

DO WE NEED DOMAIN-SPECIFIC TIME-SERIES MODELS? INSIGHTS FROM EEG CLASSIFICATION BENCHMARKS

Liuyin Yang*, Qiang Sun*, & Marc M. Van Hulle

Laboratory for Neuro- and Psychophysiology, Department of Neurosciences
 KU Leuven, Belgium
 {liuyin.yang, qiang.sun, marc.vanhulle}@kuleuven.be

ABSTRACT

Generic time-series foundation models (TS FMs) are now trained on large-scale collections of heterogeneous time series, but whether this broad pre-training actually helps in specialized domains remains an open question. We take EEG classification as a case study, benchmarking two generic TS FMs—Mantis (~8M parameters) and MOMENT (>300M)—on four public EEG datasets (TUEV, FACED, BCI-IV-2A, and Error) under linear probing and fine-tuning, and comparing them against EEG foundation models and classic neural baselines. Fine-tuning often matches EEG-specific models, while linear probing transfers poorly. Interestingly, randomly initialized Mantis performs comparably to its pre-trained version, suggesting that its architecture, rather than pre-training, may be driving much of its performance. These results illustrate both the promise and the limits of generic TS FMs for specialized domains.

Track: Research

1 INTRODUCTION

Electroencephalography (EEG) is a non-invasive brain recording technique widely used in brain-computer interfaces (BCIs). Yet EEG datasets are often small, differ in channel layouts, and exhibit strong subject/session variability. These factors limit the reliability of conventional deep learning models and frequently require task-specific training.

Motivated by the success of large-scale pre-training in vision and language, recent work has explored EEG FMs. Approaches such as BENDR (Kostas et al., 2021), BIOT (Yang et al., 2023), and LaBraM (Jiang et al., 2024) show that self-supervised pre-training can help, although gains over classical decoders are sometimes modest. Most EEG-specific models design pre-training tasks tailored to neural signals (e.g., reconstructing spectral structure in LaBraM), or with architectures adapted to multichannel inputs (e.g., spatial CNN in CBraMod (Wang et al., 2025)).

In parallel, a growing family of generic TS FMs, such as TimesFM (Das et al., 2024), Chronos (Ansari et al., 2024), UniTS (Gao et al., 2024), Mantis (Feofanov et al., 2025), MOMENT (Goswami et al., 2024), have demonstrated strong cross-domain results on forecasting and classification. Some of these models appear to transfer to EEG as well: Mantis, pre-trained mostly on non-neural time series, was shown to outperform EEGNet and CBraMod on motor imagery and sleep staging (Gnassounou et al., 2025). However, direct comparisons with EEG-specific foundation models on challenging BCI tasks are still rare. Recent EEG FM benchmarks (Yang et al., 2026; Xiong et al., 2026; Liu et al., 2026) report that fine-tuned foundation models work well in population-level settings but show diminishing returns under leave-one-out or low-data regimes, with linear probing consistently weak. An important question remains: when a generic TS FM does well on EEG, is it because of what it learned during pre-training, or because of its architecture? These open questions motivate a closer look at how well generalist models actually transfer to EEG.

*Equal contribution and shared correspondence.

In this work, we evaluate Mantis and MOMENT on four public EEG datasets: TUEV (6-class event classification) (Obeid & Picone, 2016), FACED (9-class emotion recognition) (Chen et al., 2023), BCI-IV-2A (4-class motor imagery) (Tangermann et al., 2012), and Error (error-related negativity) (Kueper et al., 2024). Following recent EEG FM benchmarks, we test cross-subject decoding on TUEV and FACED and compare population-level decoding with leave-one-out fine-tuning on BCI-IV-2A and Error. Complementing EEG-specific FM benchmarks (Yang et al., 2026), we focus on cross-domain transfer from generic TS FMs and include a random-initialization ablation to separate the effects of architecture and pre-training. We ask: (1) how well do generic TS FMs transfer to EEG, (2) is linear probing sufficient, (3) when does fine-tuning provide consistent gains, and (4) does Mantis’s competitiveness stem from its pre-trained representations or its architecture?

2 METHODOLOGY

2.1 BENCHMARK MODELS

Classic EEG Neural Network Models We benchmarked two well-established convolutional neural network (CNN)-based EEG decoders: DeepConvNet (Schirrmester et al., 2017) and EEGNet (Lawhern et al., 2018). Additionally, we included more recent, transformer-based architectures, EEG Conformer (Song et al., 2023) and CTNet (Zhao et al., 2024), which have demonstrated state-of-the-art performance yet remain computationally simpler compared to larger foundation models.

EEG Foundation Models For all benchmark datasets, we evaluated a range of EEG FMs, including BENDR (Kostas et al., 2021), BIOT (Yang et al., 2023), LaBraM (Jiang et al., 2024), EEGPT (Wang et al., 2024), CBraMod (Wang et al., 2025), and ST-EEGFormer (Yang et al., 2026). ST-EEGFormer is available in three capacity variants: small (ST-EEGFormer-s), base (ST-EEGFormer-b), and large (ST-EEGFormer-l), which share the same spatio-temporal transformer design but differ in width/depth and parameter count to trade off accuracy and compute. For all foundation models, we evaluated both linear probing and fine-tuning performance.

Time-Series Foundation Models We benchmarked Mantis (Feofanov et al., 2025), a lightweight contrastive model (~ 8 M parameters) pre-trained on ~ 7 M time series with a small fraction of EEG data (SleepEEG and Epilepsy), and MOMENT (Goswami et al., 2024), a large masked-reconstruction transformer (> 300 M parameters). We additionally evaluated a randomly initialized Mantis (same architecture, no pre-trained weights) to disentangle the contributions of architecture and pre-training. A summary of model parameters can be found in Table 1.

2.2 EXPERIMENTAL SETUP

We consider three evaluation settings. **Population decoding** (BCI-IV-2A, Error) pools training data from all subjects to learn a single model, which is then evaluated separately on each subject. **Leave-one-subject-out** (LOO; BCI-IV-2A, Error) trains the population model on all but one held-out subject and reports its performance on the held-out subject (LOO zero-shot). We then fine-tune the same model on the held-out subject and re-evaluate it (LOO fine-tune); the change in performance on the non-held-out subjects after this adaptation measures the **generalization drop**. **Cross-subject** evaluation (TUEV, FACED) follows the standard zero-shot protocol used in prior work.

For all tasks, we report top-1 accuracy (Acc1), top-2 accuracy (Acc2), balanced

Table 1: Model parameter counts (BCI-IV-2A).

Model	Scratch	Fine-tune	Linear-prob
<i>Classic neural networks</i>			
EEGNet	5.8K	–	–
DeepConvNet	0.28M	–	–
Conformer	1.5M	–	–
CTNet	0.15M	–	–
<i>EEG foundation models</i>			
BENDR	–	3.9M	22.9K
BIOT	–	3.1M	1.4K
LaBraM	–	5.8M	0.8K
EEGPT	–	25.3M	34.3K
CBraMod	–	19.1M	14.2M
ST-EEGFormer-s	–	25.4M	2.0K
ST-EEGFormer-b	–	85.3M	3.0K
ST-EEGFormer-l	–	302.7M	4.1K
<i>Time-series foundation models</i>			
MOMENT	–	341.2M	4.1K
Mantis	–	8.1M	22.5K

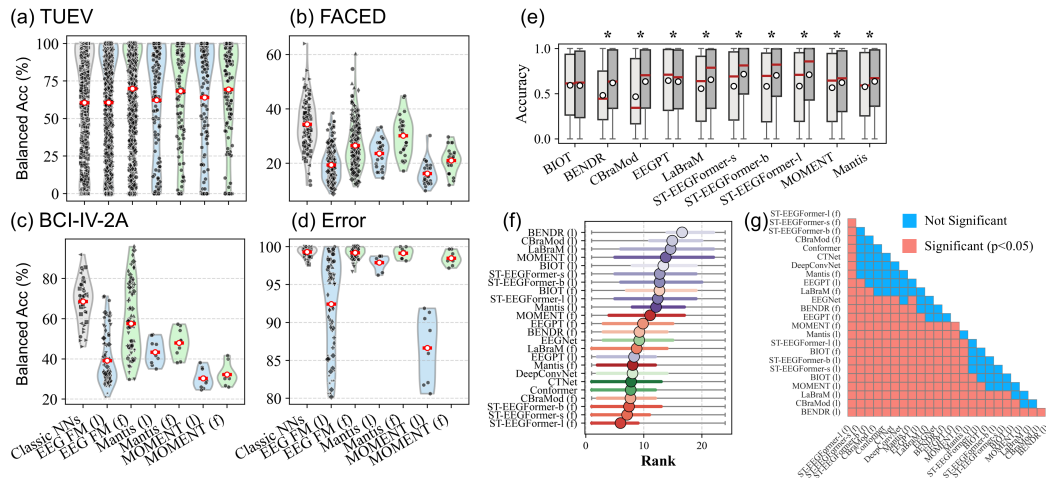


Figure 1: **Per-dataset balanced-accuracy breakdown and cross-model comparisons.** (a–d) Per-subject balanced accuracy on TUEV, FACED, BCI-IV-2A, and Error. TUEV and FACED use cross-subject evaluation, while BCI-IV-2A and Error report population decoding results. Models are grouped into classic neural-network decoders and EEG FMs, where (f) and (l) denote fine-tuning and linear probing, respectively; we compare against the generalist TS FMs Mantis and MOMENT under the same training regimes. Each dot corresponds to a subject and the red dot denotes the mean. (e) Aggregated accuracy across foundation models under linear probing vs. fine-tuning, where * indicates significance using the Wilcoxon signed-rank test with $p < 0.05$.

accuracy (BAcc), AUC, and Cohen’s κ . A more detailed description of the training setup can be found in the Appendix.

3 RESULTS AND DISCUSSION

We summarize the main findings in Figure 1 and Figure 2. Full metric tables are provided in the Appendix. We organize the discussion around four questions.

(1) How well do generic TS FMs transfer to EEG? Across datasets and evaluation protocols, the TS FMs are competitive, but they do not generally dominate. In particular, fine-tuned Mantis consistently performs in the upper tier, whereas fine-tuned MOMENT tends to lag behind most fine-tuned EEG FMs. The rank-based summary and significance testing in Figure 1(e) further support this: Mantis is among the stronger generalist baselines after fine-tuning, while MOMENT is often significantly worse than the best-performing EEG FMs.

(2) Is linear probing sufficient? Linear probing is consistently weaker than end-to-end fine-tuning for both EEG FMs and TS FMs. Figure 1(e) shows a clear aggregate advantage for fine-tuning over linear probing, indicating that the representations learned during pre-training are not directly applicable to downstream EEG tasks and typically require task-specific adaptation.

(3) When does fine-tuning provide consistent gains? Fine-tuning generally improves performance on the target subject/task, but the magnitude of the gain depends on the dataset and evaluation setting. In population or cross-subject decoding, fine-tuned foundation models can match or exceed classic neural network decoders trained from scratch on some datasets. For example, Mantis and MOMENT have higher average performance than classic neural network decoders on TUEV, but not on the others, especially on BCI-IV-2A, where both Mantis and MOMENT show much lower decoding performance than classic neural network decoders and the average of the EEG FMs. The LOO analysis in Figure 2 highlights an additional trade-off: adapting to the held-out subject can induce a generalization drop on the previously seen training subjects. Compared with scratch-trained

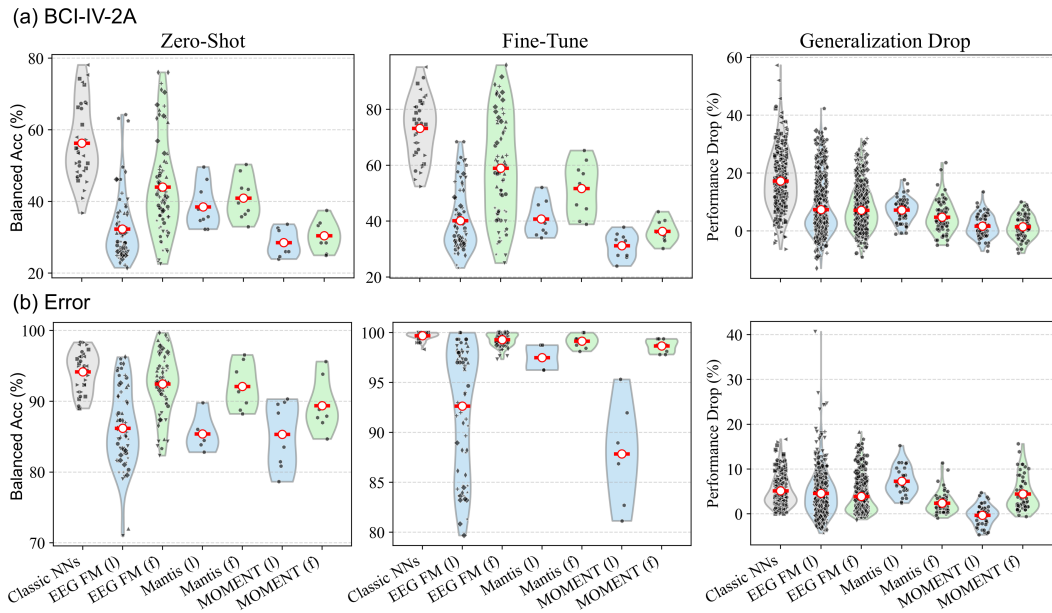


Figure 2: **Leave-one-subject-out fine-tuning analysis.** (a) BCI-IV-2A leave-one-subject-out (LOO) fine-tuning results. (b) Error dataset LOO fine-tuning results. From left to right, we show the zero-shot decoding results on the held-out (LOO) subject, the performance after fine-tuning on that held-out subject, and the generalization drop on the seen training subjects (before vs. after fine-tuning).

models, foundation models often show smaller drops, suggesting that pre-training yields representations that are more stable under subject-specific adaptation (i.e., less catastrophic forgetting).

(4) Comparing Mantis and MOMENT, and the role of pre-training. Despite having far fewer parameters, Mantis outperforms MOMENT under both linear probing and fine-tuning in our benchmarks. To examine the role of pre-training more directly, we fine-tuned a randomly initialized Mantis with the same architecture on all four datasets. The pattern is informative. On BCI-IV-2A, the randomly initialized model outperforms the pre-trained model in population decoding (balanced accuracy $p < 0.01$, AUC $p < 0.05$, paired Wilcoxon test), whereas on TUEV and FACED the two perform similarly ($p > 0.5$). By contrast, on Error, pre-training provides a modest but significant advantage in the leave-one-out fine-tuning setting (Cohen’s κ 0.971 vs. 0.936, $p < 0.01$). Overall, these findings suggest that Mantis’s EEG performance is driven more by its architecture than by its pre-trained representations. Consistent with this interpretation, CTNet, a much smaller model with only 0.15M parameters trained from scratch, still outperforms both TS FMs on BCI-IV-2A. This observation underscores that model scale alone does not guarantee better transfer when the architecture is not well matched to the task. Full results are reported in Table 12 in the Appendix.

4 CONCLUSION

This work evaluated whether generic TS FMs can transfer to EEG classification across four public datasets and multiple evaluation protocols. Overall, our results indicate that using TS FMs for EEG decoding is feasible: after end-to-end adaptation, models such as Mantis can reach competitive performance and sometimes match or surpass EEG-specific baselines.

However, linear probing is consistently weak across datasets, suggesting that off-the-shelf TS FM representations are not directly applicable to downstream EEG tasks and typically require task-specific adaptation. Moreover, even after fine-tuning, TS FMs are often not the top-performing approach and can lag behind small, task-trained neural networks, highlighting a remaining gap between current generic TS FMs and the requirements of robust EEG decoding.

Taken together, these findings suggest that domain-specific EEG FMs remain important in the near term. From a representation-learning perspective, this is plausible: generic TS FMs are trained to capture structure that transfers across many time-series domains, whereas EEG decoding may additionally require representations that are more closely aligned with domain-specific signal characteristics. The random-initialization results are consistent with this interpretation: in our benchmarks, architectural choices can matter as much as, or more than, generic pre-training, and the benefits of pre-training are not uniform across datasets. At the same time, the strong performance of Mantis—even without pre-training—suggests that lightweight architectures provide a promising starting point for EEG decoding.

Limitations and future work. This study has several limitations. First, we evaluate only two generic TS FMs, so it remains unclear how broadly our conclusions extend to other general-purpose time-series models such as TimesFM, Chronos, Timer, or UniTS. Second, our analysis is primarily performance-based and does not directly compare the representations learned by generic TS FMs and EEG-specific FMs. A more detailed representational analysis could help clarify whether the observed performance gap reflects differences in the features encoded before adaptation or differences in how readily those features can be adapted to EEG tasks.

Several directions follow from this. An immediate next step is to expand the benchmark to a broader set of generic TS FMs and EEG datasets, including settings with stronger domain shift and lower data availability. It would also be valuable to characterize learned representations more directly, for example through probing analyses, cross-model similarity measures, or frequency-aware evaluations that test sensitivity to EEG-relevant signal structure. Another promising direction is to investigate intermediate adaptation strategies, such as continued pre-training on unlabeled EEG, which may help bridge the gap between generic pre-training and domain-specific decoding. More broadly, an important goal is to develop stronger EEG representations that can consistently outperform compact task-specific models such as CTNet, rather than merely approaching their performance with substantially larger models.

ACKNOWLEDGMENTS

*L.Y. is supported by the Research Foundation – Flanders (FWO) grant 1S65622N.

*Q.S. is supported by the China Scholarship Council (no. 202206050022).

*M.M.V.H. is supported by research grants received from Horizon Europe’s Marie Skłodowska-Curie Action (grant agreement No. 101118964), Horizon 2020 research and innovation programme under grant agreement No. 857375, the special research fund of the KU Leuven (C24/18/098), the Belgian Fund for Scientific Research – Flanders (G0A4321N, G0C1522N, G031426N), and the Hercules Foundation (AKUL 043).

*The resources and services used in this work were provided by the VSC (Flemish Supercomputer Center), funded by the Research Foundation - Flanders (FWO) and the Flemish Government.

REFERENCES

- Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al. Chronos: Learning the language of time series. In *Advances in Neural Information Processing Systems*, 2024.
- Jingjing Chen, Xiaobin Wang, Chen Huang, Xin Hu, Xinke Shen, and Dan Zhang. A large finer-grained affective computing EEG dataset. *Sci. Data*, 10(1):740, October 2023.
- Abhimanyu Das, Weihao Kong, Andrew Leber, Rajat Mathews, and Ravikumar Sen. A decoder-only foundation model for time-series forecasting. In *Proceedings of the 41st International Conference on Machine Learning*, 2024.
- Vasilii Feofanov, Songkang Wen, Marius Alonso, Romain Ilbert, Hongbo Guo, Malik Tiomoko, Lujia Pan, Jianfeng Zhang, and Ievgen Redko. Mantis: Lightweight calibrated foundation model for user-friendly time series classification, 2025. URL <https://arxiv.org/abs/2502.15637>.
- Shanghua Gao, Teddy Koker, Owen Queen, Thomas Hartvigsen, Theodoros Tsiligkaridis, and Marinka Zitnik. Units: A unified multi-task time series model. In *Advances in Neural Information Processing Systems*, 2024.
- Théo Gnassounou, Yessin Moakher, Shifeng Xie, Vasilii Feofanov, and Ievgen Redko. Leveraging generic time series foundation models for eeg classification, 2025. URL <https://arxiv.org/abs/2510.27522>.
- Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski. Moment: A family of open time-series foundation models, 2024. URL <https://arxiv.org/abs/2402.03885>.
- Wei-Bang Jiang, Li-Ming Zhao, and Bao-Liang Lu. Large brain model for learning generic representations with tremendous eeg data in bci, 2024. URL <https://arxiv.org/abs/2405.18765>.
- Demetres Kostas, Stephane Aroca-Ouellette, and Frank Rudzicz. Bendr: using transformers and a contrastive self-supervised learning task to learn from massive amounts of eeg data, 2021. URL <https://arxiv.org/abs/2101.12037>.
- Niklas Kueper, Kartik Chari, Judith Bütefür, Julia Habenicht, Tobias Rossol, Su Kyoung Kim, Marc Tabie, Frank Kirchner, and Elsa Andrea Kirchner. EEG and EMG dataset for the detection of errors introduced by an active orthosis device. *Front. Hum. Neurosci.*, 18:1304311, January 2024.
- Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and Brent J Lance. EEGNet: a compact convolutional neural network for EEG-based brain-computer interfaces. *J. Neural Eng.*, 15(5):056013, October 2018.
- Dingkun Liu, Yuheng Chen, Zhu Chen, Zhenyao Cui, Yaozhi Wen, Jiayu An, Jingwei Luo, and Dongrui Wu. Eeg foundation models: Progresses, benchmarking, and open problems, 2026. URL <https://arxiv.org/abs/2601.17883>.
- Iyad Obeid and Joseph Picone. The temple university hospital EEG data corpus. *Front. Neurosci.*, 10:196, May 2016.
- Robin Tibor Schirrmester, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin Glasstetter, Katharina Eggensperger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. Deep learning with convolutional neural networks for EEG decoding and visualization. *Hum. Brain Mapp.*, 38(11):5391–5420, November 2017.
- Yonghao Song, Qingqing Zheng, Bingchuan Liu, and Xiaorong Gao. EEG conformer: Convolutional transformer for EEG decoding and visualization. *IEEE Trans. Neural Syst. Rehabil. Eng.*, 31:710–719, February 2023.

- Michael Tangermann, Klaus-Robert Müller, Ad Aertsen, Niels Birbaumer, Christoph Braun, Clemens Brunner, Robert Leeb, Carsten Mehring, Kai J Miller, Gernot R Müller-Putz, Guido Nolte, Gert Pfurtscheller, Hubert Preissl, Gerwin Schalk, Alois Schlögl, Carmen Vidaurre, Stephan Waldert, and Benjamin Blankertz. Review of the BCI competition IV. *Front. Neurosci.*, 6:55, July 2012.
- Guagnyu Wang, Wenchao Liu, Yuhong He, Cong Xu, Lin Ma, and Haifeng Li. Eegpt: Pre-trained transformer for universal and reliable representation of eeg signals. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 39249–39280. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/4540d267ehec4e5dbd9dae9448f0b739-Paper-Conference.pdf.
- Jiquan Wang, Sha Zhao, Zhiling Luo, Yangxuan Zhou, Haiteng Jiang, Shijian Li, Tao Li, and Gang Pan. Cbramod: A criss-cross brain foundation model for eeg decoding, 2025. URL <https://arxiv.org/abs/2412.07236>.
- Wei Xiong, Jiangtong Li, Jie Li, Kun Zhu, and Changjun Jiang. Eeg-fm-bench: A comprehensive benchmark for the systematic evaluation of eeg foundation models, 2026. URL <https://arxiv.org/abs/2508.17742>.
- Chaoqi Yang, M. Brandon Westover, and Jimeng Sun. Biot: Cross-data biosignal learning in the wild, 2023. URL <https://arxiv.org/abs/2305.10351>.
- Liuyin Yang, Qiang Sun, Ang Li, and Marc M. Van Hulle. Are EEG foundation models worth it? comparative evaluation with traditional decoders in diverse BCI tasks. In *The Fourteenth International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=5Xwm8e6vbh>.
- Wei Zhao, Xiaolu Jiang, Baocan Zhang, Shixiao Xiao, and Sujun Weng. CTNet: a convolutional transformer network for EEG-based motor imagery classification. *Sci. Rep.*, 14(1):20237, August 2024.

A APPENDIX

A.1 DATA PRE-PROCESSING

We apply minimal, standardized pre-processing across all datasets to improve comparability while preserving the raw signal characteristics. We band-pass filter EEG signals between 0.1 and 128 Hz and apply a notch filter at the local power-line frequency using the `mne.filter` module. We then downsample all datasets to a common baseline of 256 Hz to match EEGPT (which operates natively at 256 Hz). For foundation models that require lower sampling rates, we additionally resample the data on the fly using `mne.resample` and then apply each model’s default normalization. For the TS FMs, we follow ST-EEGFormer and use a sampling rate of 128 Hz.

A.2 MODEL TRAINING

All experiments were run on an HPC cluster with NVIDIA H100 GPUs. We consider two standard adaptation strategies for foundation models: linear probing and end-to-end fine-tuning. For linear probing, we freeze the pre-trained backbone and train only a lightweight classification head, which assesses the quality of the learned representations. For fine-tuning, we jointly update all parameters (backbone and head) to adapt the representation to the downstream dataset.

Unless otherwise stated, we follow the training protocol of the EEG benchmark study (Yang et al., 2026), using the AdamW optimizer, a cosine decay learning-rate schedule, and 10 warmup epochs. For Mantis, we adopt the fine-tuning strategy of (Gnassounou et al., 2025): we concatenate tokens across EEG channels and fine-tune with a small backbone learning rate (3×10^{-4}) for up to 50 epochs. For MOMENT, we use a multichannel adapter as recommended for multivariate time series and train with the same overall procedure as Mantis.

A.3 BENCHMARK RESULTS

Table 2: FACED dataset results (all metrics as fractions; mean \pm standard deviation). Best results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
CTNet	Scratch	0.206 \pm 0.075	0.313 \pm 0.075	0.294 \pm 0.067	0.732 \pm 0.051
DeepConvNet	Scratch	0.239 \pm 0.100	0.340 \pm 0.093	0.325 \pm 0.089	0.769 \pm 0.063
EEGNet	Scratch	0.216 \pm 0.062	0.305 \pm 0.058	0.301 \pm 0.053	0.732 \pm 0.052
Conformer	Scratch	0.328 \pm 0.094	0.416 \pm 0.090	0.405 \pm 0.085	0.771 \pm 0.061
EEG foundation models					
BENDR	Fine-tune	0.242 \pm 0.083	0.333 \pm 0.082	0.326 \pm 0.073	0.728 \pm 0.056
BENDR	Linear Probe	0.079 \pm 0.038	0.187 \pm 0.034	0.181 \pm 0.032	0.596 \pm 0.028
BIOT	Fine-tune	0.078 \pm 0.045	0.186 \pm 0.042	0.181 \pm 0.041	0.580 \pm 0.029
BIOT	Linear Probe	0.075 \pm 0.038	0.193 \pm 0.046	0.178 \pm 0.032	0.612 \pm 0.044
CBraMod	Fine-tune	0.282 \pm 0.094	0.361 \pm 0.087	0.363 \pm 0.084	0.751 \pm 0.060
CBraMod	Linear Probe	0.110 \pm 0.038	0.205 \pm 0.040	0.212 \pm 0.035	0.654 \pm 0.040
EEGPT	Fine-tune	0.100 \pm 0.046	0.206 \pm 0.042	0.200 \pm 0.040	0.627 \pm 0.039
EEGPT	Linear Probe	0.185 \pm 0.061	0.279 \pm 0.056	0.275 \pm 0.055	0.707 \pm 0.047
LaBraM	Fine-tune	0.143 \pm 0.051	0.233 \pm 0.045	0.239 \pm 0.045	0.645 \pm 0.044
LaBraM	Linear Probe	0.031 \pm 0.034	0.153 \pm 0.043	0.137 \pm 0.029	0.582 \pm 0.034
ST-EEGFormer-s	Fine-tune	0.170 \pm 0.075	0.269 \pm 0.073	0.264 \pm 0.067	0.675 \pm 0.058
ST-EEGFormer-s	Linear Probe	0.063 \pm 0.035	0.172 \pm 0.044	0.163 \pm 0.029	0.641 \pm 0.044
ST-EEGFormer-b	Fine-tune	0.132 \pm 0.073	0.227 \pm 0.066	0.230 \pm 0.065	0.639 \pm 0.057
ST-EEGFormer-b	Linear Probe	0.062 \pm 0.052	0.170 \pm 0.059	0.163 \pm 0.044	0.636 \pm 0.063
ST-EEGFormer-l	Fine-tune	0.219 \pm 0.088	0.303 \pm 0.075	0.306 \pm 0.079	0.691 \pm 0.057
ST-EEGFormer-l	Linear Probe	0.078 \pm 0.055	0.192 \pm 0.054	0.178 \pm 0.048	0.636 \pm 0.054
Time-series foundation models					
Mantis	Fine-tune	0.223 \pm 0.084	0.301 \pm 0.074	0.311 \pm 0.076	0.725 \pm 0.063
Mantis	Linear Probe	0.135 \pm 0.052	0.236 \pm 0.052	0.233 \pm 0.047	0.659 \pm 0.055
MOMENT	Fine-tune	0.128 \pm 0.059	0.210 \pm 0.046	0.225 \pm 0.053	0.654 \pm 0.052
MOMENT	Linear Probe	0.050 \pm 0.037	0.162 \pm 0.043	0.153 \pm 0.032	0.591 \pm 0.045

Table 3: TUEV dataset results (all metrics as fractions; mean \pm standard deviation). Best results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
CTNet	Scratch	0.351 \pm 0.075	0.632 \pm 0.345	0.635 \pm 0.233	0.776 \pm 0.197
DeepConvNet	Scratch	0.333 \pm 0.297	0.608 \pm 0.393	0.568 \pm 0.194	0.756 \pm 0.215
EEGNet	Scratch	0.317 \pm 0.239	0.560 \pm 0.390	0.544 \pm 0.195	0.734 \pm 0.222
Conformer	Scratch	0.371 \pm 0.191	0.616 \pm 0.388	0.627 \pm 0.217	0.749 \pm 0.225
EEG foundation models					
BENDR	Fine-tune	0.279 \pm 0.364	0.645 \pm 0.374	0.611 \pm 0.217	0.741 \pm 0.222
BENDR	Linear Probe	0.033 \pm 0.116	0.498 \pm 0.300	0.400 \pm 0.142	0.613 \pm 0.190
BIOT	Fine-tune	0.265 \pm 0.354	0.636 \pm 0.360	0.553 \pm 0.222	0.715 \pm 0.240
BIOT	Linear Probe	0.259 \pm 0.346	0.638 \pm 0.344	0.560 \pm 0.215	0.739 \pm 0.226
CBraMod	Fine-tune	0.302 \pm 0.370	0.643 \pm 0.377	0.624 \pm 0.188	0.751 \pm 0.241
CBraMod	Linear Probe	0.092 \pm 0.238	0.473 \pm 0.386	0.435 \pm 0.150	0.630 \pm 0.254
EEGPT	Fine-tune	0.357 \pm 0.374	0.696 \pm 0.331	0.641 \pm 0.203	0.767 \pm 0.243
EEGPT	Linear Probe	0.363 \pm 0.375	0.676 \pm 0.373	0.658 \pm 0.182	0.773 \pm 0.239
LaBraM	Fine-tune	0.381 \pm 0.389	0.725 \pm 0.333	0.670 \pm 0.203	0.759 \pm 0.219
LaBraM	Linear Probe	0.125 \pm 0.256	0.627 \pm 0.380	0.461 \pm 0.171	0.725 \pm 0.221
ST-EEGFormer-s	Fine-tune	0.418 \pm 0.449	0.765 \pm 0.304	0.700 \pm 0.248	0.756 \pm 0.237
ST-EEGFormer-s	Linear Probe	0.107 \pm 0.258	0.645 \pm 0.404	0.456 \pm 0.166	0.749 \pm 0.207
ST-EEGFormer-b	Fine-tune	0.413 \pm 0.414	0.745 \pm 0.324	0.687 \pm 0.250	0.750 \pm 0.235
ST-EEGFormer-b	Linear Probe	0.284 \pm 0.284	0.645 \pm 0.396	0.466 \pm 0.181	0.751 \pm 0.216
ST-EEGFormer-l	Fine-tune	0.398 \pm 0.303	0.738 \pm 0.331	0.690 \pm 0.240	0.750 \pm 0.243
ST-EEGFormer-l	Linear Probe	0.273 \pm 0.273	0.643 \pm 0.398	0.473 \pm 0.170	0.741 \pm 0.217
Time-series foundation models					
Mantis	Fine-tune	0.424 \pm 0.423	0.682 \pm 0.356	0.667 \pm 0.240	0.732 \pm 0.256
Mantis	Linear Probe	0.331 \pm 0.331	0.622 \pm 0.377	0.561 \pm 0.225	0.722 \pm 0.242
MOMENT	Fine-tune	0.286 \pm 0.347	0.693 \pm 0.323	0.589 \pm 0.227	0.745 \pm 0.229
MOMENT	Linear Probe	0.079 \pm 0.225	0.637 \pm 0.385	0.436 \pm 0.172	0.670 \pm 0.265

Table 4: BCI-IV-2A population results (all metrics as fractions; mean \pm standard deviation). Best results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
DeepConvNet	Scratch	0.620 \pm 0.118	0.715 \pm 0.089	0.715 \pm 0.089	0.909 \pm 0.051
EEGNet	Scratch	0.614 \pm 0.104	0.711 \pm 0.078	0.711 \pm 0.078	0.911 \pm 0.045
Conformer	Scratch	0.423 \pm 0.098	0.567 \pm 0.073	0.567 \pm 0.073	0.806 \pm 0.054
CTNet	Scratch	0.665 \pm 0.127	0.749 \pm 0.096	0.749 \pm 0.096	0.929 \pm 0.045
EEG foundation models					
BIOT	Fine-tune	0.255 \pm 0.128	0.441 \pm 0.096	0.441 \pm 0.096	0.687 \pm 0.085
BIOT	Linear Probe	0.240 \pm 0.171	0.430 \pm 0.128	0.430 \pm 0.128	0.692 \pm 0.116
BENDR	Fine-tune	0.287 \pm 0.076	0.465 \pm 0.057	0.465 \pm 0.057	0.725 \pm 0.053
BENDR	Linear Probe	0.095 \pm 0.035	0.321 \pm 0.026	0.321 \pm 0.026	0.573 \pm 0.019
CBraMod	Fine-tune	0.490 \pm 0.102	0.618 \pm 0.077	0.617 \pm 0.077	0.851 \pm 0.051
CBraMod	Linear Probe	0.092 \pm 0.044	0.319 \pm 0.033	0.319 \pm 0.033	0.634 \pm 0.049
EEGPT	Fine-tune	0.203 \pm 0.061	0.402 \pm 0.046	0.402 \pm 0.046	0.674 \pm 0.042
EEGPT	Linear Probe	0.475 \pm 0.104	0.606 \pm 0.078	0.606 \pm 0.078	0.841 \pm 0.054
LaBraM	Fine-tune	0.226 \pm 0.085	0.419 \pm 0.064	0.419 \pm 0.064	0.677 \pm 0.050
LaBraM	Linear Probe	0.048 \pm 0.056	0.286 \pm 0.041	0.286 \pm 0.042	0.572 \pm 0.056
ST-EEGFormer-s	Fine-tune	0.586 \pm 0.154	0.689 \pm 0.116	0.689 \pm 0.116	0.863 \pm 0.069
ST-EEGFormer-s	Linear Probe	0.162 \pm 0.108	0.371 \pm 0.081	0.371 \pm 0.081	0.690 \pm 0.098
ST-EEGFormer-b	Fine-tune	0.700 \pm 0.131	0.775 \pm 0.098	0.775 \pm 0.098	0.908 \pm 0.058
ST-EEGFormer-b	Linear Probe	0.215 \pm 0.128	0.412 \pm 0.096	0.411 \pm 0.096	0.713 \pm 0.092
ST-EEGFormer-l	Fine-tune	0.728 \pm 0.147	0.796 \pm 0.110	0.796 \pm 0.110	0.928 \pm 0.052
ST-EEGFormer-l	Linear Probe	0.188 \pm 0.146	0.390 \pm 0.110	0.391 \pm 0.109	0.716 \pm 0.103
Time-series foundation models					
MOMENT	Fine-tune	0.096 \pm 0.063	0.322 \pm 0.048	0.322 \pm 0.048	0.588 \pm 0.047
MOMENT	Linear Probe	0.070 \pm 0.059	0.304 \pm 0.044	0.303 \pm 0.044	0.586 \pm 0.054
Mantis	Fine-tune	0.307 \pm 0.087	0.480 \pm 0.065	0.480 \pm 0.065	0.749 \pm 0.063
Mantis	Linear Probe	0.246 \pm 0.077	0.434 \pm 0.057	0.434 \pm 0.058	0.695 \pm 0.054

Table 5: BCI-IV-2A leave-one-out zero-shot results (all metrics as fractions; mean \pm standard deviation). Best results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
DeepConvNet	Scratch	0.418 \pm 0.108	0.563 \pm 0.081	0.563 \pm 0.081	0.838 \pm 0.069
EEGNet	Scratch	0.454 \pm 0.129	0.591 \pm 0.097	0.591 \pm 0.097	0.858 \pm 0.063
Conformer	Scratch	0.308 \pm 0.092	0.481 \pm 0.069	0.481 \pm 0.069	0.738 \pm 0.062
CTNet	Scratch	0.489 \pm 0.150	0.617 \pm 0.112	0.617 \pm 0.112	0.876 \pm 0.068
EEG foundation models					
BIOT	Fine-tune	0.084 \pm 0.094	0.313 \pm 0.071	0.313 \pm 0.071	0.578 \pm 0.064
BIOT	Linear Probe	0.086 \pm 0.071	0.315 \pm 0.053	0.315 \pm 0.053	0.597 \pm 0.071
BENDR	Fine-tune	0.217 \pm 0.063	0.412 \pm 0.048	0.412 \pm 0.047	0.676 \pm 0.046
BENDR	Linear Probe	0.065 \pm 0.026	0.299 \pm 0.021	0.299 \pm 0.020	0.555 \pm 0.024
CBraMod	Fine-tune	0.359 \pm 0.097	0.519 \pm 0.073	0.519 \pm 0.073	0.775 \pm 0.062
CBraMod	Linear Probe	0.114 \pm 0.039	0.335 \pm 0.030	0.335 \pm 0.029	0.623 \pm 0.050
EEGPT	Fine-tune	0.129 \pm 0.049	0.347 \pm 0.037	0.347 \pm 0.037	0.618 \pm 0.035
EEGPT	Linear Probe	0.330 \pm 0.137	0.498 \pm 0.103	0.498 \pm 0.103	0.774 \pm 0.064
LaBraM	Fine-tune	0.081 \pm 0.056	0.311 \pm 0.042	0.311 \pm 0.042	0.573 \pm 0.040
LaBraM	Linear Probe	0.040 \pm 0.044	0.279 \pm 0.031	0.280 \pm 0.033	0.557 \pm 0.050
ST-EEGFormer-s	Fine-tune	0.349 \pm 0.166	0.512 \pm 0.125	0.512 \pm 0.125	0.746 \pm 0.083
ST-EEGFormer-s	Linear Probe	0.031 \pm 0.074	0.273 \pm 0.055	0.273 \pm 0.055	0.642 \pm 0.087
ST-EEGFormer-b	Fine-tune	0.377 \pm 0.151	0.533 \pm 0.114	0.533 \pm 0.114	0.780 \pm 0.085
ST-EEGFormer-b	Linear Probe	0.043 \pm 0.089	0.282 \pm 0.066	0.282 \pm 0.066	0.660 \pm 0.095
ST-EEGFormer-l	Fine-tune	0.434 \pm 0.197	0.575 \pm 0.148	0.575 \pm 0.148	0.800 \pm 0.103
ST-EEGFormer-l	Linear Probe	0.074 \pm 0.104	0.305 \pm 0.078	0.306 \pm 0.078	0.672 \pm 0.105
Time-series foundation models					
MOMENT	Fine-tune	0.073 \pm 0.052	0.304 \pm 0.039	0.304 \pm 0.039	0.557 \pm 0.034
MOMENT	Linear Probe	0.047 \pm 0.045	0.285 \pm 0.034	0.285 \pm 0.034	0.571 \pm 0.055
Mantis	Fine-tune	0.212 \pm 0.075	0.409 \pm 0.056	0.409 \pm 0.056	0.688 \pm 0.048
Mantis	Linear Probe	0.181 \pm 0.078	0.385 \pm 0.059	0.386 \pm 0.058	0.653 \pm 0.044

Table 6: BCI-IV-2A leave-one-out fine-tuning results (all metrics as fractions; mean \pm standard deviation). Best results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
DeepConvNet	Scratch	0.665 \pm 0.118	0.748 \pm 0.089	0.748 \pm 0.089	0.920 \pm 0.049
EEGNet	Scratch	0.681 \pm 0.113	0.760 \pm 0.085	0.760 \pm 0.085	0.929 \pm 0.043
Conformer	Scratch	0.501 \pm 0.091	0.626 \pm 0.068	0.626 \pm 0.068	0.846 \pm 0.052
CTNet	Scratch	0.725 \pm 0.136	0.794 \pm 0.102	0.794 \pm 0.102	0.947 \pm 0.040
EEG foundation models					
BIOT	Fine-tune	0.254 \pm 0.165	0.440 \pm 0.124	0.440 \pm 0.124	0.689 \pm 0.098
BIOT	Linear Probe	0.306 \pm 0.149	0.480 \pm 0.111	0.480 \pm 0.111	0.736 \pm 0.100
BENDR	Fine-tune	0.347 \pm 0.119	0.510 \pm 0.089	0.510 \pm 0.089	0.769 \pm 0.064
BENDR	Linear Probe	0.101 \pm 0.032	0.326 \pm 0.024	0.326 \pm 0.024	0.589 \pm 0.020
CBraMod	Fine-tune	0.528 \pm 0.105	0.646 \pm 0.079	0.646 \pm 0.079	0.871 \pm 0.052
CBraMod	Linear Probe	0.131 \pm 0.047	0.348 \pm 0.035	0.348 \pm 0.035	0.634 \pm 0.048
EEGPT	Fine-tune	0.245 \pm 0.095	0.434 \pm 0.071	0.434 \pm 0.071	0.698 \pm 0.054
EEGPT	Linear Probe	0.428 \pm 0.121	0.571 \pm 0.091	0.571 \pm 0.091	0.815 \pm 0.067
LaBraM	Fine-tune	0.164 \pm 0.075	0.373 \pm 0.056	0.373 \pm 0.056	0.617 \pm 0.061
LaBraM	Linear Probe	0.085 \pm 0.071	0.313 \pm 0.053	0.314 \pm 0.053	0.581 \pm 0.061
ST-EEGFormer-s	Fine-tune	0.634 \pm 0.130	0.725 \pm 0.098	0.725 \pm 0.098	0.896 \pm 0.060
ST-EEGFormer-s	Linear Probe	0.186 \pm 0.117	0.390 \pm 0.088	0.390 \pm 0.088	0.713 \pm 0.092
ST-EEGFormer-b	Fine-tune	0.707 \pm 0.119	0.780 \pm 0.089	0.780 \pm 0.089	0.929 \pm 0.049
ST-EEGFormer-b	Linear Probe	0.218 \pm 0.102	0.413 \pm 0.077	0.414 \pm 0.077	0.714 \pm 0.089
ST-EEGFormer-l	Fine-tune	0.744 \pm 0.143	0.808 \pm 0.108	0.808 \pm 0.108	0.937 \pm 0.057
ST-EEGFormer-l	Linear Probe	0.161 \pm 0.117	0.371 \pm 0.088	0.371 \pm 0.088	0.713 \pm 0.099
Time-series foundation models					
MOMENT	Fine-tune	0.152 \pm 0.051	0.363 \pm 0.039	0.364 \pm 0.038	0.627 \pm 0.038
MOMENT	Linear Probe	0.083 \pm 0.059	0.312 \pm 0.044	0.312 \pm 0.044	0.592 \pm 0.057
Mantis	Fine-tune	0.355 \pm 0.122	0.517 \pm 0.091	0.517 \pm 0.091	0.768 \pm 0.071
Mantis	Linear Probe	0.210 \pm 0.079	0.408 \pm 0.059	0.408 \pm 0.059	0.675 \pm 0.051

Table 7: BCI-IV-2A leave-one-out generalization drop after fine-tuning on seen subjects (all metrics as fractions; mean \pm standard deviation). Smaller is better; best (smallest) results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
DeepConvNet	Scratch	0.256 \pm 0.048	0.192 \pm 0.036	0.192 \pm 0.036	0.098 \pm 0.020
EEGNet	Scratch	0.245 \pm 0.034	0.183 \pm 0.025	0.183 \pm 0.025	0.086 \pm 0.011
Conformer	Scratch	0.097 \pm 0.021	0.073 \pm 0.016	0.073 \pm 0.016	0.047 \pm 0.013
CTNet	Scratch	0.323 \pm 0.071	0.242 \pm 0.054	0.242 \pm 0.054	0.113 \pm 0.019
EEG foundation models					
BIOT	Fine-tune	0.089 \pm 0.050	0.067 \pm 0.037	0.067 \pm 0.037	0.050 \pm 0.022
BIOT	Linear Probe	0.129 \pm 0.063	0.097 \pm 0.047	0.097 \pm 0.047	0.068 \pm 0.034
BENDR	Fine-tune	0.068 \pm 0.029	0.051 \pm 0.022	0.051 \pm 0.022	0.047 \pm 0.014
BENDR	Linear Probe	0.007 \pm 0.013	0.005 \pm 0.009	0.005 \pm 0.009	0.005 \pm 0.007
CBraMod	Fine-tune	0.103 \pm 0.029	0.077 \pm 0.022	0.077 \pm 0.022	0.054 \pm 0.011
CBraMod	Linear Probe	0.043 \pm 0.029	0.033 \pm 0.022	0.032 \pm 0.022	0.043 \pm 0.013
EEGPT	Fine-tune	0.031 \pm 0.029	0.023 \pm 0.022	0.023 \pm 0.022	0.018 \pm 0.017
EEGPT	Linear Probe	0.267 \pm 0.048	0.200 \pm 0.036	0.200 \pm 0.036	0.146 \pm 0.038
LaBraM	Fine-tune	0.012 \pm 0.016	0.009 \pm 0.012	0.009 \pm 0.012	0.003 \pm 0.008
LaBraM	Linear Probe	0.017 \pm 0.019	0.013 \pm 0.014	0.013 \pm 0.015	0.021 \pm 0.013
ST-EEGFormer-s	Fine-tune	0.176 \pm 0.050	0.132 \pm 0.038	0.132 \pm 0.037	0.063 \pm 0.009
ST-EEGFormer-s	Linear Probe	0.094 \pm 0.042	0.070 \pm 0.032	0.070 \pm 0.032	0.021 \pm 0.016
ST-EEGFormer-b	Fine-tune	0.145 \pm 0.034	0.108 \pm 0.025	0.108 \pm 0.025	0.048 \pm 0.012
ST-EEGFormer-b	Linear Probe	0.118 \pm 0.053	0.088 \pm 0.039	0.088 \pm 0.039	0.019 \pm 0.015
ST-EEGFormer-l	Fine-tune	0.146 \pm 0.035	0.109 \pm 0.026	0.109 \pm 0.026	0.044 \pm 0.017
ST-EEGFormer-l	Linear Probe	0.115 \pm 0.051	0.086 \pm 0.038	0.086 \pm 0.038	0.020 \pm 0.015
Time-series foundation models					
MOMENT	Fine-tune	0.020 \pm 0.024	0.015 \pm 0.018	0.015 \pm 0.018	0.008 \pm 0.016
MOMENT	Linear Probe	0.022 \pm 0.023	0.017 \pm 0.018	0.017 \pm 0.017	0.000 \pm 0.003
Mantis	Fine-tune	0.064 \pm 0.046	0.048 \pm 0.034	0.048 \pm 0.034	0.027 \pm 0.011
Mantis	Linear Probe	0.097 \pm 0.020	0.073 \pm 0.015	0.073 \pm 0.015	0.062 \pm 0.015

Table 8: Error population results (all metrics as fractions; mean \pm standard deviation). Best results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
DeepConvNet	Scratch	0.979 \pm 0.013	0.993 \pm 0.004	0.990 \pm 0.009	0.998 \pm 0.004
EEGNet	Scratch	0.976 \pm 0.021	0.993 \pm 0.007	0.986 \pm 0.015	0.999 \pm 0.003
Conformer	Scratch	0.974 \pm 0.021	0.992 \pm 0.006	0.985 \pm 0.014	0.997 \pm 0.007
CTNet	Scratch	0.981 \pm 0.020	0.994 \pm 0.006	0.988 \pm 0.014	0.998 \pm 0.002
EEG foundation models					
BIOT	Fine-tune	0.944 \pm 0.028	0.983 \pm 0.009	0.964 \pm 0.016	0.992 \pm 0.009
BIOT	Linear Probe	0.880 \pm 0.036	0.962 \pm 0.011	0.937 \pm 0.021	0.979 \pm 0.012
BENDR	Fine-tune	0.976 \pm 0.012	0.992 \pm 0.004	0.985 \pm 0.010	0.998 \pm 0.002
BENDR	Linear Probe	0.939 \pm 0.032	0.981 \pm 0.010	0.963 \pm 0.019	0.988 \pm 0.009
CBraMod	Fine-tune	0.976 \pm 0.021	0.992 \pm 0.007	0.988 \pm 0.012	0.996 \pm 0.006
CBraMod	Linear Probe	0.898 \pm 0.033	0.967 \pm 0.011	0.951 \pm 0.015	0.987 \pm 0.008
EEGPT	Fine-tune	0.967 \pm 0.016	0.990 \pm 0.005	0.975 \pm 0.012	0.999 \pm 0.001
EEGPT	Linear Probe	0.982 \pm 0.014	0.994 \pm 0.004	0.987 \pm 0.009	0.999 \pm 0.002
LaBraM	Fine-tune	0.979 \pm 0.015	0.993 \pm 0.005	0.986 \pm 0.011	0.999 \pm 0.003
LaBraM	Linear Probe	0.197 \pm 0.101	0.823 \pm 0.016	0.570 \pm 0.039	0.859 \pm 0.058
ST-EEGFormer-s	Fine-tune	0.988 \pm 0.015	0.996 \pm 0.005	0.991 \pm 0.011	0.998 \pm 0.005
ST-EEGFormer-s	Linear Probe	0.625 \pm 0.134	0.898 \pm 0.030	0.778 \pm 0.079	0.951 \pm 0.024
ST-EEGFormer-b	Fine-tune	0.979 \pm 0.017	0.994 \pm 0.005	0.985 \pm 0.012	0.993 \pm 0.011
ST-EEGFormer-b	Linear Probe	0.410 \pm 0.159	0.852 \pm 0.035	0.667 \pm 0.077	0.929 \pm 0.037
ST-EEGFormer-l	Fine-tune	0.987 \pm 0.015	0.996 \pm 0.004	0.992 \pm 0.010	0.998 \pm 0.005
ST-EEGFormer-l	Linear Probe	0.703 \pm 0.120	0.918 \pm 0.027	0.817 \pm 0.071	0.965 \pm 0.020
Time-series foundation models					
MOMENT	Fine-tune	0.949 \pm 0.028	0.985 \pm 0.008	0.971 \pm 0.017	0.997 \pm 0.002
MOMENT	Linear Probe	0.488 \pm 0.191	0.866 \pm 0.041	0.710 \pm 0.093	0.931 \pm 0.040
Mantis	Fine-tune	0.971 \pm 0.022	0.991 \pm 0.007	0.979 \pm 0.016	0.997 \pm 0.005
Mantis	Linear Probe	0.932 \pm 0.030	0.979 \pm 0.009	0.966 \pm 0.014	0.992 \pm 0.007

Table 9: Error leave-one-out zero-shot results (all metrics as fractions; mean \pm standard deviation). Best results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
DeepConvNet	Scratch	0.824 \pm 0.104	0.946 \pm 0.029	0.905 \pm 0.068	0.964 \pm 0.040
EEGNet	Scratch	0.764 \pm 0.114	0.932 \pm 0.028	0.863 \pm 0.082	0.961 \pm 0.037
Conformer	Scratch	0.791 \pm 0.107	0.938 \pm 0.030	0.873 \pm 0.066	0.961 \pm 0.035
CTNet	Scratch	0.836 \pm 0.077	0.950 \pm 0.021	0.905 \pm 0.056	0.973 \pm 0.031
EEG foundation models					
BIOT	Fine-tune	0.539 \pm 0.112	0.862 \pm 0.027	0.761 \pm 0.079	0.857 \pm 0.066
BIOT	Linear Probe	0.400 \pm 0.122	0.817 \pm 0.049	0.693 \pm 0.066	0.820 \pm 0.062
BENDR	Fine-tune	0.737 \pm 0.117	0.922 \pm 0.032	0.851 \pm 0.071	0.934 \pm 0.046
BENDR	Linear Probe	0.536 \pm 0.164	0.868 \pm 0.042	0.747 \pm 0.090	0.855 \pm 0.068
CBraMod	Fine-tune	0.733 \pm 0.171	0.925 \pm 0.037	0.848 \pm 0.101	0.943 \pm 0.060
CBraMod	Linear Probe	0.533 \pm 0.162	0.855 \pm 0.061	0.753 \pm 0.077	0.877 \pm 0.069
EEGPT	Fine-tune	0.709 \pm 0.152	0.918 \pm 0.037	0.826 \pm 0.080	0.935 \pm 0.067
EEGPT	Linear Probe	0.767 \pm 0.114	0.933 \pm 0.028	0.860 \pm 0.069	0.950 \pm 0.051
LaBraM	Fine-tune	0.767 \pm 0.141	0.933 \pm 0.035	0.861 \pm 0.086	0.959 \pm 0.036
LaBraM	Linear Probe	0.160 \pm 0.094	0.817 \pm 0.016	0.556 \pm 0.035	0.831 \pm 0.070
ST-EEGFormer-s	Fine-tune	0.837 \pm 0.104	0.952 \pm 0.028	0.907 \pm 0.073	0.967 \pm 0.023
ST-EEGFormer-s	Linear Probe	0.550 \pm 0.165	0.877 \pm 0.037	0.750 \pm 0.099	0.926 \pm 0.046
ST-EEGFormer-b	Fine-tune	0.774 \pm 0.135	0.937 \pm 0.033	0.862 \pm 0.085	0.958 \pm 0.033
ST-EEGFormer-b	Linear Probe	0.323 \pm 0.192	0.837 \pm 0.038	0.633 \pm 0.097	0.909 \pm 0.045
ST-EEGFormer-l	Fine-tune	0.811 \pm 0.156	0.948 \pm 0.040	0.886 \pm 0.101	0.964 \pm 0.039
ST-EEGFormer-l	Linear Probe	0.581 \pm 0.193	0.892 \pm 0.044	0.756 \pm 0.109	0.946 \pm 0.033
Time-series foundation models					
MOMENT	Fine-tune	0.647 \pm 0.122	0.894 \pm 0.034	0.814 \pm 0.072	0.926 \pm 0.042
MOMENT	Linear Probe	0.427 \pm 0.213	0.853 \pm 0.044	0.682 \pm 0.100	0.920 \pm 0.050
Mantis	Fine-tune	0.717 \pm 0.131	0.921 \pm 0.030	0.834 \pm 0.086	0.952 \pm 0.034
Mantis	Linear Probe	0.503 \pm 0.103	0.854 \pm 0.024	0.743 \pm 0.075	0.839 \pm 0.063

Table 10: Error leave-one-out fine-tuning results (all metrics as fractions; mean \pm standard deviation). Best results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
DeepConvNet	Scratch	0.990 \pm 0.013	0.997 \pm 0.004	0.996 \pm 0.006	1.000 \pm 0.000
EEGNet	Scratch	0.991 \pm 0.006	0.997 \pm 0.002	0.994 \pm 0.005	1.000 \pm 0.001
Conformer	Scratch	0.989 \pm 0.011	0.997 \pm 0.003	0.993 \pm 0.009	0.999 \pm 0.004
CTNet	Scratch	0.989 \pm 0.017	0.997 \pm 0.005	0.994 \pm 0.011	0.999 \pm 0.002
EEG foundation models					
BIOT	Fine-tune	0.955 \pm 0.029	0.986 \pm 0.009	0.972 \pm 0.021	0.991 \pm 0.010
BIOT	Linear Probe	0.925 \pm 0.021	0.976 \pm 0.007	0.956 \pm 0.014	0.993 \pm 0.004
BENDR	Fine-tune	0.979 \pm 0.014	0.993 \pm 0.004	0.987 \pm 0.009	1.000 \pm 0.000
BENDR	Linear Probe	0.916 \pm 0.037	0.973 \pm 0.012	0.957 \pm 0.015	0.995 \pm 0.005
CBraMod	Fine-tune	0.984 \pm 0.015	0.995 \pm 0.005	0.989 \pm 0.012	0.998 \pm 0.004
CBraMod	Linear Probe	0.933 \pm 0.026	0.978 \pm 0.008	0.965 \pm 0.011	0.992 \pm 0.009
EEGPT	Fine-tune	0.968 \pm 0.018	0.990 \pm 0.006	0.979 \pm 0.011	0.999 \pm 0.002
EEGPT	Linear Probe	0.986 \pm 0.012	0.995 \pm 0.004	0.989 \pm 0.008	1.000 \pm 0.000
LaBraM	Fine-tune	0.982 \pm 0.013	0.994 \pm 0.004	0.987 \pm 0.011	0.998 \pm 0.004
LaBraM	Linear Probe	0.368 \pm 0.158	0.851 \pm 0.027	0.643 \pm 0.073	0.905 \pm 0.044
ST-EEGFormer-s	Fine-tune	0.983 \pm 0.019	0.995 \pm 0.006	0.990 \pm 0.013	0.997 \pm 0.005
ST-EEGFormer-s	Linear Probe	0.693 \pm 0.164	0.916 \pm 0.034	0.821 \pm 0.097	0.967 \pm 0.016
ST-EEGFormer-b	Fine-tune	0.980 \pm 0.016	0.994 \pm 0.005	0.985 \pm 0.013	0.996 \pm 0.006
ST-EEGFormer-b	Linear Probe	0.259 \pm 0.168	0.834 \pm 0.026	0.598 \pm 0.071	0.930 \pm 0.035
ST-EEGFormer-l	Fine-tune	0.985 \pm 0.020	0.996 \pm 0.006	0.990 \pm 0.015	0.998 \pm 0.003
ST-EEGFormer-l	Linear Probe	0.554 \pm 0.163	0.887 \pm 0.033	0.732 \pm 0.087	0.963 \pm 0.022
Time-series foundation models					
MOMENT	Fine-tune	0.955 \pm 0.021	0.987 \pm 0.006	0.974 \pm 0.009	0.998 \pm 0.001
MOMENT	Linear Probe	0.533 \pm 0.202	0.879 \pm 0.046	0.732 \pm 0.102	0.938 \pm 0.040
Mantis	Fine-tune	0.971 \pm 0.022	0.991 \pm 0.006	0.979 \pm 0.015	0.998 \pm 0.005
Mantis	Linear Probe	0.918 \pm 0.041	0.975 \pm 0.013	0.963 \pm 0.016	0.980 \pm 0.015

Table 11: Generalization drop on the Error dataset after fine-tuning on seen subjects (all metrics as fractions; mean \pm standard deviation). Smaller is better; best (smallest) results are shown in bold.

Model	Training	Kappa	Acc1	Balanced Acc	AUC
Classic neural networks					
DeepConvNet	Scratch	0.168 \pm 0.061	0.048 \pm 0.015	0.101 \pm 0.040	0.022 \pm 0.010
EEGNet	Scratch	0.207 \pm 0.103	0.055 \pm 0.024	0.130 \pm 0.063	0.018 \pm 0.007
Conformer	Scratch	0.195 \pm 0.049	0.055 \pm 0.010	0.116 \pm 0.040	0.027 \pm 0.012
CTNet	Scratch	0.169 \pm 0.095	0.047 \pm 0.022	0.107 \pm 0.061	0.019 \pm 0.011
EEG foundation models					
BIOT	Fine-tune	0.184 \pm 0.048	0.053 \pm 0.010	0.101 \pm 0.041	0.039 \pm 0.024
BIOT	Linear Probe	0.410 \pm 0.085	0.132 \pm 0.048	0.196 \pm 0.065	0.087 \pm 0.021
BENDR	Fine-tune	0.307 \pm 0.046	0.082 \pm 0.010	0.187 \pm 0.029	0.059 \pm 0.020
BENDR	Linear Probe	0.224 \pm 0.065	0.063 \pm 0.017	0.135 \pm 0.036	0.045 \pm 0.016
CBraMod	Fine-tune	0.265 \pm 0.146	0.069 \pm 0.031	0.155 \pm 0.093	0.028 \pm 0.021
CBraMod	Linear Probe	0.189 \pm 0.054	0.053 \pm 0.018	0.123 \pm 0.039	0.038 \pm 0.015
EEGPT	Fine-tune	0.156 \pm 0.091	0.043 \pm 0.023	0.096 \pm 0.056	0.015 \pm 0.011
EEGPT	Linear Probe	0.165 \pm 0.051	0.046 \pm 0.013	0.109 \pm 0.032	0.015 \pm 0.004
LaBraM	Fine-tune	0.090 \pm 0.105	0.025 \pm 0.025	0.055 \pm 0.067	0.004 \pm 0.004
LaBraM	Linear Probe	-0.002 \pm 0.099	0.006 \pm 0.009	-0.006 \pm 0.041	0.049 \pm 0.030
ST-EEGFormer-s	Fine-tune	0.050 \pm 0.080	0.014 \pm 0.021	0.032 \pm 0.053	0.002 \pm 0.002
ST-EEGFormer-s	Linear Probe	0.086 \pm 0.135	0.030 \pm 0.024	0.017 \pm 0.085	0.018 \pm 0.004
ST-EEGFormer-b	Fine-tune	0.039 \pm 0.037	0.012 \pm 0.011	0.027 \pm 0.025	0.005 \pm 0.012
ST-EEGFormer-b	Linear Probe	0.085 \pm 0.133	0.012 \pm 0.017	0.032 \pm 0.058	0.007 \pm 0.004
ST-EEGFormer-l	Fine-tune	0.043 \pm 0.042	0.013 \pm 0.012	0.027 \pm 0.031	0.004 \pm 0.005
ST-EEGFormer-l	Linear Probe	0.124 \pm 0.157	0.024 \pm 0.029	0.059 \pm 0.087	0.003 \pm 0.004
Time-series foundation models					
MOMENT	Fine-tune	0.169 \pm 0.152	0.044 \pm 0.034	0.098 \pm 0.096	0.019 \pm 0.019
MOMENT	Linear Probe	-0.022 \pm 0.117	-0.003 \pm 0.018	-0.020 \pm 0.069	0.001 \pm 0.001
Mantis	Fine-tune	0.081 \pm 0.060	0.024 \pm 0.016	0.051 \pm 0.044	0.003 \pm 0.003
Mantis	Linear Probe	0.265 \pm 0.059	0.073 \pm 0.010	0.158 \pm 0.048	0.052 \pm 0.009

Table 12: Pre-trained vs. randomly initialized Mantis (fine-tuning only). We report balanced accuracy and Cohen’s κ (mean \pm std). * $p < 0.05$, ** $p < 0.01$ (paired Wilcoxon signed-rank test). Bold indicates the significantly better result.

Dataset	Setting	Balanced Acc			Kappa		
		Pre-trained	Random Init	Sig.	Pre-trained	Random Init	Sig.
TUEV	Cross-sub ($n=61$)	0.667 \pm 0.239	0.649 \pm 0.231	n.s.	0.400 \pm 0.420	0.382 \pm 0.407	n.s.
FACED	Cross-sub ($n=25$)	0.311 \pm 0.075	0.308 \pm 0.074	n.s.	0.223 \pm 0.083	0.220 \pm 0.082	n.s.
BCI-IV-2A	Pop ($n=9$)	0.480 \pm 0.069	0.535 \pm 0.060	**	0.307 \pm 0.091	0.381 \pm 0.081	**
BCI-IV-2A	LOO ZS ($n=9$)	0.409 \pm 0.056	0.462 \pm 0.060	n.s.	0.212 \pm 0.075	0.282 \pm 0.081	*
BCI-IV-2A	LOO FT ($n=9$)	0.517 \pm 0.091	0.510 \pm 0.064	n.s.	0.355 \pm 0.122	0.346 \pm 0.086	n.s.
Error	Pop ($n=8$)	0.979 \pm 0.017	0.964 \pm 0.037	n.s.	0.971 \pm 0.024	0.934 \pm 0.060	*
Error	LOO ZS ($n=8$)	0.834 \pm 0.086	0.881 \pm 0.062	*	0.717 \pm 0.131	0.778 \pm 0.094	n.s.
Error	LOO FT ($n=8$)	0.979 \pm 0.015	0.968 \pm 0.021	n.s.	0.971 \pm 0.022	0.936 \pm 0.034	**