# CoMET: A Contrastive-Masked Brain Foundation Model for Universal EEG Representation

**Ang Li**[1, 2, 3*], **Zikai Wang**[1, 3*], **Liuyin Yang**[2], **Zhenyu Wang**[1], **Tianheng Xu**[1, 3],
**Honglin Hu**[1, 3], **Marc M. Van Hulle**[2]

[1]Shanghai Advanced Research Institute, Chinese Academy of Sciences
[2]Faculty of Medicine, KU Leuven,
[3]University of Chinese Academy of Sciences
(huhl@sari.ac.cn)

## Abstract

Electroencephalography (EEG) is a non-invasive technique for recording brain activity, widely used in brain-computer interfaces, clinic, and healthcare. Traditional EEG deep models typically focus on specific dataset and task, limiting model size and generalization. Recently, self-supervised brain foundation models have emerged and been applied to various downstream tasks. Nevertheless, these models still have limitations: current SOTA models typically rely on masked reconstruction strategy; however, EEG features of adjacent channels are highly correlated, which causes the pre-training to overly focus on low-dimensional signal-similarity features in local regions and neglect the global discriminative patterns vital for downstream tasks. To address these limitations, we propose a brain foundation model called CoMET. Specifically, we employ the masked autoencoder with re-designed patching and embedding for EEG as backbone and devise a novel contrastive learning framework with mirror-scale augmentation to strengthen the global discrimination ability. CoMET is pre-trained on mixed EEG datasets over 3000 subjects with over one million samples. It is evaluated on ten different downstream datasets, and the SOTA results demonstrate CoMET's superior ability in extracting universal EEG representations and strong clinical potential.

## Introduction

Electroencephalography (EEG) is a non-invasive technique that captures the amplitude of electrical signals through electrodes placed on the scalp (Teplan et al. 2002). Due to its non-invasive nature and portability, EEG plays a vital role in brain-computer interfaces (BCIs) and healthcare (Edelman et al. 2019). EEG signals can be formulated as a matrix $S \in \mathbb{R}^{C \times T}$ (Jiang, Zhao, and Lu 2024), where $C$ represents the number of channels that could vary across devices, and $T$ represents the sample length that depends on the sampling rate and the collection time. By analyzing the features of EEG, the BCI system can obtain a person's intention, reaction and mental state, such as motor imagery (Al-Saegh, Dawwd, and Abdul-Jabbar 2021), visual evoked potential classification (Norcia et al. 2015), emotion recognition (Li et al. 2022), and seizure detection (Chen et al. 2025), etc.

---

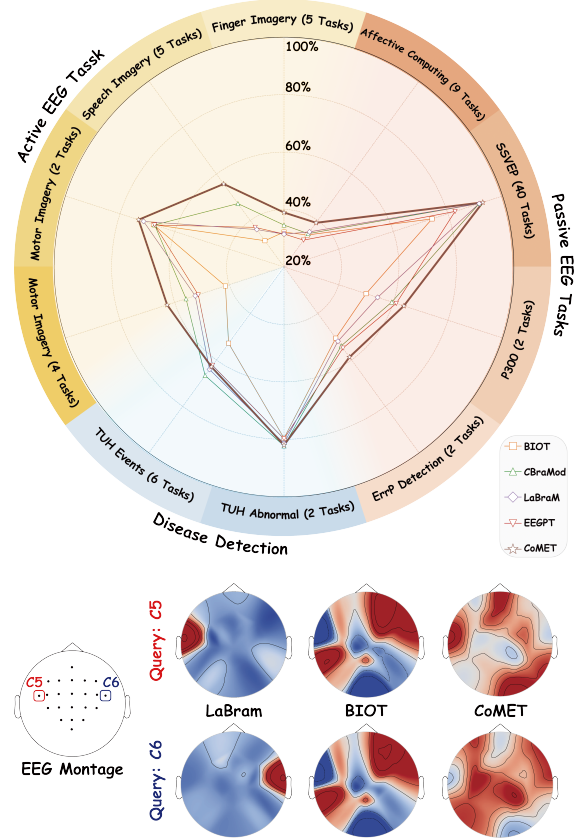*These authors contributed equally.



Figure 1: **Top**: Balanced accuracy on ten downstream datasets. **Bottom**: Self-attention maps on channel-level scalp plots. LabraM (MEM, ICLR'2024) mainly focuses on local area with similar signal of query channel. BIOT (CL, NeurIPS'2023) captures global brain information but collapses into homogeneous attention for different query channels. CoMET (ours) focus on whole-brain in a balanced way with attention diversity.

The EEG signal originates from the cerebral cortex, but the signal is demoted by the intracranial fluid, skull and scalp, causing the signal-to-noise ratio to be low (Nunez and Srinivasan 2006). To decode EEG, researchers have de-

signed various machine learning and deep learning models, such as FBCSP (Ang et al. 2008), EEGNet (Lawhern et al. 2018), EEGConformer (Song et al. 2022), etc. These models typically focus on EEG samples from specific tasks and devices. When the task and/or device changes, the model needs to be retrained. Moreover, the size of a single dataset is limited when originating from a targeted EEG experiment. It necessitates constrained parameters to avoid overfitting, hindering its capacity to learn robust and generalizable EEG representations (Wang et al. 2024).

Recently, inspired by transformer-based (Vaswani et al. 2017) self-supervised learning (SSL) on natural language processing (NLP) (Devlin et al. 2019) and computer vision (CV) (Bao et al. 2021; He et al. 2022), some studies have proposed brain foundation models that are pre-trained on a vast amount of heterogeneous EEG data and further adapted to downstream datasets of different tasks. These methods typically slice EEG sample into patches according to the channel-time dimension, add positional embeddings and feed them into a transformer to learn generalizable representations of EEG signals. The SSL strategies of these models fall into two categories. The first comprises contrastive learning (CL), such as BENDR (Kostas, Aroca-Ouellette, and Rudzicz 2021) and BIOT (Yang, Westover, and Sun 2023), which learn view-invariant EEG representations by minimizing the distance between positive pairs and maximizing the separation from negative pairs in the representation space. The second is masked EEG modeling (MEM), such as LabraM (Jiang, Zhao, and Lu 2024), EEGPT (Wang et al. 2024), CBraMod (Wang et al. 2025), which randomly masks a portion of EEG patches and learn the shallow features by reconstructing the original signals or features through visible patches. Although both strategies successfully distill generic EEG representations across diverse datasets and tasks, several issues remain as follows.

**Problems:** As discussed in the existing study (Park et al. 2023), CL mainly operates at "sample-level" and captures global patterns, but the homogeneous token representations make the encoder suffer from attention collapse, resulting the requirement of high-quality data to converge. On the other hand, MEM functions at the "token level" and emphasizes local shallow features. Although this makes it easier to converge during pre-training, it struggles to model global discriminative representations. These issues are even more pronounced when working with EEG data. The EEG signal is temporally non-stationary, making it easier for CL to overfit to task-irrelevant features due to limited-data-induced attention collapse. In contrast, because of volume conduction in EEG, signals recorded at neighbouring electrodes are sourced from similar neurons. Electrodes spaced less than 10 cm apart exhibit high zero-lag correlations across virtually all frequency bands (Brunner et al. 2016), which leads the MEM-based methods always focus on neighbouring visible channels' tokens, learn low-dimensional local similarity features to complete the pre-training tasks, and reduce the ability to represent the global discrimination that is important for downstream tasks. **Figure 1** provides an intuitive visualization of the attention differences between the two strategies (more visualization comparison please refer to Appendix 8).

Based on the above analysis, a plausible idea is to leverage CL to enhance MEM's ability to capture global patterns. Some studies of vision models (Huang et al. 2023) have explored constructing positive and negative pairs through pixel shifting, and the integration of CL and mask image modeling shows better performance than using either alone. However, EEG data differ from images in many aspects, making it more challenging to construct contrastive pairs:

- EEG has fixed channel locations. Spatial shifting will disrupt the spatial channel dependence.

- In contrast to images, cross-sample channel similarity in EEG is consistent. Temporal shifting or shuffle operations rarely change the similarity.

- The channel dependency of EEG signals leads to performance degradation when a single linear transformation is applied across different channel combinations for building contrastive pairs.

To address these issues, we propose a new EEG pre-training framework that integrates strengths of each strategy. Our approach integrates the following components: mirror-scale augmentation, contrastive global representation, channel-time decoupling embedding, and masked token reconstruction inspired by MAE (He et al. 2022).

Based on this framework, we train a large brain foundation model, called CoMET (Contrastive Masked Encoding Transformer) for universal EEG feature extraction. CoMET is pre-trained on over one million EEG samples from different BCI datasets. For the downstream task, we adopt the linear probing strategy that freezes the encoder and evaluate the CoMET on ten popular BCI and medical datasets. Our experiment demonstrates CoMET's state-of-the-art performance in extracting local and global discriminative EEG features for downstream tasks. The contributions of this paper are as follows:

- Providing a 151-million-parameter brain foundation model for various BCIs and clinical applications.

- Recognizing the challenge of EEG neighboring channels similarity for MEM, and designing mirror-scale augmentation and contrastive global representation, effectively to strengthen the global discrimination ability.

- Adoption of MAE-inspired asymmetric encoder and decoder framework that redesigns patching and embedding for EEG, further enhancing the model's ability to capture local features, and demonstration on ten downstream datasets of different tasks and heterogeneous data formats, including different fine-tuning strategies. The results demonstrate that CoMET outperforms other SOTA models, and the performance improvement complies with the scaling law.

## Related Work

The brain foundation model is trained on a large amount of EEG data, typically through self-supervised learning, and can be applied on various downstream tasks. BENDR (Kostas, Aroca-Ouellette, and Rudzicz 2021) trains
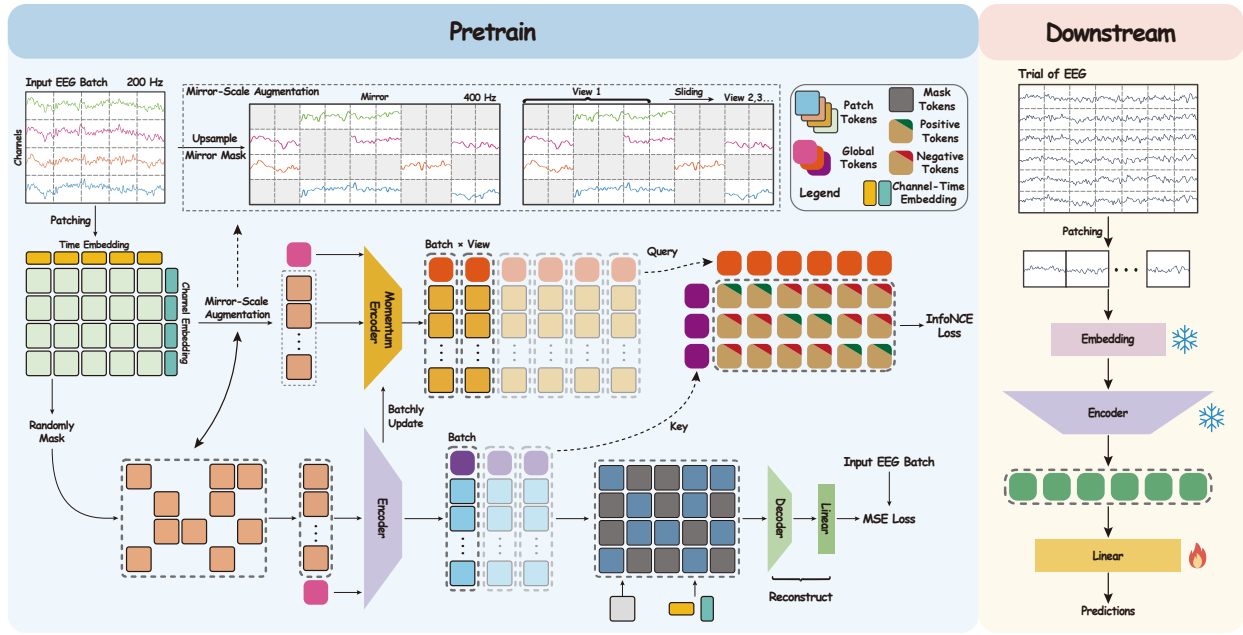
Figure 2: The structure of CoMET. (1) **Pre-training Stage**: There are two branches during the pre-training stage. **Upper** is the contrastive branch. The momentum encoder receives mirror-scale augmented views (mirror masked compared to input of reconstruction branch) and outputs global tokens. **Lower** is the reconstruction branch. Raw EEG signals are patched, added channel-temporal embeddings and randomly masked 50% channels tokens in each time step. Then the visible tokens are sent to encoder together with global token. The output patch tokens are appended with learnable mask tokens and forwarded into decoder to compute the MSE Loss with original signals. The InfoNCE loss is computed between the global tokens of the MEM and momentum encoders, where positive pairs are from the same sample and negative pairs are from different samples in the same batch. (2) **Downstream Stage**: The pre-trained encoder is frozen and adapted to different tasks via linear probing.

a Transformer encoder with an InfoNCE contrastive objective on temporally cropped views of the same recording, yielding sample-invariant representations that generalise well across tasks but suffering attention collapse on fine-grained patterns. BIOT (Yang, Westover, and Sun 2023) extended this paradigm to a multi-biosignal scenario, enriching inputs with explicit channel & time-position embeddings so the model can align unseen heterogeneous datasets, but it still relies solely on view consistency and faces attention collapse. Afterwards, researchers introduced the masked image modeling (Bao et al. 2021; He et al. 2022) in CV, to make the model converge faster and avoid attention collapse. LabraM (Jiang, Zhao, and Lu 2024), inspired by BEiT (Bao et al. 2021), first learns the EEG tokenizer by predicting the frequency domain information and subsequently reconstructs the features of occluded patches via a masked reconstruction strategy. EEGPT (Wang et al. 2024) adopts dual self-supervision by combining spatio-temporal alignment and masked reconstruction, but its pre-trained model only captures spatial features and the temporal dependencies are not utilized for downstream. CBraMod (Wang et al. 2025), also based on MEM, devises a criss-cross transformer as the backbone and use the asymmetric convolutional positional encoding scheme to encode spatial-temporal positional information.

## Methodology

In this section, we give a detailed description of the proposed CoMET model (**Figure 2**). The pre-training stage contains two components. **Masked EEG Modeling**. Departing from previous brain foundation models LabraM, EEGPT and CBraMod, our method employs an asymmetric encoder–decoder design inspired by MAE, which introduces mask tokens only in the decoder phase. The pipeline consists of patching, channel-time embedding, randomly masking, the transformer encoder, double embedding, and the decoder. The goal is to recover the original masked signal. **Contrastive Learning**. Generating positive views of input samples through mirror-scale augmentation, and a learnable global token is used to aggregate discriminative representations from visible patches, followed by the InfoNCE loss to optimize contrastive learning between positive and negative features. For the downstream stage, the linear probing is used to effectively evaluate the cross-domain feature extraction capability of the pre-trained model.

### Patching and Channel-Time Embedding

Given an EEG sample represented as $S \in \mathbb{R}^{C \times T}$, where $C$ is the number of channels and $T$ the sample length. We first partition $S$ into channel-temporal patches. This is achieved by applying a one-dimensional convolution with same kernel size and a stride of length $l$. It produces non-overlapping

vectors as follows:

$$\{x_{i,j} \in \mathbb{R}^d | i \in (1,2,...,C), j \in (1,2,...,N)\}, \quad (1)$$

where $d$ denotes the embedding dimension and $N = T/l$ the number of temporal patches. For channel $i$ and temporal patch $j$, the token embedding is

$$e_{i,j} = x_{i,j} + e_i^c + e_j^t, \quad (2)$$

where the learnable channel embeddings $\{e_i^c\}_{i=1}^C \subset \mathbb{R}^d$ are looked up from channel names, and the learnable temporal embeddings $\{e_j^t\}_{j=1}^N \subset \mathbb{R}^d$ are added in temporal order. The complete set of patch tokens is

$$E = \{e_{i,j} | i \in (1,2,...,C), j \in (1,2,...,N)\}. \quad (3)$$

## Masking

During the encoding phase, we randomly keep half of the patches at each time position and remove the remaining patches rather than replacing them with masked tokens (mask ratio=0.5, according to CBraMod). The sampling patches keep the original channel and temporal embedding, and each time position has the same number of patches from different channels:

$$E_v = \{e_{i,j} | i \in C_j, j \in (1,2,...,N)\}, \quad (4)$$

where $C_j$ is the visible channel at each time position.

## MEM Encoder and Decoder

Before feeding visible patches $E_v$ to the encoder, we append a learnable global token $e_g \in \mathbb{R}^d$ to summarize global discriminative representations. Unlike vision model, which derives linear projection after the final layer, global token participates in every self-attention layer. The reorganized paches is defined as $\widetilde{E} = E_v \cup \{e_g \in \mathbb{R}^d\}$.

Then, visible patches and the global token are fed into the encoder:

$$F^{(n)} = Attention(\widetilde{E}W_n^Q, \widetilde{E}W_n^K, \widetilde{E}W_n^V), \quad (5)$$

where $F_n$ is the $n$-th head attention. The final output of encoder is denoted as:

$$E_R = \{e_{i,j}^r \mid i \in C_j, j = 1, \ldots, N\} \cup \{e_g^r\}. \quad (6)$$

We use output patch embeddings to reconstruct masked patches, and use the global token embedding $e_g^r$ for discriminative training through CL in the following part. We insert the learnable masked tokens $e_m \in \mathbb{R}^d$ at every masked position. Because the channel-time embeddings of the visible tokens have already been updated by the MEM encoder, we re-add the learnable channel ($e_i^c$) and temporal ($e_j^t$) positional embeddings to each token. Then all tokens are fed into the transformer decoder:

$$\{rec_{i,j}\} = \text{Decoder}\Big(\{e_{i,j}^r + e_i^c + e_j^t \mid (i,j) \in \mathcal{V}\}$$
$$\cup \{e_m + e_i^c + e_j^t \mid (i,j) \in \mathcal{M}\}\Big), \quad (7)$$

where $\mathcal{V}$ and $\mathcal{M}$ denote the index sets of visible and masked tokens, respectively, and $rec_{c,t}$ the reconstruction of the patch located at channel $i$ and time step $j$.

## Momentum Encoder

The momentum encoder (He et al. 2020) is introduced to generate contrastive targets for the MEM encoder $\mathcal{F}_o$ to learn global EEG representations. The momentum encoder shares the same architecture as $\mathcal{F}_o$ but operates on the mirror-scale augmentation signal.

The inputs of the MEM encoder and momentum encoder can be considered as two perspectives of the brain state. Contrastive learning on these paired views strengthen the encoder's ability to derive global brain-state features and prepares the encoder for downstream tasks that involve multiple channel combinations.

To focus the contrastive loss on global semantics, we retain only the global tokens from the momentum encoder and discard its patch tokens, thereby reducing the similarity among adjacent patches. The momentum encoder $\mathcal{F}_m$ is updated through exponential moving average (EMA). Denoting the parameters of $\mathcal{F}_m$ and $\mathcal{F}_o$ as $\theta_m$ and $\theta_o$, the update rule is

$$\theta_m \leftarrow \mu\theta_o + (1-\mu)\theta_m. \quad (8)$$

## Mirror-Scale Augmentation

Typically, two distinct views of the same input are required in contrastive learning. Because the MEM branch already provides one view, we construct a complementary view to form positive and negative pairs.

We propose a mirror-scale augmentation that (i) preserves global correlations with the original view and avoids token-level similarity that could cause feature collapse and (ii) Generates more paired samples. Let the visible channels for MEM be

$$V_a = \{v_{i,j} \mid i \in V^j, j = 1, \ldots, N\}. \quad (9)$$

And the visible channels fed to the momentum encoder are defined as:

$$V_b = \{v_{i,j} | i \in \{C_{all} \setminus V^j\}, j = 1, \ldots, N\}. \quad (10)$$

Then the input signal is upsampled from the original sampling rate $f$ to $2f$, doubling the number of temporal patches from $N$ to $2N$ while keeping the patch length and visible-channel list unchanged. And the visible patch set after upsampling becomes:

$$V_b^{(2f)} = \left\{ v_{i,j} \mid i \in \{C_{all} \setminus V^{\lceil j/2 \rceil}\}, j = 1, \ldots, 2N \right\}. \quad (11)$$

Afterwards, a sliding window of length $N$ with stride 1 is applied along the temporal axis, which generates a group of different augmentation views for each sample. The augmented views from the same sample as the MEM encoder are treated as positive pairs, and the augmented view of different samples in the same batch are treated as negative pairs.

# Training objective

## Pre-training Stage

Our pre-training objective combines a reconstructive loss that restores masked EEG patches and a contrastive loss that

aligns global representations across augmented views. Given the set of masked indices $\mathcal{M}$, the reconstructive loss is calculated through mean squared error (MSE):

$$\mathcal{L}_R = \frac{1}{|\mathcal{M}|} \sum_{(i,j) \in \mathcal{M}} \|rec_{i,j} - raw_{i,j}\|_2^2. \qquad (12)$$

For the contrastive loss, we employ the InfoNCE loss (Oord, Li, and Vinyals 2018), which simultaneously pulls positive view from the same sample closer while pushing away negative views from different samples. Let $e_g^1$ and $e_g^2$ denote the global tokens produced by the MEM encoder and the momentum encoder, respectively. Their cosine similarity is

$$\rho = \frac{\langle e_g^1, e_g^2 \rangle}{\|e_g^1\|_2 \, \|e_g^2\|_2}. \qquad (13)$$

For each sample we treat $\rho^+$ (similarity between the two views of the same sample) as the positive pair, and $\rho_i^-$ as the negative pair obtained by matching the MEM-view of one sample with the momentum-view of other samples in the same batch. The loss is defined as:

$$\mathcal{L}_C = -\log \frac{\exp(\rho^+/\tau)}{\exp(\rho^+/\tau) + \sum_{i=1}^{K-1} \exp(\rho_i^-/\tau)} \qquad (14)$$

where $\tau$ is the temperature. The total loss is calculated by combining $\mathcal{L}_R$ and $\mathcal{L}_C$.

## Downstream Stage

For downstream stage, we adopt the linear-probing strategy (Alain and Bengio 2016): the parameters of the pre-trained encoder are kept frozen, and only a lightweight linear head is updated. All the patch embeddings produced by the encoder are fed to task-specific linear layer. In this way, the downstream performance directly reflects the quality of the representations learned during pre-training. In addition, we also conduct full parameter fine-tuning experiment to show adaptability, with results shown in the Appendix 9.

## Experiment

### Datasets and Model Settings

CoMET is pre-trained on mixed datasets, including HBN-EEG (Shirazi et al. 2024), Stieger21 (Stieger, Engel, and He 2021), M3CV (Huang et al. 2022) and SEED (Zheng and Lu 2015). For downstream tasks, ten widely used datasets are evaluated across multiple BCIs and clinical tasks. More details about the datasets are introduced in the Appendix 2.

### Pre-training

For pre-training, We randomly divide 90% of the training data into the training set and the remaining 10% into the validation set. The mask ratio is set to 50% according to CbraMod. We use the batchsize of 256 and the epochs of 100. The model is optimized using the AdamW optimizer. The max learning rate is set to 5e-4 and scheduled by CosineAnnealingLR. The pre-training experiment is conducted on a single node, equipped with 4 NVIDIA H100 (80GB) GPUs. The main software versions are Python

3.11.4, PyTorch 2.0.1 and CUDA 11.8. Subsequently, only simple and essential preprocessing steps are applied, including 0.5-70 Hz bandpass filtering, resampling to 200 Hz, 4 seconds segmenting without overlap, converting units to 0.1 mV. We use three configurations for CoMET with different depths of the encoder and hidden sizes: CoMET-Tiny (5M), CoMET-Base (19M), CoMET-Large (151M). More detailed settings please refer to the Appendix 3.

### Evaluation and Metrics

We compare the CoMET with four SOTA brain foundation models (BIOT, LabraM, EEGPT and CBraMod). We re-implement baseline models using their official codes and released model weights. For data preprocessing, we follow the strategy of BIOT, EEGPT and CBraMod. All models use the same processed data on each dataset to ensure fairness. Downstream experiments were conducted on single node, equipped with 4 NVIDIA A100 (40GB) GPUs. Results shown in **Table 1** are obtained following the original baseline's downstream strategy, and additional evaluations of alternative adaptation strategies are provided in the Appendix 9.

## Results

We conducted a comprehensive evaluation of our proposed CoMET models against strong baselines across diverse EEG tasks. The full results are summarised in **Table 1**. On the widely used motor imagery datasets BCIC IV 2A, CoMET-Large attains balanced accuracy of 62.75% ± 1.62 on 2A, exceeding the previous SOTA CBraMod (55.85% ± 0.97) by an absolute margin of 6.9%. Similarly, on BCIC IV 2B, CoMET-Large reaches 73.22%±0.81, outperforming the strongest baseline LaBraM (71.39% ± 0.28) by 1.83%. For Large-5F (five-finger imagery), CoMET-Large achieves balanced accuracy of 38.97% ± 0.84, surpassing the best baseline by 4.37%. On the BCIC-2020-3 (imagined speech), CoMET-Large achieves the highest performance of 55.72% ± 1.42, outperforming the baseline by 5.81%. For emotion recognition dataset KaggleERN and FACED, CoMET-Large surpasses the leading baseline by 2.67% and 3.14% in balanced accuracy, respectively. As for BCI tasks that required visual stimuli (THUBenchmark, and PhysioP300), CoMET-Large also achieved the best banlanced accuracies of 92.74% and 63.86%, respectively, which are 1.29% and 2.79% higher than the best baseline CBraMod and EEGPT.

On TUAB and TUEV, CoMET-Large marginally underperforms LabraM and CBraMod. However, it is important to note that the TUH dataset is used by CBraMod in pre-training, which includes both TUAB and TUEV. LabraM employs TUEG from TUH for pre-training, which shares a similar distribution with TUAB and TUEV. None of the pre-training data of CoMET originates from the TUH family (TUH, TUAB, TUEV, etc.). When the comparison is restricted to models pre-trained on different source datasets, CoMET still achieves the best performance. Moreover, CoMET's performance increases monotonically from Tiny to Base, and Large, adhering to the scaling law and indicating headroom for further improvements as model capacity grows.

| Model | BCIC IV 2A | | | BCIC IV 2B | | |
|---|---|---|---|---|---|---|
| | **B. Acc** | **Kappa** | **F1** | **B. Acc** | **Kappa** | **F1** |
| BIOT (NeurIPS'23) | 41.44 ± 0.58*** | 21.90 ± 1.19 | 37.37 ± 0.56 | 67.78 ± 0.18*** | 35.56 ± 0.36 | 66.41 ± 2.48 |
| LaBraM (ICLR'24) | 52.49 ± 1.34*** | 36.60 ± 2.12 | 52.32 ± 0.97 | 71.39 ± 0.38* | 42.76 ± 0.76 | 70.68 ± 0.66 |
| EEGPT (NeurIPS'24) | 51.37 ± 0.96*** | 35.17 ± 1.26 | 49.73 ± 0.41 | 67.33 ± 0.61*** | 34.65 ± 1.23 | 67.35 ± 1.03 |
| CBraMod (ICLR'25) | 55.85 ± 0.97** | 41.13 ± 1.30 | 55.08 ± 1.02 | 67.35 ± 0.98*** | 34.70 ± 1.96 | 71.06 ± 4.38 |
| CoMET-Tiny | 58.53 ± 0.93 | 44.70 ± 0.50 | 58.00 ± 0.34 | 71.42 ± 0.16 | 42.82 ± 0.31 | 70.48 ± 0.25 |
| CoMET-Base | 61.66 ± 1.81 | 48.89 ± 2.42 | 60.52 ± 1.85 | 72.71 ± 0.93 | 45.42 ± 1.86 | 72.52 ± 1.64 |
| CoMET-Large | **62.75 ± 1.62** | **51.70 ± 1.84** | **63.37 ± 1.34** | **73.22 ± 0.81** | **46.32 ± 1.92** | **73.36 ± 1.74** |

| Model | Large-5F | | | BCIC2020-3 | | |
|---|---|---|---|---|---|---|
| | **B. Acc** | **Kappa** | **F1** | **B. Acc** | **Kappa** | **F1** |
| BIOT | 32.16 ± 0.04*** | 15.05 ± 0.05 | 31.02 ± 0.41 | 31.31 ± 0.31*** | 14.14 ± 0.39 | 30.66 ± 1.03 |
| LaBraM | 31.41 ± 0.24*** | 14.44 ± 0.27 | 31.21 ± 0.57 | 36.10 ± 0.53*** | 20.13 ± 0.66 | 35.30 ± 0.55 |
| EEGPT | 31.48 ± 0.47*** | 14.80 ± 0.62 | 29.78 ± 0.74 | 37.05 ± 0.54*** | 21.31 ± 0.68 | 36.22 ± 0.98 |
| CBraMod | 34.60 ± 0.72*** | 18.38 ± 0.88 | 33.86 ± 0.95 | 47.29 ± 2.03*** | 34.11 ± 2.54 | 46.66 ± 2.71 |
| CoMET-Tiny | 34.66 ± 0.25 | 18.60 ± 0.33 | 33.75 ± 0.26 | 52.97 ± 1.52 | 41.76 ± 2.19 | 52.80 ± 1.86 |
| CoMET-Base | 37.03 ± 0.68 | 21.66 ± 0.82 | 35.93 ± 1.03 | 54.43 ± 1.33 | 43.04 ± 1.66 | 54.12 ± 1.35 |
| CoMET-Large | **38.97 ± 0.84** | **24.42 ± 0.94** | **39.24 ± 0.85** | **55.72 ± 1.42** | **44.84 ± 1.66** | **55.22 ± 1.03** |

| Model | KaggleERN | | | FACED | | |
|---|---|---|---|---|---|---|
| | **B. Acc** | **Kappa** | **F1** | **B. Acc** | **Kappa** | **F1** |
| BIOT | 50.64 ± 0.01*** | 1.62 ± 0.03 | 71.53 ± 0.01 | 33.77 ± 0.39*** | 25.32 ± 0.46 | 33.32 ± 0.3 |
| LaBraM | 51.95 ± 0.81*** | 4.38 ± 1.75 | 77.21 ± 1.82 | 35.13 ± 1.76** | 26.87 ± 2.00 | 35.05 ± 1.9 |
| EEGPT | 54.92 ± 0.04*** | 11.89 ± 0.12 | 76.77 ± 0.04 | 31.56 ± 0.11*** | 22.93 ± 0.14 | 31.35 ± 0.19 |
| CBraMod | 53.92 ± 0.02*** | 8.15 ± 0.05 | 76.49 ± 0.04 | 34.35 ± 4.34*** | 25.83 ± 4.86 | 33.51 ± 5.1 |
| CoMET-Tiny | 55.45 ± 0.12 | 13.28 ± 0.41 | 77.74 ± 0.78 | 33.28 ± 0.31 | 24.78 ± 0.36 | 33.12 ± 0.35 |
| CoMET-Base | 57.08 ± 0.12 | 14.63 ± 0.33 | 78.92 ± 0.53 | 38.40 ± 0.55 | 30.49 ± 0.64 | 38.02 ± 0.75 |
| CoMET-Large | **58.78 ± 0.88** | **15.24 ± 0.92** | **79.66 ± 1.43** | **39.02 ± 0.38** | **31.47 ± 0.47** | **39.13 ± 0.59** |

| Model | THUBenchmark | | | PhysioP300 | | |
|---|---|---|---|---|---|---|
| | **B. Acc** | **Kappa** | **F1** | **B. Acc** | **Kappa** | **F1** |
| BIOT | 74.17 ± 0.06*** | 73.50 ± 0.07 | 74.24 ± 0.05 | 50.04 ± 0.12*** | 0.09 ± 0.25 | 6.31 ± 0.88 |
| LaBraM | 91.43 ± 0.13 | 91.21 ± 0.13 | 91.43 ± 0.13 | 54.26 ± 0.22*** | 8.38 ± 0.43 | 56.86 ± 1.38 |
| EEGPT | 82.59 ± 0.09*** | 82.15 ± 0.09 | 82.57 ± 0.08 | 61.07 ± 0.27 | 22.15 ± 0.52 | 54.24 ± 0.48 |
| CBraMod | 91.45 ± 0.24 | 91.23 ± 0.25 | 91.43 ± 0.24 | 59.46 ± 0.01*** | 18.85 ± 0.03 | 55.68 ± 0.03 |
| CoMET-Tiny | 91.72 ± 1.85 | 90.17 ± 0.67 | 89.66 ± 1.49 | 61.36 ± 0.50 | 22.58 ± 1.00 | 63.44 ± 0.58 |
| CoMET-Base | 92.01 ± 1.09 | 91.81 ± 1.10 | 92.00 ± 1.11 | 62.36 ± 0.71 | 24.55 ± 1.41 | 60.52 ± 3.63 |
| CoMET-Large | **92.74 ± 1.61** | **92.65 ± 2.46** | **93.50 ± 0.81** | **63.86 ± 1.78** | **25.05 ± 2.01** | **61.72 ± 2.28** |

| Model | TUAB | | | TUEV | | |
|---|---|---|---|---|---|---|
| | **B. Acc** | **AUC** | **–** | **B. Acc** | **Kappa** | **F1** |
| BIOT | 79.59 ± 0.57** | 88.15 ± 0.43 | – | 52.81 ± 2.25** | 52.73 ± 2.49 | 74.92 ± 0.82 |
| LaBraM△ | 81.40 ± 0.19 | 90.22 ± 0.09 | – | 64.09 ± 0.65 | 66.37 ± 0.93 | 83.12 ± 0.52 |
| EEGPT | 79.83 ± 0.30* | 87.18 ± 0.50 | – | 62.32 ± 1.14* | 63.51 ± 1.34 | 81.87 ± 0.63 |
| CBraMod△ | **82.29 ± 0.22** | **92.27 ± 0.11** | – | **66.71 ± 1.07** | **67.72 ± 0.96** | **83.42 ± 0.64** |
| CoMET-Tiny | 80.02 ± 1.34 | 87.58 ± 1.08 | – | 60.31 ± 2.40 | 62.83 ± 0.84 | 80.32 ± 0.85 |
| CoMET-Base | 81.87 ± 0.19 | 89.16 ± 0.24 | – | 62.31 ± 2.40 | 64.05 ± 1.51 | 80.47 ± 2.11 |
| CoMET-Large | 82.02 ± 1.05 | 91.04 ± 0.97 | – | 62.97 ± 1.31 | 66.95 ± 0.31 | 82.88 ± 1.93 |

Table 1: Performance comparison of models on ten downstream datasets with Wilcoxon Signed-Rank Test (CoMET-Large vs. others) (*:p<0.05, **:p<0.01, ***:p<0.001).

## Ablation Study

In order to evaluate the effectiveness of key components in CoMET, we conduct an ablation study in the pre-training stage as follows: 1) w/o $\mathcal{L}_C$: only use masked reconstruction in pre-training; 2) w/o $\mathcal{L}_R$: only use contrastive learning in pre-training; 3) w/o mirror-scale augmentation: use random cropped masking to construct contrastive view augmentation; 4) w/o global representation: Append linear projection after encoder to get global features. The ablation results in **Figure 3** demonstrate that each component of CoMET is indispensable and mutually reinforcing: MEM serves as the backbone, guaranteeing the model's ability to capture local, shallow features, whereas CL further re-
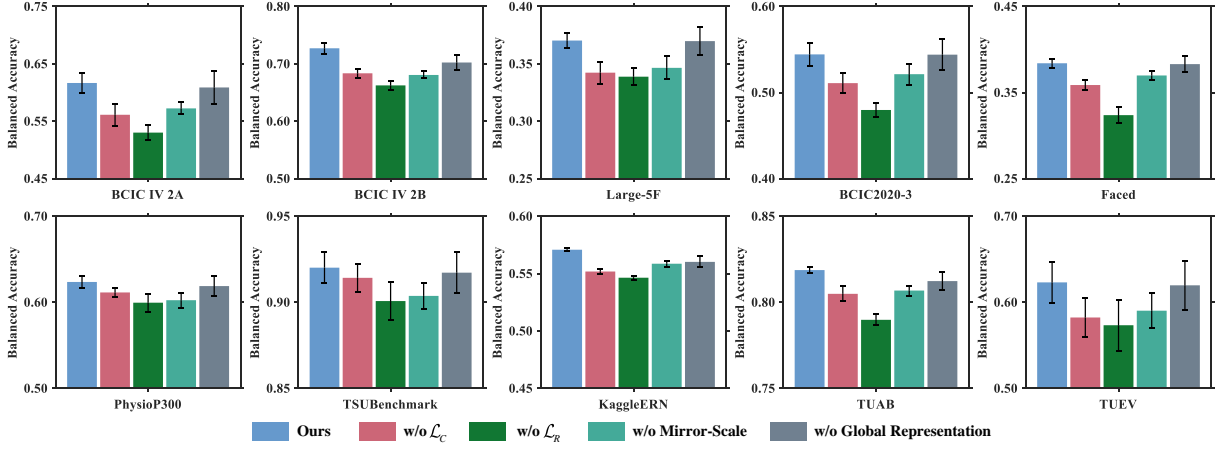
Figure 3: Ablation Study. We study the effort of two losses (w/o $\mathcal{L}_R$ and w/o $\mathcal{L}_C$), augmentation method (random masking w/o Mirror-Scale) and contrastive feature extraction module (w/o Global Representation) on CoMET-Base.

inforces its capacity for global modelling. Removing either the $\mathcal{L}_R$ or the $\mathcal{L}_C$—instead of employing both concurrently—reduces downstream classification accuracy by up to 6.98% (TUEV) and 5.54% (BCIC IV 2A), demonstrating that MEM and CL emphasize different aspects during pre-training yet complement one another effectively. $\mathcal{L}_C$ contributes less on THUBenchmark, possibly because SSVEP focus more on the frequency domain features of specific brain regions (Lin et al. 2006). Mirror-Scale Augmentation ensures that CL aligns global invariant features across different views of the same sample; when the augmented views are instead generated by random masking, the CL branch contributes little and can even be detrimental. The global token contributes little improvement on mean accuracy, but it markedly reduces variance, thereby enhancing the stability of the pre-trained model.

## Discussion

To understand how the reconstructive loss ($\mathcal{L}_R$) and contrastive loss ($\mathcal{L}_C$) contribute to the model's receptive fields, we analyze the average distance between the query token weights and key token weights of different channels (Dosovitskiy et al. 2020). Channels use relative coordinates normalized to (0, 1). Figure 4 (a) shows $\mathcal{L}_C$ captures wider global information and stabilizes between layers. While $\mathcal{L}_R$ focus on local channel's relationships, $\mathcal{L}_R + \mathcal{L}_C$ (ours) can effectively increase the receptive fields. We also use normalized mutual information (NMI) (Strehl and Ghosh 2002) to measure the attention collapse. NMI computes the mutual information between query and key tokens $I(q, k)$ from different channels, and normalize it by their marginal entropies: $\frac{I(q,k)}{\sqrt{H(q)H(k)}}$. Lower NMI indicates attention map are less dependent on query tokens and attention collapse into homogeneity. As shown in Figure 4 (b), $\mathcal{L}_C$ exhibits homogeneous attention collapse across all layers, whereas $\mathcal{L}_R + \mathcal{L}_C$ (ours) markedly enhances attention diversity, even outperforming $\mathcal{L}_R$ in final three layers. This analysis visually confirms the motivation of our article: MEM and CL focus on different EEG features and they are able to gain their respective ad-
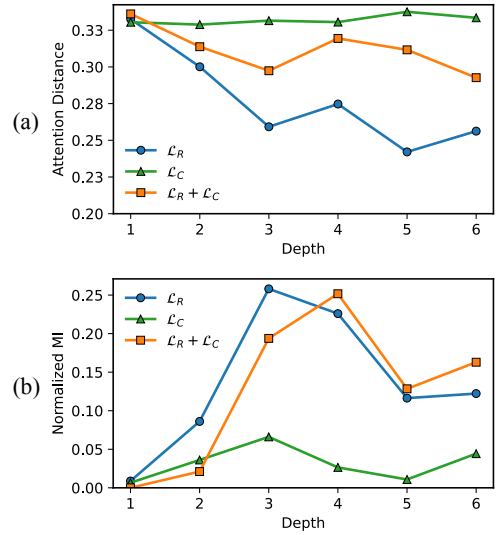


Figure 4: Analysis of (a) attention distance and (b) normalized mutual information across layers. Larger attention distance indicates wider effective receptive fields, and lower Normalized MI indicates attention homogeneous collapse.

vantages non-destructively through aggregation.

## Conclusion

In this study, we propose a contrastive-masked EEG self-supervised framework, training the 151 million parameter brain foundation model CoMET. We employ the masked autoencoder with redesigned patching and embedding strategies for EEG as the backbone, and devise a novel contrastive learning framework with mirror-scale augmentation to effectively to merge the local fast convergence ability of masked EEG reconstruction and the global discrimination ability of contrastive learning. Experiments on ten downstream datasets demonstrate CoMET's SOTA capacity to extract both local context–sensitive and globally discriminative EEG representations, outperforming prior foundation models based solely on contrastive learning or masked re-

construction, especially on these tasks that relies on connection between multi brain regions. The proposed foundation model CoMET has been designed with full consideration of the characteristics of EEG signals, and the new global modeling idea provides a meaningful contribution to the study of brain foundation models, promoting the AI utilization in real-world BCIs, clinical, and healthcare applications.

# References

Al-Saegh, A.; Dawwd, S. A.; and Abdul-Jabbar, J. M. 2021. Deep learning for motor imagery EEG-based classification: A review. *Biomedical Signal Processing and Control*, 63: 102172.

Alain, G.; and Bengio, Y. 2016. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*.

Ang, K. K.; Chin, Z. Y.; Zhang, H.; and Guan, C. 2008. Filter bank common spatial pattern (FBCSP) in brain-computer interface. In *2008 IEEE international joint conference on neural networks (IEEE world congress on computational intelligence)*, 2390–2397. IEEE.

Bao, H.; Dong, L.; Piao, S.; and Wei, F. 2021. BEiT: BERT Pre-Training of Image Transformers. In *International Conference on Learning Representations*.

Brunner, C.; Billinger, M.; Seeber, M.; Mullen, T. R.; and Makeig, S. 2016. Volume conduction influences scalp-based connectivity estimates. *Frontiers in computational neuroscience*, 10: 121.

Brunner, C.; Leeb, R.; Müller-Putz, G.; Schlögl, A.; and Pfurtscheller, G. 2008. BCI Competition 2008–Graz data set A. *Institute for knowledge discovery (laboratory of brain-computer interfaces), Graz University of Technology*, 16(1-6): 34.

Chen, J.; Wang, X.; Huang, C.; Hu, X.; Shen, X.; and Zhang, D. 2023. A large finer-grained affective computing EEG dataset. *Scientific Data*, 10(1): 740.

Chen, Z.; Matsubara, Y.; Sakurai, Y.; and Sun, J. 2025. Long-term eeg partitioning for seizure onset detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 14221–14229.

Citi, L.; Poli, R.; and Cinel, C. 2010. Documenting, modelling and exploiting P300 amplitude changes due to variable target delays in Donchin's speller. *Journal of Neural Engineering*, 7(5): 056006.

Committee, B. C. 2022. 2020 International BCI Competition.

Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. 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*, 4171–4186.

Donahue, J.; Jia, Y.; Vinyals, O.; Hoffman, J.; Zhang, N.; Tzeng, E.; and Darrell, T. 2014. Decaf: A deep convolutional activation feature for generic visual recognition. In *Int. Conf. Mach. Learn.*, 647–655.

Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. 2020. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*.

Edelman, B. J.; Meng, J.; Suma, D.; Zurn, C.; Nagarajan, E.; Baxter, B. S.; Cline, C. C.; and He, B. 2019. Noninvasive neuroimaging enhances continuous neural tracking for robotic device control. *Science robotics*, 4(31): eaaw6844.

He, K.; Chen, X.; Xie, S.; Li, Y.; Dollar, P.; and Girshick, R. 2022. Masked Autoencoders Are Scalable Vision Learners. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15979–15988.

He, K.; Fan, H.; Wu, Y.; Xie, S.; and Girshick, R. 2020. Momentum contrast for unsupervised visual representation learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9729–9738.

Huang, G.; Hu, Z.; Chen, W.; Zhang, S.; Liang, Z.; Li, L.; Zhang, L.; and Zhang, Z. 2022. M3CV: A multi-subject, multi-session, and multi-task database for EEG-based biometrics challenge. *NeuroImage*, 264: 119666.

Huang, Z.; Jin, X.; Lu, C.; Hou, Q.; Cheng, M.-M.; Fu, D.; Shen, X.; and Feng, J. 2023. Contrastive masked autoencoders are stronger vision learners. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(4): 2506–2517.

Jiang, W.; Zhao, L.; and Lu, B.-l. 2024. Large Brain Model for Learning Generic Representations with Tremendous EEG Data in BCI. In *The Twelfth International Conference on Learning Representations*.

Kaya, M.; Binli, M. K.; Ozbay, E.; Yanar, H.; and Mishchenko, Y. 2018. A large electroencephalographic motor imagery dataset for electroencephalographic brain computer interfaces. *Scientific data*, 5(1): 1–16.

Kostas, D.; Aroca-Ouellette, S.; and Rudzicz, F. 2021. BENDR: Using transformers and a contrastive self-supervised learning task to learn from massive amounts of EEG data. *Frontiers in Human Neuroscience*, 15: 653659.

Lawhern, V. J.; Solon, A. J.; Waytowich, N. R.; Gordon, S. M.; Hung, C. P.; and Lance, B. J. 2018. EEGNet: a compact convolutional neural network for EEG-based brain–computer interfaces. *Journal of neural engineering*, 15(5): 056013.

Leeb, R.; Brunner, C.; Müller-Putz, G.; Schlögl, A.; and Pfurtscheller, G. 2008. BCI Competition 2008–Graz data set B. *Graz University of Technology, Austria*, 16: 1–6.

Li, X.; Zhang, Y.; Tiwari, P.; Song, D.; Hu, B.; Yang, M.; Zhao, Z.; Kumar, N.; and Marttinen, P. 2022. EEG based emotion recognition: A tutorial and review. *ACM Computing Surveys*, 55(4): 1–57.

Lin, Z.; Zhang, C.; Wu, W.; and Gao, X. 2006. Frequency recognition based on canonical correlation analysis for SSVEP-based BCIs. *IEEE transactions on biomedical engineering*, 53(12): 2610–2614.

Margaux, P.; Emmanuel, M.; Sébastien, D.; Olivier, B.; and Jérémie, M. 2012. Objective and Subjective Evaluation of

Online Error Correction during P300-Based Spelling. *Advances in Human-Computer Interaction*, 2012(1): 578295.

Norcia, A. M.; Appelbaum, L. G.; Ales, J. M.; Cottereau, B. R.; and Rossion, B. 2015. The steady-state visual evoked potential in vision research: A review. *Journal of vision*, 15(6): 4–4.

Nunez, P. L.; and Srinivasan, R. 2006. *Electric fields of the brain: the neurophysics of EEG*. Oxford university press.

Obeid, I.; and Picone, J. 2016. The temple university hospital EEG data corpus. *Frontiers in neuroscience*, 10: 196.

Oord, A. v. d.; Li, Y.; and Vinyals, O. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*.

Park, N.; Kim, W.; Heo, B.; Kim, T.; and Yun, S. 2023. What Do Self-Supervised Vision Transformers Learn? In *The Eleventh International Conference on Learning Representations*.

Shirazi, S. Y.; Franco, A.; Scopel Hoffmann, M.; Esper, N. B.; Truong, D.; Delorme, A.; Milham, M. P.; and Makeig, S. 2024. HBN-EEG: The FAIR implementation of the Healthy Brain Network (HBN) electroencephalography dataset. *bioRxiv*, 2024–10.

Song, Y.; Zheng, Q.; Liu, B.; and Gao, X. 2022. EEG conformer: Convolutional transformer for EEG decoding and visualization. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 31: 710–719.

Stieger, J. R.; Engel, S. A.; and He, B. 2021. Continuous sensorimotor rhythm based brain computer interface learning in a large population. *Scientific Data*, 8(1): 98.

Strehl, A.; and Ghosh, J. 2002. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *Journal of machine learning research*, 3(Dec): 583–617.

Teplan, M.; et al. 2002. Fundamentals of EEG measurement. *Measurement science review*, 2(2): 1–11.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Wang, G.; Liu, W.; He, Y.; Xu, C.; Ma, L.; and Li, H. 2024. Eegpt: Pretrained transformer for universal and reliable representation of eeg signals. *Advances in Neural Information Processing Systems*, 37: 39249–39280.

Wang, J.; Zhao, S.; Luo, Z.; Zhou, Y.; Jiang, H.; Li, S.; Li, T.; and Pan, G. 2025. CBraMod: A Criss-Cross Brain Foundation Model for EEG Decoding. In *Proceedings of the International Conference on Learning Representations*.

Wang, Y.; Chen, X.; Gao, X.; and Gao, S. 2016. A benchmark dataset for SSVEP-based brain–computer interfaces. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 25(10): 1746–1752.

Yang, C.; Westover, M.; and Sun, J. 2023. Biot: Biosignal transformer for cross-data learning in the wild. *Advances in Neural Information Processing Systems*, 36: 78240–78260.

Zheng, W.; and Lu, B. 2015. Investigating Critical Frequency Bands and Channels for EEG-based Emotion Recognition with Deep Neural Networks. *IEEE Transactions on Autonomous Mental Development*, 7(3): 162–175.

# Appendix

## 1. Pre-training Datasets

1. HBN-EEG dataset (Shirazi et al. 2024) provides high-density (128-channel) EEG from more than 3 000 subjects collected as part of the Healthy Brain Network project. Participants completed six paradigms covering passive conditions—resting state, visual surround-suppression and movie watching—and active tasks of contrast-change detection, sequence learning and symbol search. Recordings were made at 500 Hz (0.1-100 Hz online band-pass, 60 Hz notch) and Cz reference. For further preprocessing we applied band-pass filter 0.5-70 Hz, segment the continuous data into non-overlapping 4 seconds windows with 1 second interval and down-sample to 200 Hz. Following three pre-training datasets are applied similar preprocessing as HBN-EEG.

2. Stieger21 dataset (Stieger, Engel, and He 2021) contains longitudinal EEG from 62 healthy adults who practiced online sensorimotor-rhythm BCI control over 7-11 sessions (598 sessions, >600 h, 269 000 trials). Each session comprised four continuous cursor-control tasks: horizontal (left/right MI), vertical (up/down MI), two-dimensional (combined MI) and rest calibration, with real-time feedback every 40 ms. Signals were recorded with a 64-channel BioSemi cap (10–5 montage) at 1 kHz and hardware-filtered 0.1–200 Hz with a 60 Hz notch.

3. M3CV dataset (Huang et al. 2022) is a large-scale "multi-subject, multi-session, multi-task" resource designed for biometric and variability studies. EEG was recorded from 106 subjects (95 returned for a second visit) while they performed six broad paradigms—resting-state, transient sensory, steady-state sensory, cognitive oddball, motor execution and selective-attention SSVEP—yielding 14 concrete task types and 120 000 labeled epochs. Signals were acquired with 64-channel EasyCap nets at 1 000 Hz (0.1–100 Hz pass-band, 50 Hz notch) with FCz reference and AFz ground.

4. SEED dataset (Zheng and Lu 2015) comprises emotion-elicitation EEG from 15 university students who viewed fifteen 4-min film clips intended to provoke positive, neutral or negative affect in three separate sessions spaced one week apart. EEG (62 channels, NeuroScan cap, 1 000 Hz, 0.05–100 Hz online band-pass, 50 Hz notch) and synchronous eye-tracking were recorded; electrodes were referenced to CPz.

## 2. Downstream Datasets

1. BCIC-IV-2A dataset (Brunner et al. 2008) comprises EEG recordings from nine subjects performing four-class motor imagery (MI) tasks: left hand (Class 1), right hand (Class 2), both feet (Class 3), and tongue (Class 4). Each subject participated in two sessions conducted on separate days, each containing six runs and a total of 288 trials. EEG signals were acquired using 22 Ag/AgCl electrodes arranged according to the international 10–20 system, referenced to the left mastoid, sampled at 250 Hz, and originally bandpass-filtered between 0.5 Hz and 100 Hz with a 50 Hz notch filter to suppress line noise. Additionally, three monopolar EOG channels were recorded to facilitate ocular artifact removal but are not utilised in classification. For preprocessing, we apply a 0.5–38 Hz bandpass filter to retain relevant sensorimotor rhythms while discarding high-frequency noise and slow drifts. Each trial segment is extracted over a 4-second window post cue onset and subsequently downsampled to 200 Hz to reduce computational overhead. We retain all 22 EEG channels for downstream processing. To ensure subject-independent evaluation, we adopt a leave-one-subject-out (LOSO) cross-validation strategy throughout all experiments.

2. BCIC-IV-2B dataset (Leeb et al. 2008) involves binary motor imagery tasks (left hand vs. right hand) performed by nine right-handed subjects. Each subject participated in two initial screening sessions without feedback (120 trials per session), followed by three feedback sessions using a smiley-face interface (80 trials per session). EEG data were recorded from three bipolar electrode channels (C3, Cz, and C4), sampled at 250 Hz, and filtered with a 0.5–100 Hz bandpass filter alongside a 50 Hz notch filter to suppress power line interference. For preprocessing, we apply a bandpass filter in the range of 0.5–38 Hz to retain key motor-related rhythms while attenuating noise. Each trial is segmented into a 4-second epoch post-cue and downsampled to 200 Hz. All three bipolar channels are retained for analysis. We employ a LOSO cross-validation strategy to ensure subject-independent evaluation.

3. PhysioP300 dataset (Citi, Poli, and Cinel 2010) consists of EEG recordings from nine subjects (8, 10, and 12 are removed for a fair comparison with the BENDR and EEGPT methdos) performing a visual P300 speller task based on the classic row-column paradigm. Participants were instructed to attend to specific target characters in a 6x6 matrix, where random sequences of row and column flashes served as stimuli to elicit P300 event-related potentials (ERPs) in response to target cues. The original recordings were acquired using 64-channel EEG caps following the international 10-20 system, downsampled to 250 Hz, and bandpass-filtered between 0.15 Hz and 5 Hz. For our experiments, we retain 58 EEG channels, apply a 120 Hz low-pass filter to reduce high-frequency noise, and resample the data to 200 Hz. Stimulus-locked epochs of 2 seconds duration (-0.7 s to +1.3 s relative to stimulus onset) are extracted for each flash event. We frame the task as a binary classification problem (target vs. non-target) and adopt a LOSO cross-validation scheme for subject-independent evaluation.

4. KaggleERN dataset (Margaux et al. 2012), which contains EEG recordings from 26 healthy participants engaged in a P300-based speller paradigm augmented with online error detection and correction. Subjects were instructed to focus on target letters within a 6x6 flashing matrix and engaged in two spelling modes: fast mode (2 sequences per trial) and slow mode (4 sequences per trial). The task was designed to elicit error-related po-

| | Datasets | Task | Subjects | Channel | Duration | Samples | Classes |
|---|---|---|---|---|---|---|---|
| **Pre-training Datasets** | HBN | Multi-Tasks | 3155 | 128 | 4 | 1497026 | - |
| | Stieger21 | Motor Imagery | 62 | 64 | 4 | 269099 | - |
| | M3CV | Multi-Tasks | 106 | 64 | 4 | 116863 | - |
| | SEED | Emotion Recognition | 15 | 62 | 4 | 30375 | - |
| **Downstream Datasets** | BCIC IV 2A | MI | 9 | 22 | 4 | 5184 | 4 |
| | BCIC IV 2B | MI | 9 | 3 | 4 | 900 | 2 |
| | PhysioP300 | P300 | 10 | 64 | 2 | 2532 | 2 |
| | KaggleERN | ERP | 26 | 56 | 2 | 4448 | 2 |
| | FACED | Emotion Recognition | 123 | 32 | 10 | 10332 | 9 |
| | THUBenchmark | SSVEP | 35 | 64 | 2 | 8400 | 40 |
| | Large-5F | Finger MI | 13 | 38 | 4.5 | 18071 | 5 |
| | BCIC2020-3 | Speech Imagery | 20 | 62 | 3 | 8000 | 5 |
| | TUAB | Clinical | 2383 | 23 | 10 | 409083 | 2 |
| | TUEV | Artifact Detection | 288 | 21 | 5 | 112237 | 6 |

Table 2: Detailed information of the pre-training and downstream datasets

tentials (ErrPs) following incorrect classifier outputs, enabling automatic correction by selecting the classifier's second-best prediction when an error was detected. For our experiments, EEG trials were resampled to 200 Hz, and binary classification (correct vs. incorrect feedback) was performed. We adopted a 4-fold cross-subject validation strategy, where each fold utilized data from 12 subjects for training and the remaining 10 subjects for testing. A total of 19 channels were used in the analysis, and data within the time window of -0.7 s to +1.3 s relative to stimulus onset were cropped and used for classification.

5. FACED dataset (Chen et al. 2023) consists of 32-channel EEG recordings from 123 subjects who viewed 28 video clips designed to induce nine distinct emotional categories: four negative emotions (anger, fear, disgust, sadness), four positive emotions (amusement, inspiration, joy, tenderness), and a neutral condition. The data were recorded at 250 Hz or 1000 Hz using the 10-20 system. Official preprocessing steps involved bandpass filtering between 0.05 and 47 Hz, independent component analysis (ICA) for artifact removal, and re-referencing to a common average reference, resulting in 30 clean channels. For our experiments, the data were segmented into 10-second clips, resampled to 200 Hz, and evaluated using a 4-fold cross-subject validation strategy, where each fold used two-thirds of the subjects for training and the remaining one-third for testing. We used 30 channels from the officially re-referenced dataset for classification.

6. THUBenchmark dataset (Wang et al. 2016) consists of EEG recordings from 35 subjects engaged in a 40-class steady-state visual evoked potential (SSVEP) task. Participants focused on target characters modulated by distinct frequency/phase combinations to elicit class-specific SSVEP responses. EEG signals were recorded using 64 electrodes at a sampling rate of 1000 Hz. For our experiments, we extracted 2-second stimulation windows from each trial, excluding the initial 140 ms visual latency. The signals were bandpass filtered between 0.5 Hz and 45 Hz, downsampled to 200 Hz, and nine occipital-

parietal channels were selected for analysis. Evaluation was conducted using 4-fold cross-subject validation, with each fold comprising 80% of the subjects for training and the remaining 20% for testing.

7. Large-5F dataset (Kaya et al. 2018) comprises EEG recordings from nine subjects performing five-class motor imagery tasks involving individual finger movements: thumb, index, middle, ring, and pinkie. Data were collected as part of a multi-paradigm EEG study using 22 electrodes placed according to the international 10–20 system, with recordings sampled at either 200 Hz or 1000 Hz. A bandpass filter of 0.53–70 Hz (or 100 Hz) and a 50 Hz notch filter were applied during acquisition. Each subject completed up to 75 sessions, resulting in over 60,000 trials across all paradigms. For our experiments, we focused on the five-finger (5F) motor imagery paradigm. Cue-aligned EEG segments were extracted and temporally interpolated to a uniform maximum duration of 4.5 seconds to ensure consistency across trials. The signals were then downsampled to 200 Hz and bandpass filtered between 0.5 Hz and 45 Hz to remove low-frequency drifts and high-frequency noise. All 22 EEG channels were retained, and evaluation was performed using a leave-one-subject-out (LOSO) cross-validation strategy to assess subject-independent performance.

8. BCIC2020-3 dataset (Committee 2022) comprises EEG recordings from 15 healthy subjects aged 20 to 30 years, performing imagined speech tasks involving five phrases: "hello," "help me," "stop," "thank you," and "yes." EEG signals were acquired using 64 electrodes arranged according to the international 10-20 system, with ground electrode at Fpz and reference at FCz. Impedances were maintained below 15 kΩ to ensure signal quality. Each trial consisted of a 2-second auditory cue followed by a 2-second imagery phase. Each phrase class included 70 trials, with 60 trials designated for training and 10 for validation. For experimental evaluation, 3-second epochs were extracted by concatenating the 1-second cue period and the 2-second imagery phase. The data were downsampled to 200 Hz to reduce computational load. All

64 EEG channels were retained. We employed a cross-session validation scheme to assess the model's ability to generalize across recording sessions.

9. TUAB dataset (Obeid and Picone 2016) facilitates binary classification of normal versus abnormal adult EEG recordings based on background brain activity. EEG signals were recorded using 64 electrodes placed according to the international 10-20 system, with ground at Fpz and reference at FCz, maintaining impedances below 15 kΩ. The raw data were sampled at either 200 Hz or 1000 Hz and bandpass filtered between 0.53 Hz and 70 Hz (or 100 Hz) alongside a 50 Hz notch filter to reduce line noise. For our experiments, we selected 23 representative channels, downsampled the signals to 200 Hz, and extracted 10-second segments. Evaluation was performed using 4-fold cross-subject validation to ensure robust generalisation across patients.

10. TUEV dataset (Obeid and Picone 2016), a subset of the Temple University Hospital EEG Corpus, supports classification of six clinically relevant EEG event types: spike and sharp waves (SPSW), periodic lateralized epileptiform discharges (PLED), generalized periodic epileptiform discharges (GPED), artifacts (ARTF), eye movements (EYEM), and background activity (BCKG). The dataset comprises 16,986 sessions from 10,874 subjects, recorded using a variable number of EEG channels (typically 31) at sampling rates of 250, 256, 400, or 512 Hz. For our experiments, we selected 23 standard EEG channels, downsampled the signals to 200 Hz, and segmented the data into 5-second epochs. A 4-fold cross-subject validation scheme was adopted to evaluate model performance across subjects.

## 3. More Details For Experimental Settings

**Data preprocessing**: For **pre-training**, we applied minimum necessary data preprocessing. We first apply a 0.5 HZ to 70 HZ band-pass filter on EEG signals and resample them to 200 HZ. Then we normalize the EEG signals by setting the unit to 0.1 mV to guarantee the value mainly between -1 to 1 (Jiang, Zhao, and Lu 2024). EEG signals are segmented into 4 seconds non-overlap samples according to LaBraM (Jiang, Zhao, and Lu 2024) and EEGPT (Wang et al. 2024). During pre-training, we select 62 EEG channels that are commonly shared across all pre-training datasets. Although CoMET supports varying channel configurations in downstream tasks, using unbalancing number of channels during pre-training can lead to inbalancing channel embeddings, ultimately degrading model performance. In contrast, 62 channel embeddings cover the majority of channel configurations in EEG experiment, and ensure consistent representation for different channels. For **downstream tasks**, we adopt dataset-specific preprocessing strategies to account for structural differences across datasets. Detailed procedures are provided in the dataset description section. In summary, downstream samples vary in terms of filtering frequency bands, sample duration, and the number of EEG channels.

To ensure the reproducibility of our experimental results, we provide more details about the hyperparameters settings on CoMET pre-training in Table 3 and downstream tasks in Table 4.
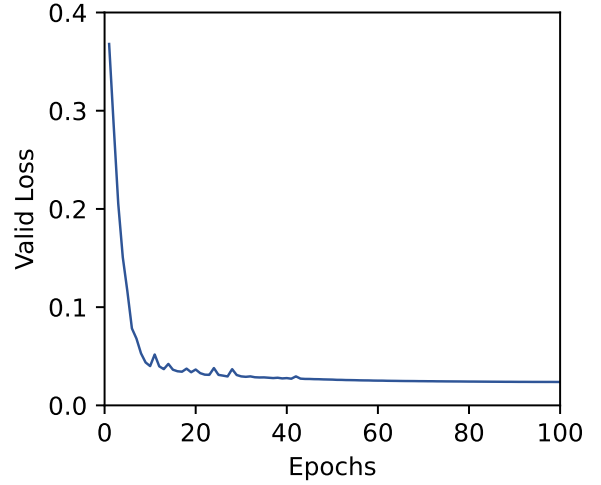
## 4. Pre-training visualization



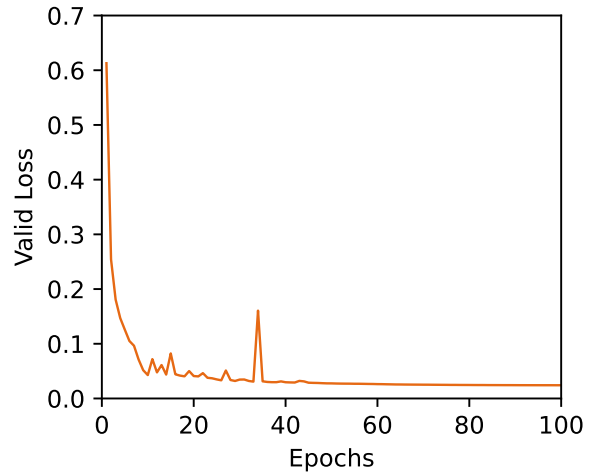Figure 5: The reconstruction loss ($\mathcal{L}_R$) of pre-training



Figure 6: The contrastive loss ($\mathcal{L}_C$) of pre-training

In Figure 5 and Figure 6, we present the reconstruction loss ($\mathcal{L}_R$) curve and contrastive loss ($\mathcal{L}_C$) curve of pre-training. It is evident that both losses decrease rapidly between epochs 1 and 20. When the epoch exceeds 20, both losses tend to stabilize, with the contrast loss exhibiting fluctuations at epoch 35. Throughout pre-training, the two loss functions exhibited no reciprocal constraint, The common downward trajectory indicates that CoMET can effectively acquires local EEG representations through the reconstruction task and global discriminative representations through the contrastive task.

| | Hyperparameters | CoMET-Tiny | CoMET-Base | CoMET-Large |
|---|---|---|---|---|
| Patching Embedding | Input dimension | 1 | 1 | 1 |
| | Output dimension | 256 | 512 | 1024 |
| | Kernel size | | (1, 50) | |
| | Stride | | (1, 50) | |
| Transformer Encoder | Layers | 6 | 6 | 12 |
| | Hidden dimension | 256 | 512 | 1024 |
| | Heads | 4 | 8 | 16 |
| | Feed-forward dimension | 1024 | 2048 | 4096 |
| Transformer Decoder | Layers | 2 | 4 | 6 |
| | Hidden dimension | 384 | 384 | 384 |
| | Heads | 4 | 8 | 16 |
| | Feed-forward dimension | 1536 | 1536 | 1536 |
| Pre-training | Epochs | | 100 | |
| | Batchsize | | 256 | |
| | Peak learning rate | | 5e-4 | |
| | Learning rate scheduler | | CosineAnnealingLR | |
| | Optimizer | | AdamW | |
| | Adam $\beta$ | | (0.9 0.999) | |
| | Weight decay scheduler | | CosineWDSchedule (1e-6) | |
| | Momentum scheduler | | (0.996, 1) | |
| | Temperature | | 0.1 | |
| | Mask ratio | | 0.5 | |

Table 3: Hyperparameters for CoMET Pre-training

| | Hyperparameters | CoMET-Tiny | CoMET-Base | CoMET-Large |
|---|---|---|---|---|
| Linear probing | Input dimension | 256, patch | 512, patch | 1024, patch |
| | Output dimension | 16, classes | 32, classes | 48, classes |
| | Epochs | | 100 | |
| | Batchsize | | 64 | |
| | Peak learning rate | | 1e-3 | |
| | Optimizer | | AdamW | |
| | Adam $\beta$ | | (0.9 0.999) | |

Table 4: Hyperparameters for CoMET Downstream

## 5. Scaling laws

**5.1 Scaling laws with model size**  To analyze how downstream performance scales with model size, we conducted experiments with two additional pre-training settings: a 0.8M-parameter model (4 layers, 128 hidden size, 4 heads) and a 51M-parameter model (12 layers, 512 hidden size, 8 heads), alongside CoMET-Tiny (5M), Base (19M), and Large (151M). We evaluated these five models on the datasets BCIC IV 2A and BCIC2020-3, with the results illustrated in Figure 7. On BCIC IV 2A, the balanced accuracy (BAC) scales with model size (M) as: $BAC = 0.013 \cdot ln(M) + 0.568, (R^2 = 0.949)$. Similarly, on BCIC2020-3, the scaling law is: $BAC = 0.009 \cdot ln(M) + 0.516, (R^2 = 0.989)$.

**5.2 Scaling laws with data size**  We conducted pre-training experiment on CoMET-Base using 20%, 40%, 60%, 80% and 100% of the training data (total 1.91M samples) and tested them on datasets BCIC IV 2A and BCIC2020-3. The results are illustrated in Figure 8. On BCIC IV 2A, the balanced accuracy (BAC) scales with model size (N) as: BAC = $0.0128 \cdot$ ln(N) + 0.567 ($R^2$ = 0.915). Similarly, on BCIC2020-3, the scaling law is: BAC = $0.0231 \cdot$ ln(N) + 0.512 ($R^2$ = 0.933).

## 6. Channel embedding similarity

After pre-training, the foundation models leverage channel embeddings to impose biases on the signals recorded on different channels. Prior work of other large brain models ((Jiang, Zhao, and Lu 2024; Wang et al. 2024, 2025)) has not, however, visualised the properties of these learned embeddings. We present an exploratory analysis of channel-embedding similarity as design instructions for future brain foundation models. Concretely, we visualised the spatial embeddings by (i) performing hierarchical clustering with cosine similarity as the distance metric and (ii) displaying the resulting channel clusters on EEG topographic maps. The visualisations for CoMET-Tiny, CoMET-Base, and CoMET-Large are reported in Figures 9, 10, and 11, respectively.

Across all three models, the embeddings consistently partition into two coarse clusters corresponding to frontal and posterior electrodes. While finer-grained clustering patterns
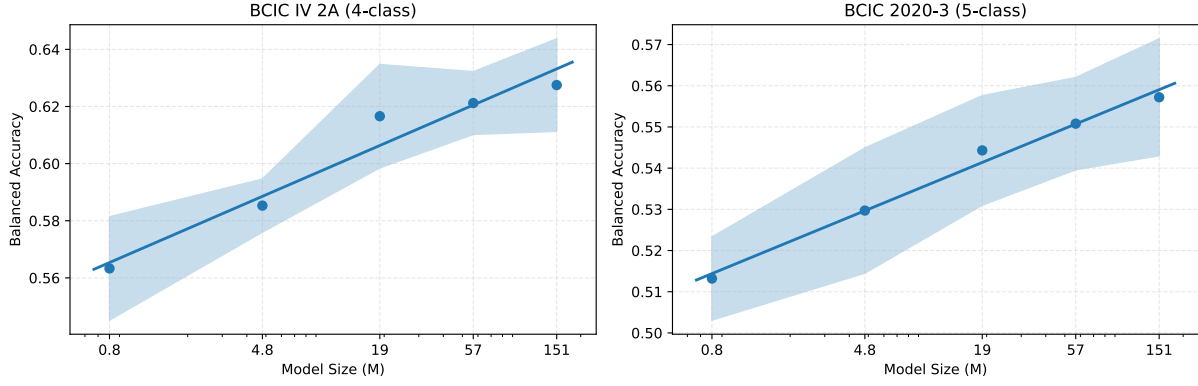
Figure 7: The scaling law with model size $(M)$ and balanced accuracy. BCIC IV 2A: $BAC = 0.013 \cdot ln(M) + 0.568, (R^2 = 0.949)$; BCIC2020-3: $BAC = 0.009 \cdot ln(M) + 0.516, (R^2 = 0.989)$
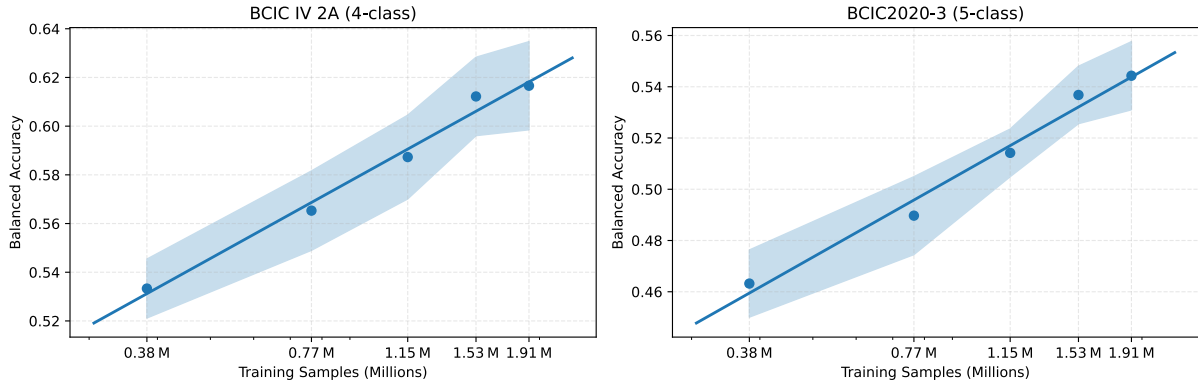


Figure 8: The scaling law with training data size $(N)$ and balanced accuracy. BCIC IV 2A: $BAC = 0.0540 \cdot ln(N) - 0.1635, (R^2 = 0.986)$; BCIC2020-3: $BAC = 0.0524 \cdot ln(N) - 0.2139, (R^2 = 0.982)$
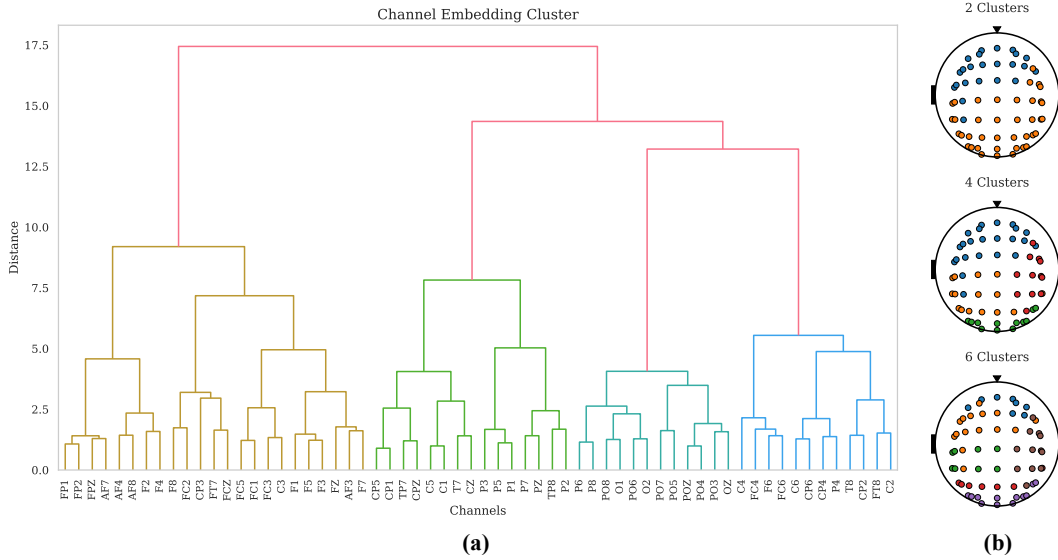


Figure 9: The learned EEG channel embeddings cluster of the model CoMET-Tiny. (a): Dendrogram shhows the cosine similarity of channel embeddings. (b): Visualization in topology of the channel embeddings based on the dendrogram results, with the same color representing the same cluster.
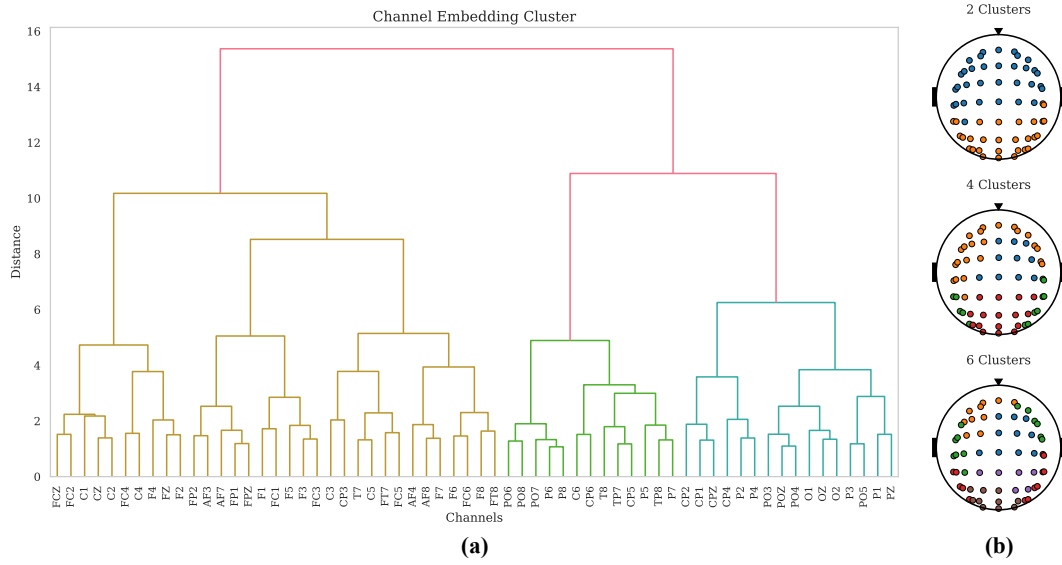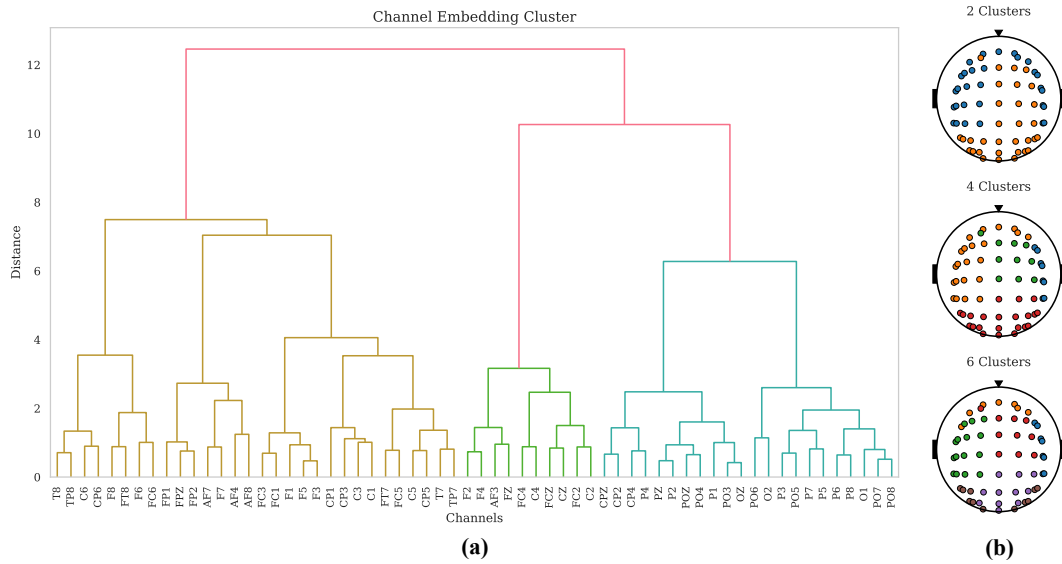
Figure 10: The learned EEG channel embeddings cluster of the model CoMET-Base. (a): Dendrogram shhows the cosine similarity of channel embeddings. (b): Visualization in topology of the channel embeddings based on the dendrogram results, with the same color representing the same cluster.



Figure 11: The learned EEG channel embeddings cluster of the model CoMET-Large. (a): Dendrogram shhows the cosine similarity of channel embeddings. (b): Visualization in topology of the channel embeddings based on the dendrogram results, with the same color representing the same cluster.

differ slightly between model sizes, they broadly respect canonical neuroanatomical regions—namely the frontal, central, occipital, and temporal lobes. These findings indicate that pre-training enables the models to internalise meaningful EEG spatial structure, and they provide embedding-level evidence that neighbouring electrodes tend to share similar feature representations.

## 7. Feature Distributions

To showcase the distribution patterns of feature representations, we performed a t-SNE analysis (Donahue et al. 2014), a widely used nonlinear dimensionality reduction technique, on the learned features across different downstream tasks. This visualization offers an intuitive perspective on how features are organized in the latent space and reveals the clustering behaviour and separability among different classes, as illustrated in Figure 12.

For datasets with high classification accuracy, such as BCIC-IV-2A (4 classes, 62.75%), BCIC-IV-2B (2 classes, 73.22%), THUBenchmark (40 classes, 92.74%), and TUAB (2 classes, 82.02%), we observe that inter-class feature distributions exhibit a high degree of separability. In contrast, for datasets with relatively lower classification accuracy, including KaggleERN (2 classes, 58.78%), Large-5F (5 classes, 38.97%), and FACED (9 classes, 39.02%), the feature representations show significant overlap among different classes. This observation is consistent with the corresponding classification accuracies, indicating that the degree of feature separability across classes aligns well with the overall recognition performance.

## 8. Brain map of attention

To facilitate a comprehensive discussion of the model's performance in the EEG paradigm, we present the attention weight distributions across different channels in the BCIC IV 2A dataset (22-channel input) under varying query conditions, as shown in Figure 13, Figure 14, and Figure 15. As observed, when the input query corresponds to different channels, BIOT's attention distribution is restricted to a limited set of paradigms, which exhibit distinct brain region activation patterns in relation to different cognitive tasks. Similarly, LaBram's attention remains confined to either the current or neighboring channels for most queries, thereby lacking extensive coverage of the brain's diverse regions. In contrast, CoMET demonstrates a much broader attention distribution across channels, regardless of the query input. This characteristic enables CoMET to maintain a high degree of flexibility across varying downstream task paradigms without the need for retraining the encoder, thus positioning it as a highly efficient and generalizable foundational model.

## 9. Different downstream strategies

In Table 5, Table 6, Table 7, Table 8, Table 9, and Table 10, we present different models alongside CoMET performance under different strategies, namely linear-probing (lb) and fine-tuning (ft), where the former freezes the encoder while the latter fully trains the encoder in the downstream tasks. Some models show significant declines when using strategies different from those in their original paper, such as

EEGPT when using the ft strategy and CbraMod when using the lb strategy. In contrast, when using different strategies, the balanced accuracy difference of CoMET is less than 6%. The lb strategy is generally considered to represent the actual performance of the pre-trained model since the model is not affected by downstream tasks, and CoMET's lb performance is superior than that of other models.
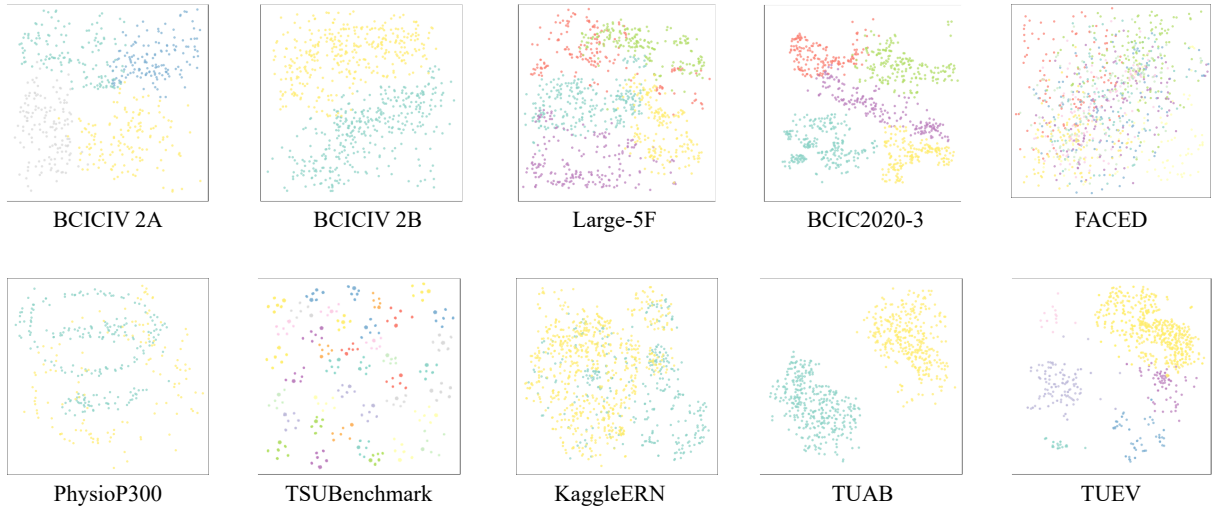
Figure 12: T-SNE visualizations of feature distributions on different datasets with CoMET-Large.
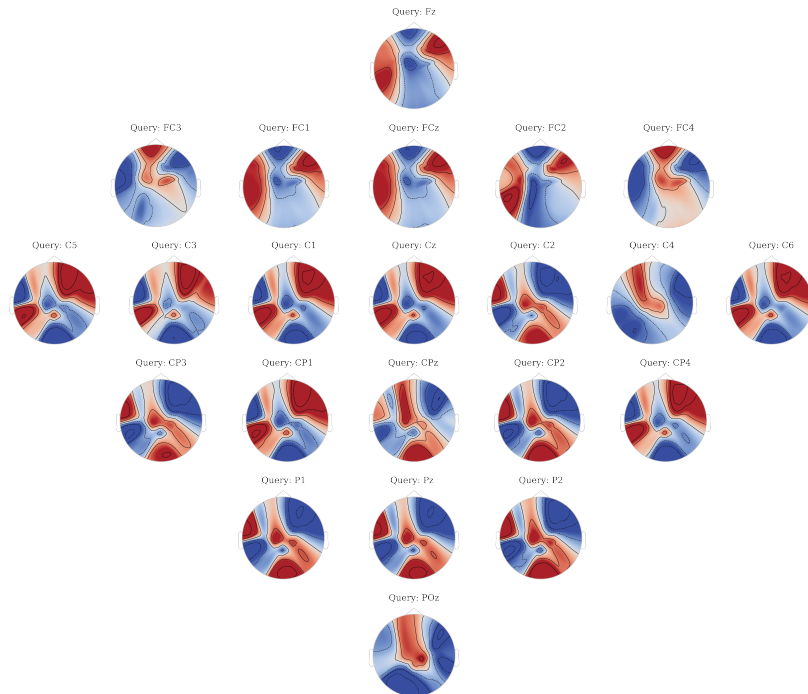


Figure 13: BIOT's attention map for 22 EEG channels (BCIC IV 2A); warmer colors mark stronger attention levels.
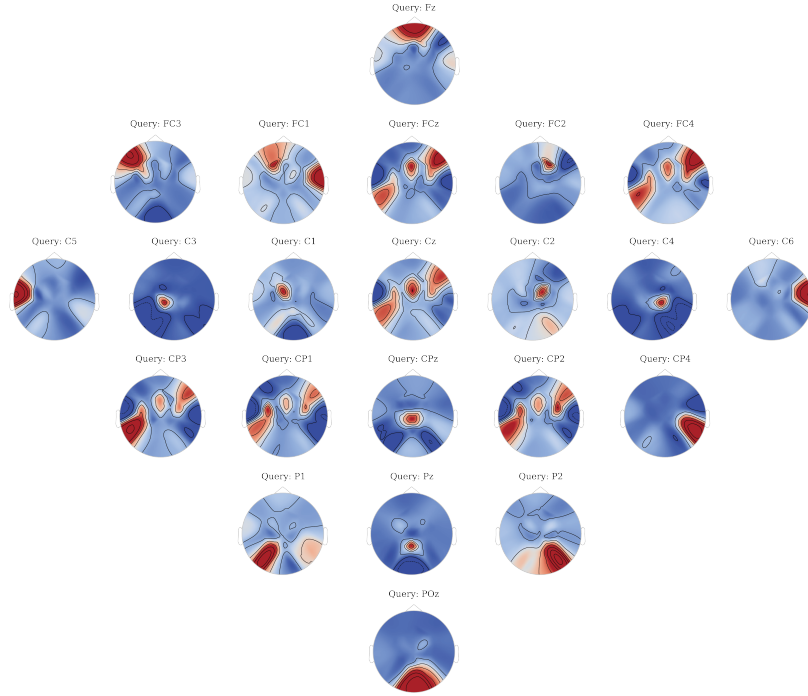
Figure 14: LaBraM's attention map for 22 EEG channels (BCIC IV 2A); warmer colors mark stronger attention levels.
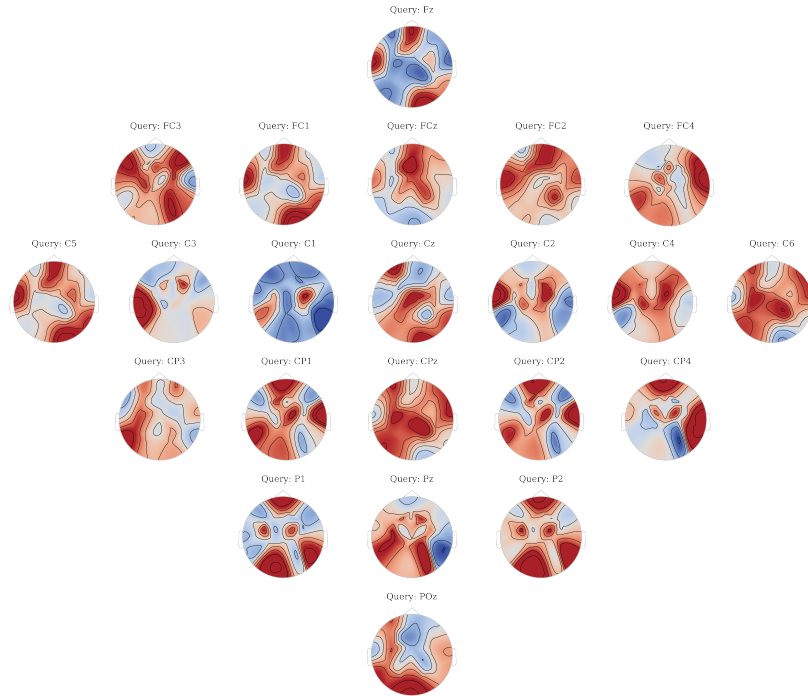


Figure 15: CoMET's attention map for 22 EEG channels (BCIC IV 2A); warmer colors mark stronger attention levels.

| Model | BCIC IV 2A | | | |
|---|---|---|---|---|
| | Strategy | B. Acc | Kappa | F1 |
| BIOT | ft | 41.44 ± 0.58 | 21.90 ± 1.19 | 37.37 ± 0.56 |
| | lb | 42.02 ± 1.17 | 22.63 ± 1.56 | 38.10 ± 1.28 |
| LaBraM | ft | 52.49 ± 1.34 | 36.60 ± 2.12 | 52.32 ± 0.97 |
| | lb | 51.37 ± 1.12 | 35.17 ± 1.48 | 49.73 ± 1.30 |
| EEGPT | ft | 38.35 ± 0.31 | 17.69 ± 0.41 | 37.97 ± 0.31 |
| | lb | 51.37 ± 0.96 | 35.17 ± 1.26 | 49.73 ± 0.41 |
| CBraMod | ft | 55.85 ± 0.97 | 41.13 ± 1.30 | 55.08 ± 1.02 |
| | lb | 29.04 ± 0.32 | 5.40 ± 0.42 | 16.69 ± 0.46 |
| CoMET | ft | 57.68 ± 0.70 | 38.25 ± 0.93 | 52.81 ± 0.72 |
| | lb | **62.75 ± 1.62** | **51.70 ± 1.84** | **63.37 ± 1.34** |

Table 5: Model performance comparisons across different strategies on the BCIC IV 2A Dataset

| Model | BCIC IV 2B | | | |
|---|---|---|---|---|
| | Strategy | B. Acc | Kappa | F1 |
| BIOT | ft | 50.04 ± 0.12 | 0.09 ± 0.25 | 6.31 ± 0.88 |
| | lb | 49.79 ± 4.13 | -0.33 ± 8.22 | 44.45 ± 19.62 |
| LaBraM | ft | 54.26 ± 0.22 | 8.38 ± 0.43 | 56.86 ± 1.38 |
| | lb | 63.93 ± 7.83 | 27.64 ± 15.81 | 61.19 ± 9.96 |
| EEGPT | ft | 59.92 ± 7.63 | 19.75 ± 15.40 | 57.45 ± 8.86 |
| | lb | 61.07 ± 0.27 | 22.15 ± 0.52 | 54.24 ± 0.48 |
| CBraMod | ft | 59.46 ± 0.01 | 18.85 ± 0.03 | 55.68 ± 0.03 |
| | lb | 50.82 ± 1.86 | 1.68 ± 3.71 | 40.83 ± 8.77 |
| CoMET | ft | **63.95 ± 0.84** | **27.76 ± 1.68** | **59.33 ± 1.47** |
| | lb | 63.86 ± 1.78 | 25.05 ± 2.01 | 61.72 ± 2.28 |

Table 6: Model performance comparisons across different strategies on the BCIC IV 2B Dataset

| Model | KaggleERN | | | |
|---|---|---|---|---|
| | Strategy | B. Acc | Kappa | F1 |
| BIOT | ft | 50.64 ± 0.01 | 1.62 ± 0.03 | 71.53 ± 0.01 |
| | lb | 52.80 ± 2.06 | 6.13 ± 4.19 | 78.11 ± 4.76 |
| LaBraM | ft | 51.95 ± 0.81 | 4.38 ± 1.75 | 77.21 ± 1.82 |
| | lb | 57.54 ± 2.55 | 17.21 ± 4.47 | 81.10 ± 1.91 |
| EEGPT | ft | 53.81 ± 1.87 | 8.85 ± 3.81 | 79.90 ± 1.69 |
| | lb | 54.92 ± 0.04 | 11.89 ± 0.12 | 76.77 ± 0.04 |
| CBraMod | ft | 53.92 ± 0.02 | 8.15 ± 0.05 | 76.49 ± 0.04 |
| | lb | 50.57 ± 0.55 | 1.51 ± 1.40 | 60.27 ± 0.98 |
| CoMET | ft | 58.18 ± 1.45 | 17.23 ± 2.63 | 77.93 ± 2.22 |
| | lb | **58.78 ± 0.88** | **15.24 ± 0.92** | **79.66 ± 1.43** |

Table 7: Model performance comparisons across different strategies on the Kaggle ERN Dataset

| Model | FACED | | | |
|-------|-------|-------|-------|-------|
| | Strategy | B. Acc | Kappa | F1 |
| BIOT | ft | 50.64 ± 0.01 | 1.62 ± 0.03 | 71.53 ± 0.01 |
| | lb | 13.56 ± 0.40 | 2.87 ± 0.46 | 12.28 ± 0.82 |
| LaBraM | ft | 51.95 ± 0.81 | 4.38 ± 1.75 | 77.21 ± 1.82 |
| | lb | 21.80 ± 1.05 | 12.09 ± 1.23 | 21.57 ± 1.14 |
| EEGPT | ft | 15.16 ± 0.06 | 4.67 ± 0.07 | 15.35 ± 0.06 |
| | lb | 54.92 ± 0.04 | 11.89 ± 0.12 | 76.77 ± 0.04 |
| CBraMod | ft | 53.92 ± 0.02 | 8.15 ± 0.05 | 76.49 ± 0.04 |
| | lb | 29.06 ± 1.18 | 19.88 ± 1.25 | 27.76 ± 0.73 |
| CoMET | ft | 40.10 ± 2.25 | 32.47 ± 2.52 | 40.13 ± 2.15 |
| | lb | **58.78 ± 0.88** | **15.24 ± 0.92** | **79.66 ± 1.43** |

Table 8: Model performance comparisons across different strategies on the FACED Dataset

| Model | THUBenchmark | | | |
|-------|-------|-------|-------|-------|
| | Strategy | B. Acc | Kappa | F1 |
| BIOT | ft | 74.17 ± 0.06 | 73.50 ± 0.07 | 74.24 ± 0.05 |
| | lb | 80.07 ± 0.89 | 78.54 ± 0.12 | 80.20 ± 0.39 |
| LaBraM | ft | 91.43 ± 0.13 | 91.21 ± 0.13 | 91.43 ± 0.13 |
| | lb | 92.02 ± 0.92 | 94.79 ± 0.45 | 95.02 ± 0.55 |
| EEGPT | ft | 30.23 ± 0.03 | 28.44 ± 0.03 | 30.87 ± 0.03 |
| | lb | 82.59 ± 0.09 | 82.15 ± 0.09 | 82.57 ± 0.08 |
| CBraMod | ft | 91.45 ± 0.24 | 91.23 ± 0.25 | 91.43 ± 0.24 |
| | lb | 33.04 ± 0.92 | 27.88 ± 0.90 | 32.19 ± 0.18 |
| CoMET | ft | **96.88 ± 0.94** | **96.63 ± 0.11** | **96.86 ± 0.14** |
| | lb | 92.74 ± 1.61 | 92.65 ± 2.46 | 93.50 ± 0.81 |

Table 9: Model performance comparisons across different strategies on the THUBenchmark dataset

| Model | PhysioP300 | | | |
|-------|-------|-------|-------|-------|
| | Strategy | B. Acc | Kappa | F1 |
| BIOT | ft | 50.04 ± 0.12 | 0.09 ± 0.25 | 6.31 ± 0.88 |
| | lb | 49.79 ± 4.13 | -0.33 ± 8.22 | 44.45 ± 19.62 |
| LaBraM | ft | 54.26 ± 0.22 | 8.38 ± 0.43 | 56.86 ± 1.38 |
| | lb | 61.93 ± 7.83 | 27.64 ± 15.81 | 61.19 ± 9.96 |
| EEGPT | ft | 59.92 ± 7.63 | 19.75 ± 15.40 | 57.45 ± 8.86 |
| | lb | 61.07 ± 0.27 | 22.15 ± 0.52 | 54.24 ± 0.48 |
| CBraMod | ft | 59.46 ± 0.01 | 18.85 ± 0.03 | 55.68 ± 0.03 |
| | lb | 50.82 ± 1.86 | 1.68 ± 3.71 | 40.83 ± 8.77 |
| CoMET | ft | **63.95 ± 0.84** | **27.76 ± 1.68** | **59.33 ± 1.47** |
| | lb | 63.86 ± 1.78 | 25.05 ± 2.01 | 61.72 ± 2.28 |

Table 10: Model performance comparisons across different strategies on the PhysioP300 Dataset