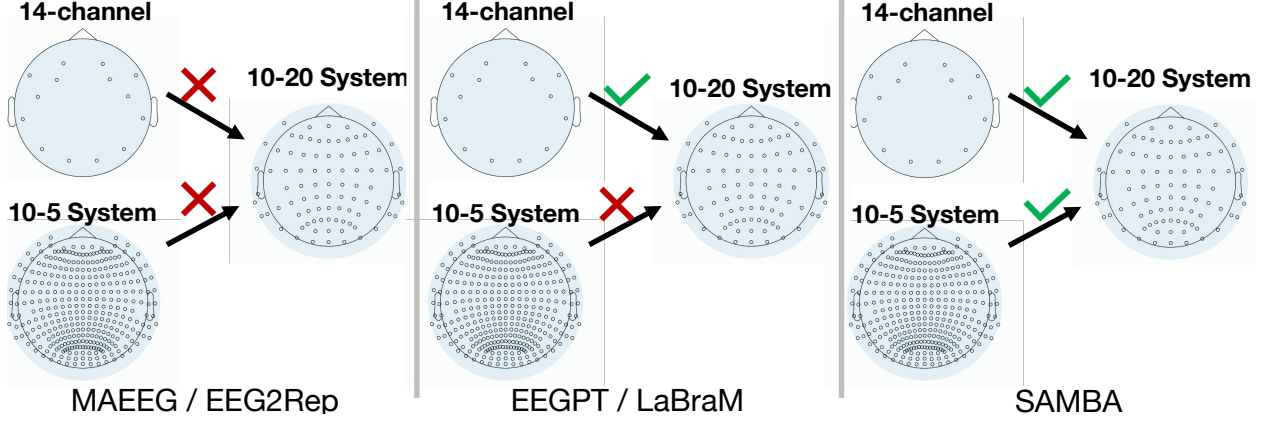


Figure 1: Spatial embedding compatibility of SAMBA and prior EEG models across heterogeneous electrode layouts.

time steps for a 100-second recording). Despite architectural variants aimed at capturing long-range dependencies, Transformers inherently struggle with such sequences due to their  $\mathcal{O}(n^2)$  complexity [11]. This not only limits the feasibility of long-sequence modeling, but also leads to excessive memory usage and slow inference. This motivates an efficient sequence modeling framework based on Structured State Spaces Models (SSMs) [11] to address this bottleneck.

(2) *Diverse EEG headsets and montages.* EEG signals are typically represented as a matrix  $\mathbf{X} \in \mathbb{R}^{C \times T}$ , where  $C$  denotes the number of channels (electrodes) and  $T$  the number of time points. In practice, both vary widely across datasets due to differences in equipment, montage configurations, and application-specific requirements. For example, P300 speller brain-computer interfaces (BCIs) may use only 8 or 16 channels [12], whereas EEG source localization analysis requires high-density 128-channel setups [13]. Even when channel names are consistent, their spatial coordinates  $(x, y, z)$  may differ between devices or cap systems. This heterogeneity limits model generalization and cross-domain transferability. For instance, EEG2Rep [8] and MAEEG [5] assume consistent channel configurations between pretraining and downstream tasks and lack mechanisms to address electrode mismatches. EEGPT [10] and LaBraM [9] use a lookup-based spatial embedding (i.e., channel name-based codebooks), but require all input channels to be part of a predefined electrode set (58 for EEGPT; 128 in the 10–20 system for LaBraM) and face challenges in handling unseen electrodes. These limitations are also visualized in Fig. 1.

(3) *Empirical fixed-length EEG segmentation.* Existing methods often constrain temporal modeling by applying empirical fixed-length window or patch segmentation. For instance, EEGPT [10] and LaBraM [9] adopt local spatio-temporal patching strategies with varied channel configurations, but both use a fixed one-second non-overlapping time window as the patch unit. These designs may overlook temporal semantics between patches and are not compatible with data segments shorter than one second. In time-locked event-related potential (ERP) paradigms, for example, the P300 component ( $\sim 300ms$  long) reflects attention-related target detection [12], while N170 ( $\sim 170ms$  long) marks early face perception over occipito-temporal regions [14]. Such components are highly time-sensitive and often span sub-second durations, making them vulnerable to distortion or loss under fixed-length segmentation.

(4) *Learning robust representations under low SNR and subject dependence.* EEG signals are characterized by a low signal-to-noise ratio (SNR) and non-stationary properties [15], making robust representation learning particularly challenging. Despite the emergence of various deep learning paradigms for EEG representation learning, many researchers still prefer to design handcrafted features for subject-specific models [15, 16, 17]. In contrast, BIOT [7] adopts a unified model across subjects but relies on spectral-domain transformations (e.g., Fast Fourier Transform), which may compromise temporal information.

To address these challenges, we introduce **SAMBA**, a U-shaped SSL framework based on Mamba, which integrates: *Temporal Semantic Random (TSR) Masking*, *Spatial-Adaptive Input Embedding (SAIE)*, and *Multi-Head Differential Mamba (MDM)* modules. The main contributions are summarized as follows.

- Memory-Efficient Long-Sequence Modeling:** SAMBA leverages CNN and Mamba2 structured SSMs[18] layers with linear time complexity, instead of the quadratic complexity inherent in attention mechanisms. This enables the model to handle long EEG recordings (e.g., 100 seconds, 12,800 time steps) in tasks such as abnormality detection without excessive memory usage.

- Spatial Compatibility with various EEG Montages:** SAMBA employs a coordinate-based embedding (SAIE) that relies only on electrode  $(x, y, z)$  positions, making it independent of electrode names or montages. This design allows pretrained SAMBA to generalize to any downstream EEG montage, including unseen electrode layouts. For example,

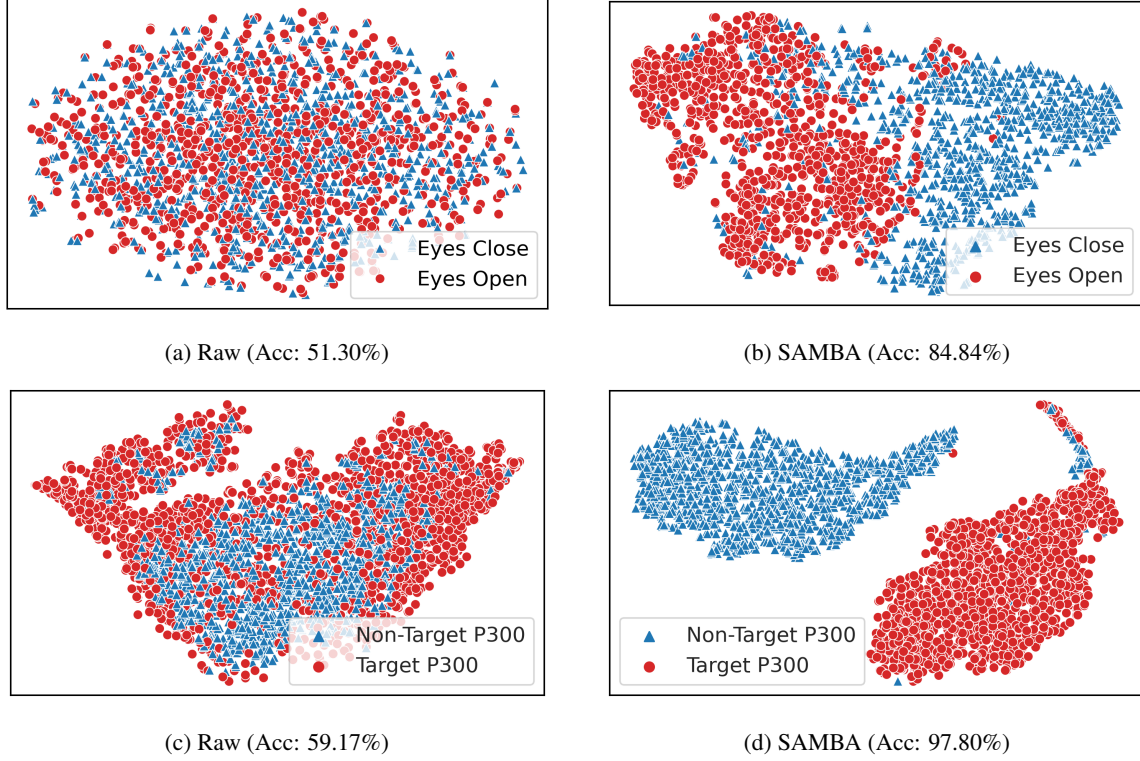


Figure 2: T-SNE plots from Crowdsourced (a-b) and P300 (c-d) datasets, comparing the distribution of raw EEG (a, c) and representations learned by SAMBA (b, d).

when SAMBA is pretrained on the 10–20 system, it can be directly applied to other configurations such as 14 channels or 10-5 systems with over 300 electrodes (Fig. 1).

•**Temporal Compatibility with various EEG durations:** SAMBA adopts a hierarchical encoder-decoder with integrated Mamba2 layers, enabling the model to capture both local and global temporal dependencies via fast and efficient inference. All components are designed to be position-agnostic and invariant to input sequence length, allowing SAMBA to scale seamlessly from short (e.g., 256 steps) to very long (e.g., 12,800 steps) EEG sequences without any architectural modifications while maintaining robust performance.

•**Effective EEG Representation Learning:** SAMBA employs TSR masking to capture semantic structures across diverse time segments and uses the MDM module to contrast dual Mamba2 pathways per head to suppress redundancy and enhance salient temporal structures. A combined temporal-spectral loss further ensures the preservation of both time and frequency domain information.

Fig. 2 visualizes EEG representations learned by SAMBA across montage and duration variations. Pretrained on the TUAB dataset (16 channels, 100 seconds, 12,800 steps at 128 Hz), SAMBA maintains strong representation quality and task performance when transferred to: (i) the Crowdsourced dataset (Appendix C.6), which contains 14 channels and 2-second sequences with 256 steps at 128 Hz (Fig. 2(b)); and (ii) the P300 ERP dataset (Appendix C.4), which includes 64 channels and 0.8-second sequences with 192 steps at 240 Hz (Fig. 2(d)). The clearly separable class clusters highlight the spatial and temporal compatibility of SAMBA.

## 2 Method

**Architecture Overview:** Fig. 3 illustrates the overall architecture of SAMBA. Given an input EEG sequence  $X \in \mathbb{R}^{B \times C_{in} \times T}$ , where  $B$  is the batch size,  $C_{in}$  the number of input channels, and  $T$  the number of time steps, the sequence is first processed by the TSR masking strategy with a 50% time mask and 0% channel mask. The masked input is then projected into a target montage (using the standard 10–20 system in this work) based on  $(x, y, z)$  coordinates of the electrodes using the 3D Spatial-Adaptive Input Embedding (SAIE), followed by a multi-branch Temporal-Receptive embedding to extract short-, mid-, and long-range temporal features at multiple resolutions.

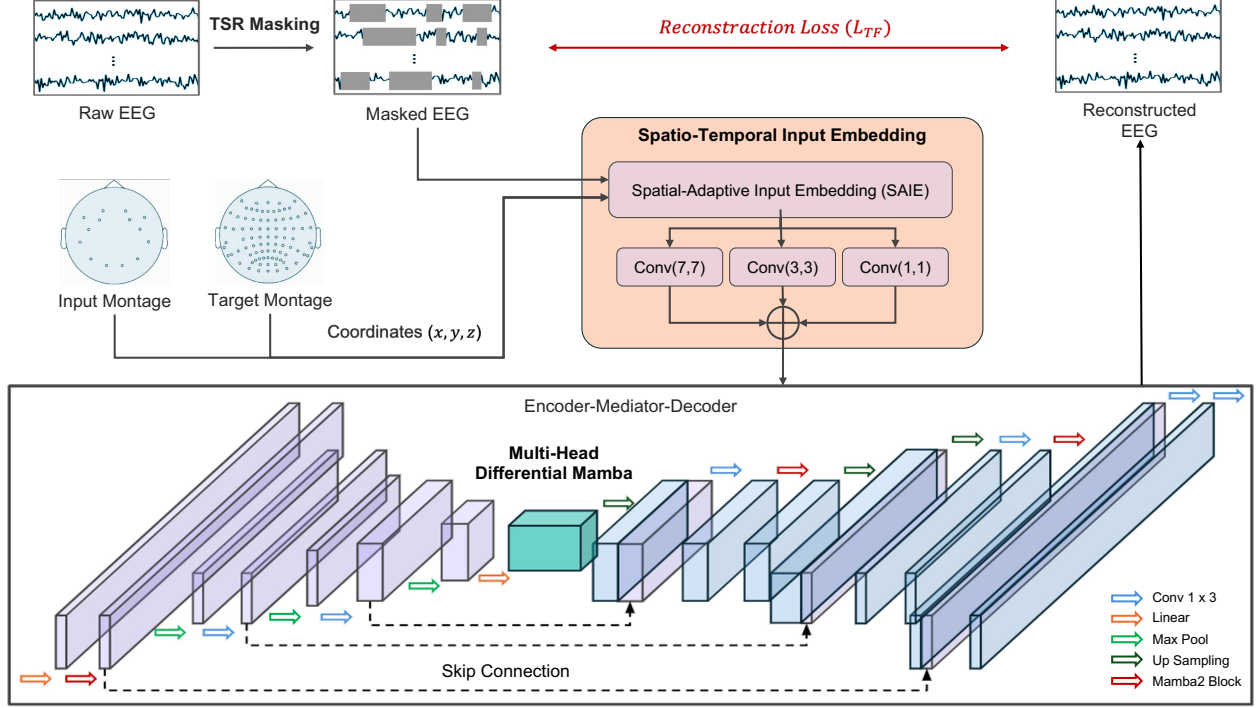


Figure 3: SAMBA Architecture.

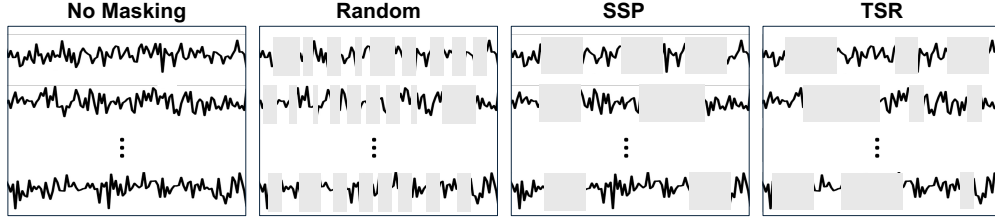


Figure 4: Comparison of proposed TSR masking with existing strategies.

The encoder consists of three stages of increasing depth and decreasing temporal resolution. The first stage applies a linear projection followed by a Mamba2 block [18] to capture fine-grained temporal dynamics. The second and third stages utilize 1D convolutions and max pooling to downsample the temporal dimension while expanding the feature dimension, enabling subsequent Mamba blocks to model longer-range dependencies over compact feature representations. This hierarchical design supports sequences of varying length without requiring fixed-size temporal windows or architectural modifications.

At the bottleneck, a Multi-head Differential Mamba (MDM) module is introduced as a mediator to contrast parallel state-space dynamics across multiple heads and suppress noise.

The decoder mirrors the encoder with symmetric structure and performs upsampling via parameter-free linear interpolation instead of transposed convolutions. This avoids checkerboard artifacts [19] and better preserves the continuity of EEG signals. Each upsampled feature map is refined by a Mamba2 block to restore long-range temporal patterns. To preserve both time-domain waveform and frequency-domain structure, we adopt a Time-Frequency loss function ( $\mathcal{L}_{TF}$ ), which combines an  $L_1$  loss in the time domain with a spectral alignment loss based on the discrete Fourier transform. Further details are provided in Appendix B. This joint objective encourages SAMBA to maintain fidelity in both temporal shape and frequency content, supporting robust EEG modeling.

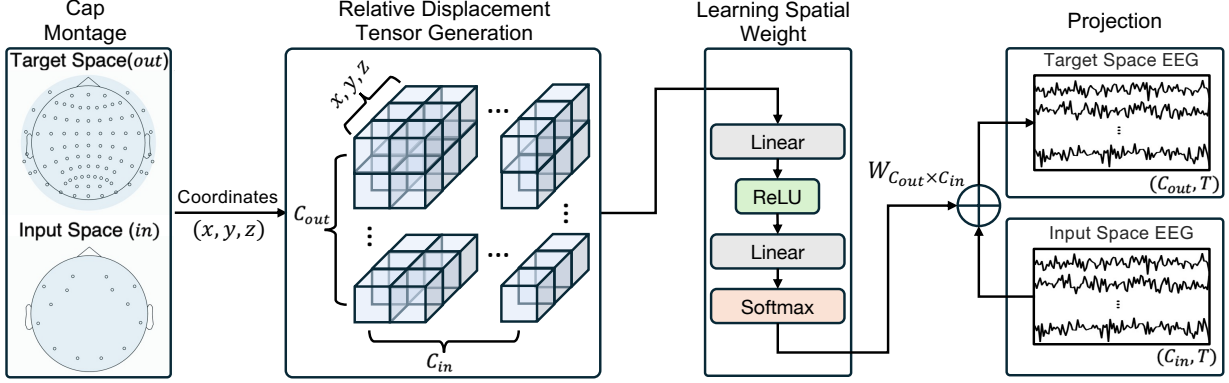


Figure 5: SAIE projects EEG from input to target space using spatial weights derived from relative 3D coordinates.

## 2.1 Temporal Semantic Random (TSR) Masking

Fig. 4 compares different masking strategies used in EEG self-supervised learning. Standard random masking [5] offers high variability but often disrupts temporal continuity. SSP [8] attempt to retain semantic coherence by preserving contiguous segments, but overlapping blocks in SSP introduce redundancy and repetitive patterns, limiting temporal diversity, especially in long EEG sequences.

To address these issues, we propose Temporal Semantic Random (TSR) masking. TSR preserves a fixed number of non-overlapping blocks with variable lengths, sampled from a scaled uniform distribution. This design promotes semantic diversity without redundancy while preserving a constant number of visible time steps.

Let  $l$  be the sequence length and  $\rho$  the masking ratio. TSR preserves  $(1-\rho) \cdot l$  time steps, divided into  $\beta$  non-overlapping blocks. Each block length is sampled as:

$$\forall i \in [1, \beta - 1], \quad \text{Block}_i \sim \mathcal{U} \left( \left[ \alpha_{\min} \cdot \frac{(1-\rho) \cdot l}{\beta}, \left[ \alpha_{\max} \cdot \frac{(1-\rho) \cdot l}{\beta} \right] \right) \quad (1)$$

$$\text{Block}_\beta = (1-\rho) \cdot l - \sum_{i=1}^{\beta-1} \text{Block}_i \quad (2)$$

All blocks are strictly disjoint to avoid overlap and ensure uniform temporal coverage. Compared to prior methods, TSR better balances semantic continuity and temporal diversity, benefiting modeling long-range EEG dependencies.

## 2.2 3D Spatial-Adaptive Input Embedding (SAIE)

Electrode layouts vary across EEG datasets due to differences in headsets, montages, and clinical objectives. To address this spatial variability, we introduce 3D Spatial-Adaptive Input Embedding (SAIE), which uses 3D electrode coordinates to align input signals into a standard brain montage. Figure 5 illustrates SAIE and shows an example that adopting the standard 10–20 system montage [20] as the target space and Emotiv’s 14-channel montage [21] as an input space. Given input EEG  $\mathbf{X} \in \mathbb{R}^{B \times C_{\text{in}} \times T}$ , let  $\mathbf{P}_{\text{in}} \in \mathbb{R}^{C_{\text{in}} \times 3}$  and  $\mathbf{P}_{\text{out}} \in \mathbb{R}^{C_{\text{out}} \times 3}$  denote the 3D coordinates of input and target electrodes, respectively. The relative displacement between each pair is:  $\Delta \mathbf{P}_{ij} = \mathbf{P}_{\text{out},i} - \mathbf{P}_{\text{in},j}$ . The unnormalized spatial weights are generated by a multilayer perceptron (MLP):

$$w_{ij} = \phi_\theta(\Delta \mathbf{P}_{ij}) = \mathbf{W}_2 \cdot \text{ReLU}(\mathbf{W}_1 \Delta \mathbf{P}_{ij} + \mathbf{b}_1) + \mathbf{b}_2, \quad (3)$$

followed by softmax normalization across input channels:

$$\tilde{w}_{ij} = \frac{\exp(w_{ij})}{\sum_{k=1}^{C_{\text{in}}} \exp(w_{ik})}. \quad (4)$$

The projected signal at target channel  $i$  is computed as a weighted sum over all input channels:

$$\mathbf{X}'_{b,i,t} = \sum_{j=1}^{C_{\text{in}}} \tilde{w}_{ij} \cdot \mathbf{X}_{b,j,t}, \quad (5)$$



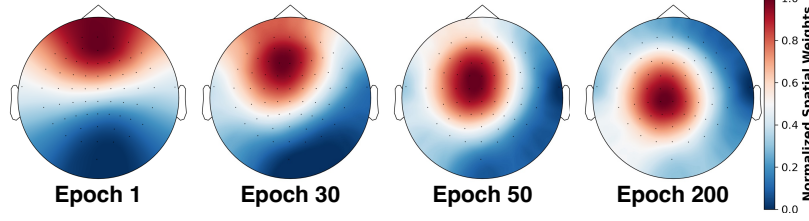


Figure 6: Topographic evolution of spatial weights.

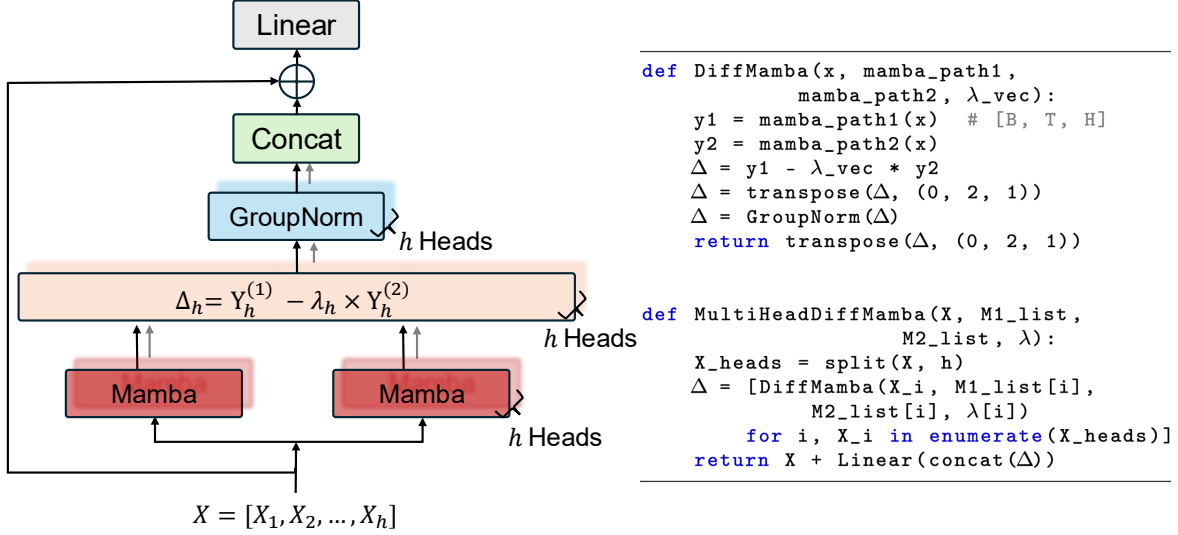


Figure 7: Multi-head Differential Mamba (MDM) block.

where  $b$  indexes the batch,  $t$  the time step,  $j$  the input channel, and  $i$  the target channel. Figure 6 illustrates the topographic evolution of learned spatial weights during pretraining on the PhysionetMI dataset (motor imagery task). Darker red regions indicate higher weights, consistently localized around the motor cortex [15].

### 2.3 Multi-head Differential Mamba (MDM)

Building on the motivation for adopting Mamba2 (Appendix A), we extend its capabilities by proposing the Multi-head Differential Mamba (MDM) module to enhance EEG representation at the encoder–decoder bottleneck. Inspired by differential attention [22], which suppresses noise by contrasting dual attention maps, MDM performs contrastive modeling in the output space of state-space models, here Mamba2, enabling continuous temporal processing without relying on softmax-based attention mechanisms.

Unlike differential attention operating in similarity space ( $QK^\top$ ), each head in MDM processes the input with two independently parameterized Mamba2 modules, enabling diverse dynamic representations. As shown in Figure 7, given input  $X \in \mathbb{R}^{B \times T \times D}$ , where  $B$ ,  $T$ , and  $D$  are batch size, sequence length, and feature dimension respectively, we divide  $X$  along the feature dimension into  $H$  heads:

$$X = [X_1, \dots, X_H], \quad X_h \in \mathbb{R}^{B \times T \times d}, \quad d = D/H. \quad (6)$$

Each head applies two separate Mamba2 modules:

$$Y_h^{(1)} = \text{Mamba}_h^{(1)}(X_h), \quad Y_h^{(2)} = \text{Mamba}_h^{(2)}(X_h). \quad (7)$$

A learnable scaling vector  $\lambda_h \in \mathbb{R}^d$  modulates the difference between the two outputs, followed by per-head Group-Norm:

$$\tilde{\Delta}_h = \text{GroupNorm}_h(Y_h^{(1)} - \lambda_h \cdot Y_h^{(2)}). \quad (8)$$

Table 1: Summary of EEG datasets used for evaluation.

Dataset Source	Task	Dataset Name	# Channels (Hz)	Time Steps	# Samples
Emotiv [8]	Emotion	DREAMER	14 (128)	256	17,246
	Eye State	Alpha	14 (128)	256	11,866
	Eye State	Crowdsourced	14 (128)	256	12,296
	Mental Workload	STEW	14 (128)	256	28,512
	Distraction	DriverDistraction	14 (128)	256	66,197
	Attention	Attention	14 (128)	256	21,894
TUAB [23, 7]	Abnormal Detection	TUAB-10s	16 (128)	1280	409,455
		TUAB-30s	16 (128)	3840	135,702
		TUAB-60s	16 (128)	7680	56,290
		TUAB-100s	16 (128)	12800	39,810
MOABB [24]	Motor Imagery	PhysionetMI	64 (160)	480	18,440
		GrosseWentrup	128 (500)	3500	3,000
		BNCI2014-001	22 (250)	750	2,592
BCIC [25, 26]	P300 ERP	P300-A	64 (240)	192	33,000
		P300-B	64 (240)	192	33,000
		P300-C	64 (240)	192	13,140

The outputs from all heads are concatenated, linearly projected, and added to the original input via a residual connection:

$$Y = X + \text{Linear}(\text{Concat}(\tilde{\Delta}_1, \dots, \tilde{\Delta}_H)). \quad (9)$$

While the original Differential Transformer [22] excludes residuals to isolate contrastive information, we retain the skip connection to preserve meaningful EEG components such as slow trends and rhythmic patterns. This design also improves training stability and convergence.

### 3 Experimental Setup

#### 3.1 Datasets

Table 1 summarizes the thirteen EEG datasets spanning eight cognitive tasks used in this study. Emotiv Research<sup>2</sup> datasets follow EEG2Rep [8] preprocessing. For TUAB [23], we follow BIOT [7] splits and preprocessing for fair comparison. For continuous recordings in TUAB and Crowdsourced EEG [21], we test sequence lengths from 10 to 100 seconds (1280 to 12,800 time steps) to evaluate SAMBA’s long-range modeling. Motor imagery datasets follow the MOABB<sup>3</sup> [24] pipeline. All datasets, except the three single-subject P300 ERP sets [16], use subject-wise splits to assess generalization to unseen subjects. Further details in Appendix C.

#### 3.2 Implementation

SAMBA is trained in two stages: (1) *Pretraining*. SAMBA is pretrained to reconstruct masked EEG sequences using TSR masking (50% time, 0% channel) without labels. Training runs for 200 epochs with a batch size of 64 for in-domain evaluation (Section 4.1), 4096 for Emotiv’s foundation training (Section 4.3) using the AdamW optimizer (weight decay:  $1 \times 10^{-2}$ ). A OneCycle learning rate schedule [27] is applied (max LR:  $5 \times 10^{-4}$ , initial LR:  $2.5 \times 10^{-4}$ , final LR:  $5 \times 10^{-6}$ ) with 10% warm-up followed by cosine annealing strategy. (2) *Downstream Tasks*. SAMBA is evaluated on supervised classification using two modes: **Linear probing**, where a logistic regression classifier is trained on representations extracted from the frozen decoder outputs right after the MDM module; and **Fine-tuning**, where the full model is initialized from the pretrained weights and jointly fine-tuned with an MLP head for a few epochs ( $< 5$ ). All experiments are conducted with 32-bit mixed precision on two NVIDIA RTX 6000 Ada GPUs. Further details in Appendix D.

<sup>2</sup>[www.emotiv.com/research](http://www.emotiv.com/research)

<sup>3</sup><https://moabb.neurotechx.com>

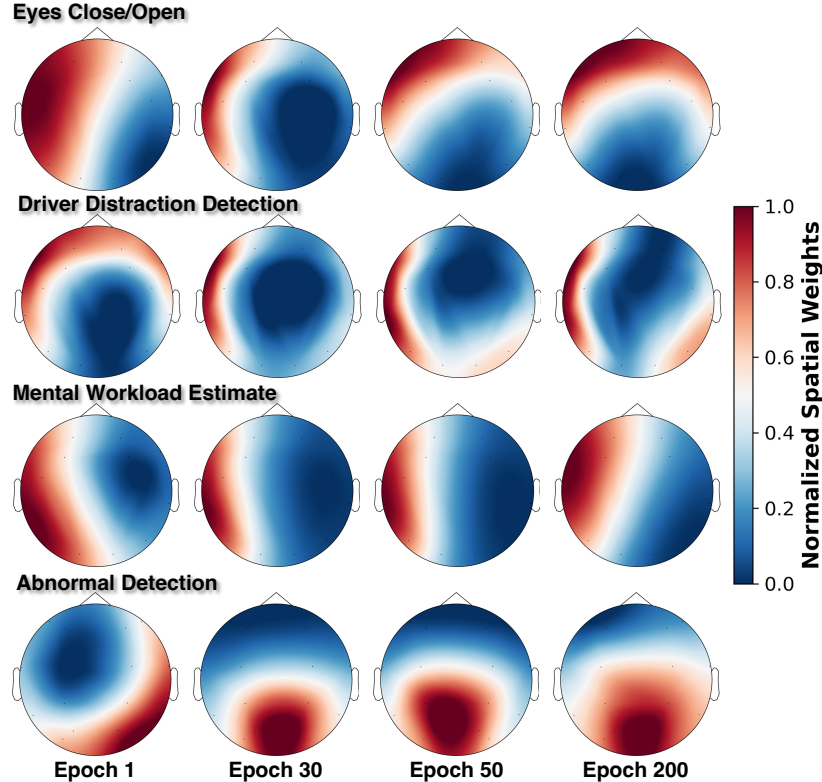


Figure 8: Evolution of learned spatial weights (input embeddings) over epochs from in-domain pretraining. Darker red regions indicate higher spatial weights.

## 4 Results

We evaluated SAMBA as a foundation model for EEG from four perspectives: (1) *Learnability* through in-domain performance, (2) *Transferability* across datasets and domains, (3) *Scalability* toward foundation models and (4) *Ablation Study*. Evaluation metrics are detailed in Appendix F.

### 4.1 In-Domain

**Spatial Learnability** Figure 8 shows the evolution of input embedding weights over training epochs across four EEG tasks. Despite differences in electrode layouts and sequence lengths, the spatial patterns converge toward physiologically relevant regions. At epoch 1, the weights are nearly random; by epoch 50–200, they become more stable and task-specific. For Crowdsourced (eyes open/closed), the model emphasizes the frontal lobe region, consistent with alpha modulation during eye closure [28]. For DriverDistraction, the weights highlight the left temporal lobe, which is associated with auditory processing, working memory, and cognitive functions essential for managing distractions during driving [29]. In the STEW dataset (mental workload), the model highlights the left frontal area, consistent with [30], which reports increased left frontal activity under acute stress. For TUAB (abnormal detection), occipital dominance aligns with [23], which shows occipital electrodes yield better performance for distinguishing abnormal EEG.

**Representation Learnability** To compare with prior work, we evaluate SAMBA’s in-domain performance on the same datasets and tasks in existing benchmarks. For Emotiv datasets (Crowdsourced, DriverDistraction, and STEW), we compare with the reported results in EEG2Rep [8]. For TUAB, we benchmark with the numbers reported in EEGPT [10] and LabraM [9].

Table 2 reports the performance of SAMBA on the Emotiv datasets under both linear probing and fine-tuning settings. For reference, results from training SAMBA from scratch (*Random*, without pretraining) are also included. The best scores for each setting are highlighted in bold. Except for BIOT, which serves as a general-purpose foundation model, all compared methods are specifically designed to enhance EEG representations. SAMBA consistently achieves the



Table 2: Performance on the Emotiv datasets.

Models	Crowdsourced		DriverDistraction		STEW	
	ACC	AUROC	ACC	AUROC	ACC	AUROC
BENDR [6]	70.46±4.14	0.7056±0.04	68.40±3.08	0.5521±0.03	63.03±1.07	0.6303±0.01
MAEEG [5]	75.21±2.11	0.7501±0.02	68.37±2.60	0.5529±0.02	67.99±1.86	0.6858±0.02
BIOT [7]	76.23±4.56	0.7633±0.04	63.93±1.28	0.6333±0.01	67.54±2.08	0.6769±0.04
EEG2Rep [8]	81.66±2.93	0.8167±0.03	76.88±2.55	0.6559±0.02	69.04±1.04	0.6910±0.01
SAMBA (Linear)	<b>86.72±0.02</b>	<b>0.9336±0.00</b>	<b>77.97±0.90</b>	<b>0.7092±0.01</b>	<b>70.62±0.03</b>	<b>0.7704±0.00</b>
BENDR [6]	83.78±2.35	0.8380±0.03	74.31±2.38	0.5986±0.03	69.74±2.11	0.6977±0.02
MAEEG [5]	86.75±3.50	0.8621±0.03	74.58±2.16	0.6079±0.03	72.46±3.67	0.7250±0.03
BIOT [7]	87.95±3.52	0.8778±0.03	74.34±3.57	0.6121±0.04	69.88±2.15	0.7011±0.03
EEG2Rep [8]	<b>94.13±2.11</b>	0.9413±0.02	80.07±2.63	0.6614±0.02	<b>73.60±1.47</b>	0.7440±0.02
SAMBA (Random)	69.66±4.70	0.7943±0.02	63.38±2.33	0.6377±0.02	57.21±1.30	0.6033±0.02
SAMBA (fine-tuned)	93.24±1.42	<b>0.9793±0.01</b>	<b>80.17±0.68</b>	<b>0.6845±0.01</b>	70.89±0.84	<b>0.7862±0.01</b>

Table 3: Performance on the TUAB datasets.

Models	Model Size	Balanced ACC	AUROC
SPaRCNet [31]	0.79M	78.96±0.18	0.8676±0.0012
ContraWR [32]	1.6M	77.46±0.41	0.8456±0.0074
CNN-LSTM [33]	2.4M	78.48±0.38	0.8569±0.0051
CNN-Transformer [34]	3.2M	77.77±0.22	0.8461±0.0013
BIOT [7]	3.2M	79.59±0.57	0.8815±0.0043
ST-Transformer [35]	3.5M	79.66±0.23	0.8707±0.0019
EEGPT [10]	25M	79.83±0.30	0.8718±0.0050
LaBraM-Base [9]	5.8M	81.40±0.16	0.9022±0.0009
SAMBA-10s	1.0 M	81.50±0.37	0.8887±0.0004
SAMBA-30s	1.0 M	81.88±0.09	0.8937±0.0025
SAMBA-60s	1.0 M	82.18±0.12	0.8962±0.0019
SAMBA-100s	1.0 M	<b>82.64±0.17</b>	<b>0.9054±0.0044</b>

highest AUROC across all tasks and settings, demonstrating superior robustness and discriminative capacity. While EEG2Rep attains competitive accuracy in some cases, its lower AUROC indicates a less stable decision boundary.

Table 3 summarizes the fine-tuning performance of SAMBA on TUAB datasets, along with model sizes. Leveraging the continuous nature of TUAB recordings [23], we pretrain SAMBA on four different sequence lengths: 10 s, 30 s, 60 s, and 100 s, denoted as SAMBA-10s through SAMBA-100s. EEGPT [10] and LabraM [9] report multiple model variants of different sizes, but only those with available checkpoints on GitHub are provided in this table. All reported results, except for SAMBA, use the default 10-second TUAB segments. LabraM-Base achieves the highest AUROC (0.9022), slightly outperforming SAMBA-10s (0.8887). Despite having only 1.0M parameters, the second smallest among all models, SAMBA-10s delivers competitive results. As the pretraining sequence length increases, SAMBA demonstrates consistent performance improvements. When pretrained on 100-second sequences, SAMBA-100s achieves the best overall performance, with a balanced accuracy of 82.64% and an AUROC of 0.9054, demonstrating its strong representation learnability for long-sequence EEG modeling.

## 4.2 Cross-Domain

We evaluate SAMBA’s across domain ability by pretraining on one dataset and transferring to downstream tasks on others. Two aspects of transferability are explored in this section: *temporal*, where sequence length mismatch, and *spatial*, where electrode number and configurations vary.

**Temporal Transferability** Enabled by SAMBA’s design for temporal transferability, we can assess how pretraining on one sequence length contributes to performance on downstream tasks with different durations. As both TUAB [23] and Crowdsourced [21] datasets contain continuous EEG recordings, we vary the sequence lengths for evaluation under a fixed sampling rate (128 Hz), as shown in Figure 9. SAMBA-100s (12800 sequence length), achieves the most stable performance across different sequence lengths on the Crowdsourced eye open/closed task, consistently yielding high accuracy and AUROC. Interestingly, from Fig. 9 (b), we can see SAMBA-30s (3840 sequence length) performs best on the 30-second task, and SAMBA-10s (1280 sequence length) performs best on the 10-second task, likely due

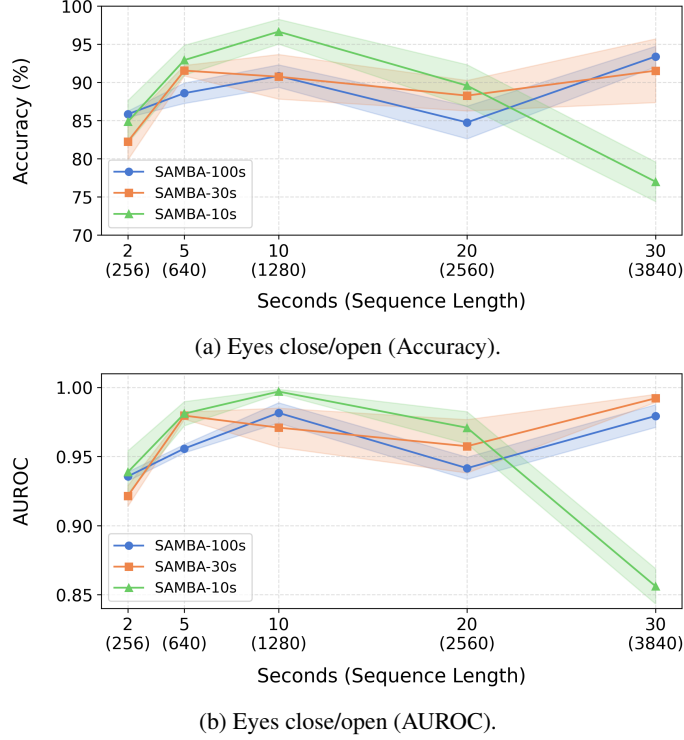


Figure 9: Temporal transferability: SAMBA pretrained on long-context data and evaluated on short-context tasks.

Table 4: Spatial Transferability: SAMBA pretrained on TUAB, tested on Emotiv datasets with a different montage.

Initialization	Crowdsourced		DriverDistraction		STEW	
	ACC	AUROC	ACC	AUROC	ACC	AUROC
Random	69.66	0.7943	63.38	0.6377	57.21	0.6033
In-domain	( $\uparrow$ 23.58)	( $\uparrow$ 0.1850)	( $\uparrow$ 16.80)	( $\uparrow$ 0.0468)	( $\uparrow$ 13.68)	( $\uparrow$ 0.1829)
Cross-domain	( $\uparrow$ 16.17)	( $\uparrow$ 0.1414)	( $\uparrow$ 10.66)	( $\uparrow$ 0.0313)	( $\uparrow$ 16.39)	( $\uparrow$ 0.1951)

to alignment between the temporal scale of pretraining and evaluation. However, as shown in Fig. 9 (a), SAMBA-30s shows higher standard deviation, indicating reduced stability. Overall, SAMBA-100s shows strong generalization to shorter sequences, indicating that pretraining on long sequences can effectively transfer to tasks with limited temporal context. In contrast, models pretrained on short sequences exhibit poor generalization to longer tasks (see Fig. 9, green line at 30s). These results highlight the benefit of long-sequence pretraining for building temporally transferable EEG representations.

**Spatial transferability** SAMBA’s spatial transferability is supported by the SAIE module, enabling transferability across different electrode montages. Table 4 reports performance on three Emotiv datasets (14 channels, 2-second inputs) using SAMBA-100s pretrained on TUAB (16 channels, 100-second sequences). Both in-domain and cross-domain pretraining outperform random initialization. In-domain pretraining achieves the highest gains due to task and distribution alignment. Notably, cross-domain pretraining also delivers consistent improvements across all tasks, despite differences in devices, montages, and objectives, highlighting the spatial and temporal transferability of SAMBA’s learned representations.

### 4.3 SAMBA as a Foundation Model

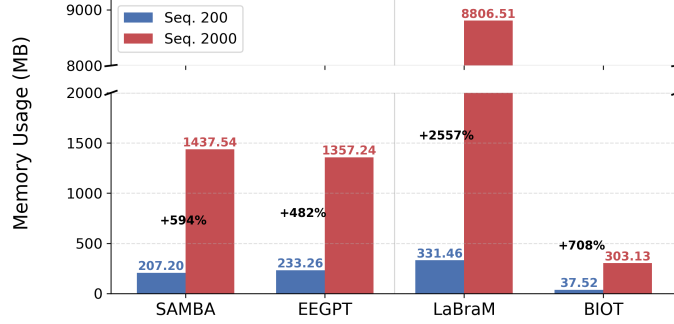
Table 5 compares SAMBA with existing foundation models across six datasets with varying montages, sampling rates, and durations. EEGPT and LaBraM offer only one checkpoint each (Large <sup>4</sup> and Base <sup>5</sup>). BIOT provides all

<sup>4</sup><https://github.com/BINE022/EEGPT>, accessed: July 01, 2025

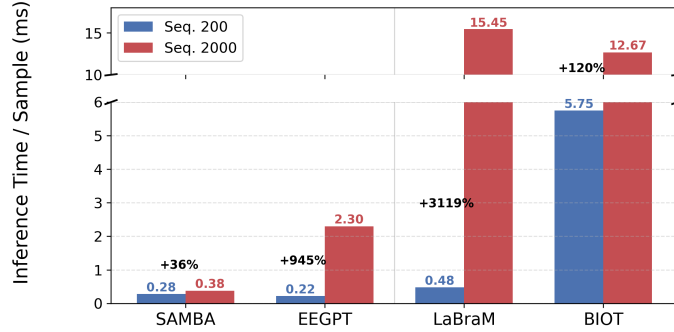
<sup>5</sup><https://github.com/935963004/LaBraM>, accessed: July 01, 2025

Table 5: SAMBA vs. existing foundation models - accuracy.

Dataset	# of Class	# of Chan.	Sampling Rate (Hz)	Time (s)	SAMBA-E		SAMBA-T		EEGPT		LaBraM		BIOT	
					Acc.	W-F1	Acc.	W-F1	Acc.	W-F1	Acc.	W-F1	Acc.	W-F1
PhysionetMI	5	64	160	3	28.29	0.3070	27.41	0.3006	23.77	0.2631	<b>35.79</b>	<b>0.3575</b>	27.14	0.2832
GrosseWentrup	2	128	500	7	53.50	0.5349	<b>57.83</b>	<b>0.5781</b>	50.67	0.5037	53.67	0.5344	50.33	0.4530
BNCI2014-001	4	22	250	4	<b>32.00</b>	<b>0.3220</b>	30.57	0.3033	25.41	0.2500	28.53	0.2807	27.31	0.2388
P300-A	2	64	240	0.8	56.26	0.6181	<b>58.49</b>	<b>0.6366</b>	55.79	0.6129	Incompatible		Incompatible	
P300-B	2	64	240	0.8	57.76	0.6303	<b>59.22</b>	<b>0.6424</b>	54.00	0.5986	Incompatible		Incompatible	
P300-C	2	64	240	0.8	<b>61.68</b>	<b>0.6615</b>	61.28	0.6593	57.36	0.6262	Incompatible		Incompatible	



(a) Memory Usage (MB).



(b) Inference Speed (ms).

Figure 10: SAMBA vs. existing foundation models - efficiency.

checkpoints<sup>6</sup> and we use the largest (PRESET+SHHS, 18 channels). SAMBA-E is pretrained on all Emotiv datasets (311k samples, 6 tasks, 170 hours) with 2 seconds (sequence length 256) duration. SAMBA-T is pretrained on TUAB-100s (sequence length 12,800) that span over 1,000 hours of EEG recordings.

All results are from linear probing on released checkpoints. Bold indicates the best model per dataset. Notably, LaBraM achieves the best accuracy on PhysionetMI, likely due to its inclusion in pretraining. However, on all other unseen datasets, SAMBA-E and SAMBA-T consistently outperform the others. On three downstream EEG datasets, SAMBA-T shows the highest performance, likely due to its pretraining to long-duration EEG recordings in TUAB, as well as its quantile-based representation strategy (see Appendix D), which enables efficient capture of long-context representations for probing. LaBraM uses a patch size of 200 samples (i.e., 1 s at 200 Hz), and BIOT performs short-time Fourier transform (STFT) with 200-point windows and a hop length of 100. To ensure compatibility, all input sequences are resampled to 200 Hz for both LaBraM and BIOT. Nevertheless, these two models remain incompatible with short-duration datasets (e.g., 800 ms P300 datasets) whose duration is shorter than the required patch or window size. Furthermore, the EEGPT and BIOT checkpoints are only compatible with input channel numbers less than 58 and 18, respectively. To address this, we adopt the channel mapping strategy shown in EEGPT GitHub<sup>1</sup> (via “Conv1dWithConstraint”) to match the required input dimensions. In contrast, SAMBA does not require channel remapping or sequence truncation, its checkpoint can generalize to arbitrary channel numbers, durations, and sampling rates, supporting SAMBA’s scalability toward an EEG foundation model.

<sup>6</sup><https://github.com/ycq091044/BIOT>, accessed: July 01, 2025

Table 6: Performance of SAMBA and its variants on TUAB-100s (sequence length: 12,800 per trial).

Model Variant	Modified Component	ACMSE	B-ACC (%)	AUROC
SAMBA	—	<b>3.42e-8</b>	<b>82.64</b>	<b>0.9054</b>
Setting 1	TSR $\rightarrow$ random mask	2.17e-5	79.79 ( $\downarrow$ 2.84)	0.8629 ( $\downarrow$ 0.04)
Setting 2	MDM $\rightarrow$ Mamba2 block	1.43e-6	80.01 ( $\downarrow$ 2.63)	0.8687 ( $\downarrow$ 0.04)
Setting 3	Remove residual in MDM	9.06e-8	81.41 ( $\downarrow$ 1.23)	0.8867 ( $\downarrow$ 0.02)
Setting 4	Remove $\mathcal{L}_T$	2.20e-7	80.18 ( $\downarrow$ 2.46)	0.8677 ( $\downarrow$ 0.04)
Setting 5	Remove $\mathcal{L}_F$	1.47e-9	80.54 ( $\downarrow$ 2.10)	0.8827 ( $\downarrow$ 0.02)
Setting 6	All Mamba $\rightarrow$ Conv	5.04e-7	78.67 ( $\downarrow$ 3.97)	0.8552 ( $\downarrow$ 0.05)
Setting 7	Mamba $\rightarrow$ Attention	OOM	OOM	OOM

**Memory Usage & Inference Speed** Figure 10 compares the memory usage and inference speed of the four EEG foundation models. We evaluate both metrics on synthetic EEG sequences with 22 channels and two sequence lengths: 200 and 2000 data points. Each measurement is repeated five times, and the average is reported. Blue bars show the results for short sequences, and red bars for long sequences. The percentage annotations indicate the relative increase when the sequence length increases. Among all models, SAMBA achieves the fastest inference and the second lowest memory growth. BIOT shows the lowest memory usage due to its STFT-based pre-processing, which compresses the input into frequency-domain features. However, this also increases the computational cost, leading to slower inference. LaBraM shows the highest memory usage, especially for long sequences. It divides each 200-point segment of every channel per second into individual tokens, so the number of tokens grows with both channel count and sequence length. Since transformer attention scales quadratically with the number of tokens, memory usage increases rapidly. LaBraM also applies full attention without early-stage compression. In contrast, EEGPT reduces memory usage for long sequences by applying a *temporal interpolation* step to resample inputs to a fixed length, as in their 30 s Sleep-EDF setting<sup>1</sup>. This strategy constrains the number of tokens and avoids attention overhead. SAMBA processes the entire sequence like LaBraM but does not rely on attention. Instead, it uses Mamba2 blocks, which scale more efficiently and support long-sequence modeling with lower memory and faster inference.

#### 4.4 Ablation Study & Long Sequence Modeling

Here, we conduct ablation studies on TUAB-100s under seven configurations: (1) replacing TSR with standard random masking, (2) removing MDM and using a single Mamba2 block, (3) removing the residual connection in MDM, (4–5) removing the time-domain or frequency-domain loss, respectively, (6) replacing all Mamba modules with convolutional layers (similar to U-Net), and (7) replacing Mamba with attention block. Evaluation metrics include Averaged Channel-wise MSE (ACMSE) during pretraining, and Balanced Accuracy and AUROC during fine-tuning (see Appendix F). As shown in Table 6, removing the frequency-domain loss (Setting 5) results in the lowest ACMSE but a large drop in AUROC, indicating the importance of spectral information for stable decision boundaries. Removing the residual connection (Setting 3) leads to a minor performance drop, validating its regularization and stabilizing effect. The convolutional baseline (Setting 6) performs worst across all metrics, while the attention-based variant (Setting 7) fails due to out-of-memory (OOM), confirming the necessity of Mamba and MDM components for efficient long-sequence modeling. Additional ablation studies can be found in Appendix E.

## 5 Conclusion

We present *SAMBA* for long-context EEG modeling. *SAMBA* effectively captures long-range temporal dependencies, handles variability in configuration and between subjects, and learns robust representations through its Mamba-based architecture, 3D spatial-adaptive input embedding, and multi-head differential Mamba module. Extensive evaluations across thirteen EEG datasets demonstrate *SAMBA*’s superior performance in both in-domain and cross-domain settings, outperforming state-of-the-art in accuracy and efficiency. We also note that pretraining on longer EEG recordings yields better generalization to shorter downstream tasks, while the reverse does not hold, highlighting the significance of long-context modeling. Through comprehensive analyses of learnability, transferability, and scalability, including comparison of memory usage and inference speed with existing EEG foundation models, we also demonstrate *SAMBA*’s potential as a foundation model for real-world BCI scenarios with heterogeneous acquisition setups. Future work will scale *SAMBA* with larger corpora and broader task to fully realize its foundation model potentials.

## References

- [1] Ruikai Li, Yixing Zhang, Guangwei Fan, Ziteng Li, Jun Li, Shiyong Fan, Cunguang Lou, and Xiuling Liu. Design and implementation of high sampling rate and multichannel wireless recorder for EEG monitoring and SSVEP response detection. *Frontiers in Neuroscience*, 17:1193950, 2023.
- [2] Ian G Campbell. EEG recording and analysis for sleep research. *Current protocols in neuroscience*, 49(1):10–2, 2009.
- [3] Wei-Long Zheng, Wei Liu, Yifei Lu, Bao-Liang Lu, and Andrzej Cichocki. Emotionmeter: A multimodal framework for recognizing human emotions. *IEEE transactions on cybernetics*, 49(3):1110–1122, 2018.
- [4] Wei-Long Zheng and Bao-Liang Lu. A multimodal approach to estimating vigilance using EEG and forehead EOG. *Journal of neural engineering*, 14(2):026017, 2017.
- [5] Hsiang-Yun Sherry Chien, Hanlin Goh, Christopher M Sandino, and Joseph Y Cheng. MAEEG: Masked auto-encoder for EEG representation learning. *arXiv preprint arXiv:2211.02625*, 2022.
- [6] Demetres Kostas, Stephane Aroca-Ouellette, and Frank Rudzicz. BENDR: Using transformers and a contrastive self-supervised learning task to learn from massive amounts of EEG data. *Frontiers in Human Neuroscience*, 15:653659, 2021.
- [7] Chaoqi Yang, M Westover, and Jimeng Sun. BIOT: Biosignal transformer for cross-data learning in the wild. *Advances in Neural Information Processing Systems*, 36, 2024.
- [8] 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.
- [9] Wei-Bang Jiang, Li-Ming Zhao, and Bao-Liang Lu. Large brain model for learning generic representations with tremendous EEG data in BCI. *arXiv preprint arXiv:2405.18765*, 2024.
- [10] Guangyu Wang, Wenchao Liu, Yuhong He, Cong Xu, Lin Ma, and Haifeng Li. EEGPT: Pretrained transformer for universal and reliable representation of EEG signals. *Advances in Neural Information Processing Systems*, 37:39249–39280, 2024.
- [11] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- [12] Jiazhen Hong and Laleh Najafizadeh. P3T: A transformer model for enhancing character recognition rates in P300 speller systems. In *2024 58th Asilomar Conference on Signals, Systems, and Computers*, pages 514–518, 2024.
- [13] Ali E Haddad et al. Source-informed segmentation: A data-driven approach for the temporal segmentation of EEG. *IEEE Transactions on Biomedical Engineering*, 66(5):1429–1446, 2018.
- [14] Stéphanie Caharel and Bruno Rossion. The N170 is sensitive to long-term (personal) familiarity of a face identity. *Neuroscience*, 458:244–255, 2021.
- [15] Jiazhen Hong, Foroogh Shamsi, and Laleh Najafizadeh. A deep learning framework based on dynamic channel selection for early classification of left and right hand motor imagery tasks. In *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 3550–3553. IEEE, 2022.
- [16] Zehui Wang, Chuanguan Chen, Junhua Li, Feng Wan, Yu Sun, and Hongtao Wang. ST-CapsNet: linking spatial and temporal attention with capsule network for P300 detection improvement. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31:991–1000, 2023.
- [17] Ronghua Ma, Tianyou Yu, Xiaoli Zhong, Zhu Liang Yu, Yuanqing Li, and Zhenghui Gu. Capsule network for ERP detection in brain-computer interface. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 29:718–730, 2021.
- [18] Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- [19] Yusuke Sugawara, Sayaka Shiota, and Hitoshi Kiya. Checkerboard artifacts free convolutional neural networks. *APSIPA Transactions on Signal and Information Processing*, 8:e9, 2019.
- [20] Alexandre Gramfort, Martin Luessi, Eric Larson, Denis A Engemann, Daniel Strohmeier, Christian Brodbeck, Roman Goj, Mainak Jas, Teon Brooks, Lauri Parkkonen, et al. MEG and EEG data analysis with MNE-Python. *Frontiers in Neuroinformatics*, 7:267, 2013.

- [21] Nikolas S Williams, William King, Geoffrey Mackellar, Roshini Randeniya, Alicia McCormick, and Nicholas A Badcock. Crowdsourced EEG experiments: A proof of concept for remote EEG acquisition using EmotivPRO builder and EmotivLABS. *Heliyon*, 9(8), 2023.
- [22] Tianzhu Ye, Li Dong, Yuqing Xia, Yutao Sun, Yi Zhu, Gao Huang, and Furu Wei. Differential transformer. *arXiv preprint arXiv:2410.05258*, 2024.
- [23] Sebas Lopez, G Suarez, D Jungreis, I Obeid, and Joseph Picone. Automated identification of abnormal adult EEGs. In *2015 IEEE signal processing in medicine and biology symposium (SPMB)*, pages 1–5. IEEE, 2015.
- [24] Vinay Jayaram and Alexandre Barachant. MOABB: trustworthy algorithm benchmarking for BCIs. *Journal of neural engineering*, 15(6):066011, 2018.
- [25] Gerwin Schalk and Benjamin Blankertz. 2nd Wadsworth BCI Dataset (P300 Evoked Potentials). <http://www.bci2000.org>, 2004. Data acquired using BCI2000 P3 Speller Paradigm. BCI Competition 2003 dataset.
- [26] Dean J. Krusienski, Benjamin Blankertz, Gerwin Schalk, Jonathan R. Wolpaw, and Klaus-Robert Müller. BCI competition iii dataset ii: P300 speller paradigm. In *Proceedings of the 3rd International BCI Meeting*, 2006. Dataset available at <http://www.bci2000.org>.
- [27] Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*, volume 11006, pages 369–386. SPIE, 2019.
- [28] Sven Hoffmann and Michael Falkenstein. The correction of eye blink artefacts in the EEG: a comparison of two prominent methods. *PloS one*, 3(8):e3004, 2008.
- [29] Guofa Li, Xiaojian Wu, Arno Eichberger, Paul Green, Cristina Olaverri-Monreal, Weiquan Yan, Yechen Qin, and Yuezhi Li. Drivers’ eeg responses to different distraction tasks. *Automotive Innovation*, 6(1):20–31, 2023.
- [30] Gesa Berretz, Julian Packheiser, Oliver T Wolf, and Sebastian Ocklenburg. Acute stress increases left hemispheric activity measured via changes in frontal alpha asymmetries. *Isience*, 25(2), 2022.
- [31] Jin Jing, Wendong Ge, Shenda Hong, Marta Bento Fernandes, Zhen Lin, Chaoqi Yang, Sungtae An, Aaron F Struck, Aline Herlopian, Ioannis Karakis, et al. Development of expert-level classification of seizures and rhythmic and periodic patterns during EEG interpretation. *Neurology*, 100(17):e1750–e1762, 2023.
- [32] Chaoqi Yang, Cao Xiao, M Brandon Westover, and Jimeng Sun. Self-supervised electroencephalogram representation learning for automatic sleep staging: model development and evaluation study. *JMIR AI*, 2(1):e46769, 2023.
- [33] Hongli Li, Man Ding, Ronghua Zhang, and Chunbo Xiu. Motor imagery eeg classification algorithm based on cnn-lstm feature fusion network. *Biomedical signal processing and control*, 72:103342, 2022.
- [34] Wei Yan Peh, Yuanyuan Yao, and Justin Dauwels. Transformer convolutional neural networks for automated artifact detection in scalp eeg. In *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 3599–3602. IEEE, 2022.
- [35] Yonghao Song, Xueyu Jia, Lie Yang, and Longhan Xie. Transformer-based spatial-temporal feature learning for EEG decoding. *arXiv preprint arXiv:2106.11170*, 2021.
- [36] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [37] Barak Lenz, Opher Lieber, Alan Arazi, Amir Bergman, Avshalom Manevich, Barak Peleg, Ben Aviram, Chen Almagor, Clara Fridman, Dan Padnos, et al. Jamba: Hybrid transformer-mamba language models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [38] Simiao Zuo, Xiaodong Liu, Jian Jiao, Denis Charles, Eren Manavoglu, Tuo Zhao, and Jianfeng Gao. Efficient long sequence modeling via state space augmented transformer. *arXiv preprint arXiv:2212.08136*, 2022.
- [39] Paolo Gloriosio, Quentin Anthony, Yury Tokpanov, James Whittington, Jonathan Pilault, Adam Ibrahim, and Beren Millidge. Zamba: A compact 7b ssm hybrid model. *arXiv preprint arXiv:2405.16712*, 2024.
- [40] Sinziana Mazilu and José Iria. L1 vs. l2 regularization in text classification when learning from labeled features. In *2011 10th international conference on machine learning and applications and workshops*, volume 1, pages 166–171. IEEE, 2011.
- [41] Clemens Brunner, Robert Leeb, Gernot Müller-Putz, Alois Schlögl, and Gert Pfurtscheller. BCI competition 2008–graz data set a. *Institute for knowledge discovery (laboratory of brain-computer interfaces)*, Graz University of Technology, 16(1-6):34, 2008.



- [42] Moritz Grosse-Wentrup, Christian Liefhold, Klaus Gramann, and Martin Buss. Beamforming in noninvasive brain–computer interfaces. *IEEE Transactions on Biomedical Engineering*, 56(4):1209–1219, 2009.
- [43] Gerwin Schalk, Dennis J McFarland, Thilo Hinterberger, Niels Birbaumer, and Jonathan R Wolpaw. BCI2000: a general-purpose brain-computer interface (BCI) system. *IEEE Transactions on biomedical engineering*, 51(6):1034–1043, 2004.
- [44] Jonathan R. Wolpaw, Gerwin Schalk, and Dean Krusienski. BCI competition III, data set II: P300 speller paradigm. Technical report, Wadsworth Center, NYS Department of Health, 2004.
- [45] Matthias Kaper, Peter Meinicke, Ulf Grossekhoefer, Thomas Lingner, and Helge Ritter. BCI competition 2003-data set iib: support vector machines for the P300 speller paradigm. *IEEE Transactions on biomedical Engineering*, 51(6):1073–1076, 2004.
- [46] Wei Lun Lim, Olga Sourina, and Lipo P Wang. STEW: Simultaneous task EEG workload data set. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 26(11):2106–2114, 2018.
- [47] Stamos Katsigiannis and Naeem Ramzan. DREAMER: A database for emotion recognition through EEG and ECG signals from wireless low-cost off-the-shelf devices. *IEEE journal of biomedical and health informatics*, 22(1):98–107, 2017.

## A State Space Models (SSMs) and Mamba

Transformer-based models exhibit  $\mathcal{O}(n^2)$  complexity in sequence length, making them inefficient for long EEG signals. State Space Models (SSMs), particularly the Mamba series [36, 18], provide a scalable alternative by enabling linear-time computation with constant memory usage. These models evolve a latent hidden state using a recurrence of the form:

$$\mathbf{h}_t = A_t \mathbf{h}_t - 1 + B_t \mathbf{x}_t, \quad \mathbf{y}_t = C_t \mathbf{h}_t, \quad (10)$$

where  $\mathbf{x}_t$  is the input at time  $t$ ,  $\mathbf{h}_t$  is the internal state, and  $A_t, B_t, C_t$  are learnable matrices. This formulation enables SSMs to process sequences efficiently while maintaining the ability to capture long-range dependencies.

Mamba2 [18] introduces a refined variant of SSMs through the Structured State Space Duality (SSD) framework. This dual formulation expresses SSMs as structured matrix multiplications, specifically, as semiseparable matrices, which admit both a linear recurrent form and a quadratic attention-like form. The structured matrix  $M \in \mathbb{R}^{T \times T}$  corresponding to an SSM transformation is given by:

$$M_{tj} = C_t^\top A_t \cdots A_{j+1} B_j, \quad \mathbf{y} = M \mathbf{x}, \quad (11)$$

revealing that SSMs share a mathematical foundation with attention mechanisms, but with improved efficiency.

To reduce hardware inefficiency in early SSM designs, Mamba2 incorporates an SSD algorithm that leverages structured matrix representations and grouped parallel projections. This design enables better utilization of matrix multiplication units on modern accelerators, making it up to  $2\text{--}8\times$  faster than prior implementations. Additionally, Mamba-2 introduces analogs to multi-head attention by decomposing the input into multiple parallel SSM channels, enhancing model expressivity without compromising efficiency.

Despite these advances, research in vision and language modeling has shown that SSMs still lag behind comparably sized Transformer models in performance [37, 38, 39]. To address this gap, hybrid architectures have been proposed that combine SSM blocks with Transformer blocks to leverage the strengths of both [37]. Other approaches, inspired by the overall architecture of Transformers, suggest stacking SSM blocks and integrating them with attention and MLP layers [18, 39], aiming to enhance the ability to capture local patterns in SSM-based architectures. While these hybrid models strive to achieve competitive performance, they often compromise the core advantage of SSMs: low model complexity.

Unlike convolutions with fixed receptive fields or recurrent models with sequential bottlenecks, Mamba2 supports both global context and dynamic temporal weighting. This makes it particularly suitable for EEG sequences, which often involve complex dependencies across varied time scales. In our framework, Mamba2 is integrated into a convolutional encoder-decoder architecture (see Section 2), where each Mamba block processes transposed inputs of shape  $(B, T, C)$  and returns to  $(B, C, T)$  to maintain spatial-temporal alignment. This integration provides the dual benefits of low-latency computation and enhanced long-term temporal modeling, essential for EEG-based applications.

## B Time-Frequency Loss Function

SAMBA adopts a reconstruction-based pretraining objective to learn informative latent EEG representations. To preserve both the temporal waveform and spectral structure of EEG signals, we utilized a *Time-Frequency Loss Function*, which combines a Mean Absolute Error (L1) loss in the time domain with a spectral loss computed in the frequency domain.

Given a predicted EEG signal  $\hat{y} \in \mathbb{R}^{B \times T}$  and the corresponding ground truth signal  $y \in \mathbb{R}^{B \times T}$ , the total loss is defined as:

$$\mathcal{L}_{\text{TF}} = \alpha \cdot \mathcal{L}_{\text{L1}} + \beta \cdot \mathcal{L}_{\text{Spec}}, \quad (12)$$

where  $\alpha$  and  $\beta$  are weighting coefficients, and we use  $\alpha = \beta = 1$  in all experiments.

### B.0.1 L1 Loss (Time Domain Reconstruction)

The first term,  $\mathcal{L}_{\text{L1}}$ , measures the Mean Absolute Error (MAE) between the predicted and ground truth waveforms:

$$\mathcal{L}_{\text{L1}} = \frac{1}{BT} \sum_{i=1}^B \sum_{t=1}^T |y_{i,t} - \hat{y}_{i,t}|. \quad (13)$$

Compared to Mean Squared Error (L2 loss), L1 loss penalizes large errors less aggressively, encouraging stable reconstruction and better generalization [40]. To compensate for its non-smooth gradients, we employ a OneCycle learning rate schedule during training.

### B.0.2 Spectral Loss (Frequency Domain Reconstruction)

The second term,  $\mathcal{L}_{\text{Spec}}$ , encourages accurate spectral reconstruction by minimizing the squared difference between the real-valued Fourier spectra:

$$\mathcal{L}_{\text{Spec}} = \frac{1}{BT} \sum_{i=1}^B \sum_{j=1}^{T/2+1} |\mathcal{F}(y_i)_j - \mathcal{F}(\hat{y}_i)_j|^2, \quad (14)$$

where  $\mathcal{F}(\cdot)$  denotes the real-valued discrete Fourier transform (rFFT) along the temporal dimension. Unlike hand-crafted frequency band loss, this term compares the full spectrum directly, allowing the model to preserve global oscillatory properties of EEG such as alpha, beta, and theta rhythms.

## C Details of Dataset

### C.1 BNCI2014-001

BNCI2014-001 dataset of MOABB, also refer to BCI Competition IV-2a [41] consists of EEG recordings from nine subjects obtained using 22 Ag/AgCl electrodes (EEG channels) according to the international 10-20 system with a sampling rate of 250 Hz. The EEG data for each subject were recorded in two sessions on different days, one for the training set and one for the testing set. Each session consisted of six runs with short breaks between them, and each run contained 48 trials. In other words, each session comprised 288 trials across four classes, with each class containing 72 trials. Each class represented a motor imagery task: imagining the movement of the left hand (L), right hand (R), both feet (F), and tongue (T). This resulted in a total of 144 trials for each MI task, with 72 from the training set and 72 from the testing set. The timeline for a trial is approximately 7.5 seconds. At the start of a trial ( $t = 0$  s), a fixation cross appears on the black screen. A cue is then displayed at  $t = 2$  s for 1.25 seconds. Upon seeing the cue, subjects perform the corresponding MI task until the fixation cross disappears from the screen at  $t = 6$  s. A short break follows each trial. Recordings from  $t = 3$  s to  $t = 6$  s are used for further analysis.

### C.2 GrosseWentrup2009

The GrosseWentrup2009 dataset is a motor imagery EEG dataset recorded from 10 healthy subjects using 128 electrodes placed according to the extended 10–20 system [42]. Each subject performed 150 trials of haptic motor imagery for both the left and right hands, totaling 300 trials per subject. During each 7-second trial, subjects imagined hand movement based on an arrow cue following a 3-second fixation. EEG was recorded at 500 Hz with Cz as reference, then referenced to the common average offline. No artifact correction or trial rejection was applied. Electrode positions were recorded in 3D using an ultrasound tracking system. The dataset includes detailed electrode coordinates and is available via MOABB[24].

### C.3 Physionet MI

The Physionet Motor Imagery dataset [43] includes EEG recordings from 109 subjects performing four motor imagery tasks. EEG was recorded from 64 channels at 160 Hz using the BCI2000 system. Each subject completed 14 runs, including both executed and imagined movements involving hands and feet. Each trial lasted 3 seconds. The dataset is publicly available via PhysioNet and is widely used for benchmarking motor imagery classification models.

### C.4 P300-A, B, C

The P300-A and P300-B datasets are from BCI Competition III dataset II [44], and the P300-C dataset is from BCI Competition II dataset IIb [45]. Each dataset contains data from a single subject. EEG signals were recorded using 64 electrodes at a sampling rate of 240 Hz. The Farwell and Donchin paradigm was used. Subjects were shown a 6×6 matrix of symbols. All rows and columns were randomly intensified at a frequency of 5.7 Hz. When the row or column containing the target symbol was intensified, a P300 evoked potential was elicited in the subject’s brain. When other rows or columns were intensified, no P300 component was present. The P300 response is elicited by rare target stimuli, while frequent non-target stimuli do not generate a P300 response. In each trial, six rows and six columns were randomly intensified, with only one row and one column corresponding to the target symbol. This results in two P300 trials and ten non-P300 trials per sequence. Each intensification lasted for 100 ms, followed by a 75 ms blank screen. Each sequence included 12 intensifications, and the sequence was repeated 15 times for each target symbol. The testing set of P300-C includes 31 symbols, while the testing sets of P300-A, P300-B each include 100 symbols. For representation analysis in Fig. 2, repetitions are aggregated to form augmented P300 datasets. Specifically, all

Table 7: We compare SAMBA’s performance across two pretraining and linear-probing settings.

Dataset	SAMBA		SAMBA(w/o quantile)		SAMBA*		SAMBA*(w/o quantile)	
	Acc.	W-F1	Acc.	W-F1	Acc.	W-F1	Acc.	W-F1
PhysionetMI	28.29	0.3070	<b>29.83</b>	<b>0.3226</b>	27.41	0.2985	29.04	0.3140
GrosseWentrup	53.50	0.5349	57.72	0.5772	57.50	0.5749	<b>57.78</b>	<b>0.5776</b>
BNCI2014-001	32.00	0.3220	<b>33.28</b>	<b>0.3285</b>	29.54	0.2954	32.25	0.3198
P300-A	56.26	0.6181	53.35	0.5927	<b>58.01</b>	<b>0.6323</b>	56.22	0.6179
P300-B	<b>57.76</b>	<b>0.6303</b>	53.37	0.5929	57.23	0.6260	54.75	0.6051
P300-C	<b>61.68</b>	<b>0.6615</b>	58.58	0.6377	60.85	0.6546	57.69	0.6327

target P300 trials are averaged across repetitions, and the same to the non-target trials. This aggregation changes the original 1:5 ratio of P300 to non-P300 to approximately 1:1, which facilitates clearer visualization of representation learning.

### C.5 Attention Dataset

The Attention Dataset was collected through an experiment where subjects completed four tasks—two visual and two auditory—designed to assess attention in classifying repeated stimuli. In visual tasks, participants viewed four-digit numbers and clicked when the same number appeared consecutively, with a total duration of 640 seconds. In auditory tasks, they listened to three words and clicked when a word was repeated in sequence, lasting 540 seconds. Each subject completed a total stimulus time of 19 minutes and 40 seconds. Data were recorded using a 14-channel Emotiv Epoc headset, generating multivariate time-series data. After preprocessing and manual labeling, data from 31 subjects were collected, with 4 excluded due to poor quality.

### C.6 Alpha & Crowdsourced

The Alpha and Crowdsourced datasets contain EEG recordings from eyes-open and eyes-closed resting-state tasks. The Alpha data are extracted from the Attention and STEW datasets during resting-state segments where subjects alternated between eyes-open and eyes-closed conditions. The Crowdsourced dataset [21] includes recordings from 60 participants, among whom 13 completed both conditions using EPOC+, EPOC X, or EPOC devices with 14 channels. EEG signals were originally sampled at 2048 Hz and later downsampled to 128 Hz. Raw EEG data, along with preprocessing and analysis scripts, are publicly available on the Open Science Framework (OSF).

### C.7 DriverDistraction

DriverDistraction contains EEG tasks from Driver Distraction Detection, which was obtained by recording EEG brain activity from 17 participants while they engaged in a driving simulation for around 40 minutes. During the simulation, participants carried out various distraction tasks, which can be categorized into three main types: (1) conversing with a passenger, (2) interacting with a mobile phone (including texting and calling), and (3) engaging in problem-solving activities. EEG signals were captured at a sampling rate of 128 Hz using the Emotiv Epoc EEG headset, which records data from 14 channels. The resulting dataset is a multivariate time series with 14 input variables and approximately 5.5 million records. Each time point in the dataset was manually labeled according to the specific activity being performed.

### C.8 STEW

The STEW dataset contains EEG tasks from Driver Distraction Detection, which is a publicly available dataset [46] that consists of raw EEG recordings collected from 48 participants who took part in a multitasking workload experiment using the SIMKAP multitasking test. Prior to the test, baseline brain activity at rest was also recorded. EEG signals were captured using a 14-channel Emotiv EPOC headset at a sampling rate of 128 Hz, resulting in 2.5 minutes of recorded data per participant. After each stage of the experiment, participants assessed their perceived mental workload on a scale from 1 to 9, with these ratings stored in a separate file. Additionally, the dataset includes binary class labels, where workload ratings greater than 4 are categorized as high, while ratings of 4 or below are classified as low. These labels are utilized for specific analytical purposes. The STEW dataset is available upon request via IEEE DataPort.

### C.9 DREAMER

The DREAMER dataset contains EEG tasks from Dataset is emotion detection [47], which consists of electroencephalogram (EEG) and electrocardiogram (ECG) recordings collected during affective stimulation using audio-visual

clips. Signals were recorded from 23 participants using a 14-channel Emotiv EPOC device. Each stimulus was followed by self-reported ratings of valence, arousal, and dominance. In this work, only EEG data and arousal labels are used for classification. ECG signals are excluded.

## D Representation Extraction for Probing

To enable efficient inference for downstream classification, especially with long EEG sequences, we design a customized representation extraction approach for linear probing. Features are captured from the encoder using a forward hook, avoiding direct flattening or global pooling. Given an input EEG sequence  $\mathbf{X} \in \mathbb{R}^{B \times C \times T}$ , we register a forward hook on a target encoder module to obtain hidden representations  $\mathbf{F} \in \mathbb{R}^{B \times C \times T'}$  during inference. We then summarize  $\mathbf{F}$  along the temporal axis using descriptive statistics: minimum, maximum, mean, standard deviation, and quantiles (0.05–0.95). The resulting vector is defined as:

$$\mathbf{z} = [\min, \max, \mu, \sigma, Q_{0.05}, Q_{0.25}, Q_{0.5}, Q_{0.75}, Q_{0.95}], \quad (15)$$

yielding a compact feature tensor  $\mathbf{z} \in \mathbb{R}^{B \times C \times 9}$ , where each of the nine dimensions encodes a temporal statistic per channel. For linear classification, it is flattened into a vector of shape  $\mathbb{R}^{B \times (9 \cdot C)}$ . For non-linear probing, the tensor  $\mathbf{z}$  is directly fed into the MLP.

## E Additional Ablation Study

SAMBA-E was pretrained on six datasets across five distinct EEG tasks, totaling 311,011 samples. We also created a variant, SAMBA\*, using five datasets across four tasks, comprising 140,765 samples to compare performance in the two setting, as shown in Table 7. Overall, SAMBA-E, trained on a more diverse corpus, outperforms SAMBA\* on most downstream tasks, except for P300-A. This exception may be attributed to the inclusion of the DREAMER dataset, which, considering the data size, may have shifted the model’s attention from frontal-lobe dominant features (common in P300-A) toward more distributed brain regions. Furthermore, quantile-based representation methods for linear probing consistently outperform their mean-based counterparts on P300 ERP tasks, whereas mean-based representations perform better on motor imagery (MI) tasks. This suggests that task-specific encoding strategies for linear probing are essential for maximizing downstream performance in EEG foundation models. In future work, we plan to further expand the pretraining EEG corpus to enhance the model’s robustness and generalizability across a broader spectrum of EEG paradigms and recording devices.

## F Evaluation Metrics

The metrics used in the experiments for evaluation are summarized below.

- Averaged Single-Channel Mean Squared Error (ACMSE)** measures the reconstruction quality by quantifying how closely the output EEG signal matches the original input. ACMSE is utilized for the pretraining phase evaluation.
- Balanced Accuracy (ACC)**: computes the average recall across all classes, making it a more reliable metric than standard accuracy, especially in imbalanced datasets. Balanced ACC is utilized for the downstream phase evaluation.
- AUROC**: Area under the ROC curve. Used primarily for binary classification. AUROC is utilized for the downstream phase evaluation.
- Weighted F1 (W-F1)**: Harmonic mean of precision and recall, weighted by class frequency. Used for multi-class evaluation. W-F1 is utilized for the downstream phase evaluation.