
Device Passport: Enabling Spatio-Temporal Pretrained Models to Generalize Across Input Layouts

Anonymous Authors¹

Abstract

New device layouts pose a challenging modeling problem due to the lack of large datasets for that specific layout. Biosignal foundation models offer a plausible solution if they are able to generalize to new layouts effectively. To improve cross-layout transfer, we study how different channel embedding techniques behave when pretraining layouts differ substantially from the downstream decoding layout. We propose Device Passport, a new channel embedding technique that learns experts and mixture models to take as input each channel’s functional activity together with its metadata. This contrasts with prior embedding methods, which typically use only functional information or only metadata to look up learned or fixed positional embeddings. Across controlled subset-transfer experiments and realistic transfer to ear-EEG, our learned-expert method is competitive overall and improves over the strongest learned baseline in the layout-transfer regimes that motivate this work. These results suggest that channel embedding design is a key consideration when reusing large-scale pretrained biosignal models on new devices.

1. Introduction

Across biosignal domains such as EMG, EKG, and EEG, large-scale biosignal foundation models have shown to improve downstream decoding, especially when pretraining and evaluation share similar sensor layouts (Kaifosh & Reardon, 2025; Abbaspourazad et al., 2023; Wang et al., 2024). However, it is common for channels to be placed in new locations when experimenting with new devices, constituting new unseen sensor layouts. In these settings, extensive pretraining data may not be available or only be available

in mismatched layouts, so it would be valuable for such pretrained models to generalize to new device layouts. This problem is central to wearable health sensing, where new form factors often collect high-frequency physiological time series before large device-specific cohorts exist. For many prototype or emerging health devices, collecting large labeled cohorts for every new montage is impractical, so the useful question is not only whether a pretrained biosignal model works, but whether its spatial knowledge can be reused when the sensor layout changes.

A core challenge in this setting is channel embedding. Transformers offer channel count flexibility, but rely on positional or channel embeddings to contextualize each input (Vaswani et al., 2017; Chau et al., 2025). Most positional or lookup-based embedding schemes assume repeated identities or enough examples per channel to relearn useful embeddings after transfer. Under strong layout shift and low-data experimental settings, these assumptions break, and the model must recover spatial relationships for sensor channels whose locations or identities were not available during pretraining. We therefore argue that channel embeddings are a core weakness in unseen-layout transfer.

Prior channel embedding methods span identity-based, coordinate-based, and activity-conditioned schemes, including strong recent approaches such as asymmetric channel positional embedding (ACPE) (Wang et al., 2024). However, it remains unclear how well these methods and ACPE transfer when the downstream montage differs or data are scarce. This motivates a focused study of how to best learn and initialize channel embeddings when layout transfer is large. In particular, we study whether channel embeddings should be derived only from functional activity, or whether they should instead be estimated from both metadata and functional activity.

In this work, we systematically compare six channel embedding strategies, including sinusoidal, channel ID, xyz-based encodings, ACPE, and two novel expert-based adapters, which we call Device Passport, that learn to leverage metadata and functional activity of new channels to estimate channel embeddings. We evaluate these methods in two settings: a controlled toy transfer problem based on pretraining and fine-tuning on different TUH/TUAB electrode subsets,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

and a realistic transfer problem from full-layout TUH pre-training to ear-EEG sleep staging on EESM17. Across both settings, Device Passport is most useful when transferring to new layouts, suggesting a practical path for reusing spatial knowledge across devices without relearning channel embeddings from scratch.

2. Approach and Methods

Models. We conduct our experiments on a single common backbone model, the CBraMod EEG model (Wang et al., 2024), which is a spatio-temporal transformer pretrained using masked reconstruction of spatiotemporal patches. A key component of CBraMod is the channel embedding technique, Asymmetric Conditional Positional Encoding (ACPE), which dynamically encodes spatial and temporal position. In our experiments, we keep the original pre-training pipeline and preprocessing (e.g., filtering and data scaling), and vary only the channel embedding technique across a broader family of alternatives, including ACPE and simpler lookup-based encodings.

Channel Embedding Techniques. We investigate 6 different positional encoding techniques. Four are baselines (APE, Channel ID, xyz, ACPE) and two are variations of our proposed expert-based channel embedding adapters (MLP experts and Cross-attention experts). The baselines span ordinal, identity, and coordinate-based encodings (Vaswani et al., 2017; Azabou et al., 2023; Jiang et al., 2024; Wang et al., 2024). Our method, named Device Passport, uses channel location (xyz) and the channel patch embedding as inputs to a mixture model for producing a channel embedding. During pretraining, Device Passport learns a bank of 10 expert embeddings together with either an MLP or a cross-attention module that maps metadata and functional activity to a mixture over those experts. This design separates reusable spatial anchors from the layout-specific rule that combines them. In contrast to randomly learning a new embedding for each unseen channel, the downstream model adapts how metadata and activity select among pretrained experts. This enables the models to derive a contextualized channel embedding given previously learned expert anchors and mapping models. At downstream transfer time, the expert embeddings are frozen, while the MLP or cross-attention module is fine-tuned so that the model can reuse pretrained spatial anchors while adapting the mixture rule to the new layout. This differs from ACPE, which derives positional information from patch embeddings and neighboring patches and continues tuning the convolutional kernel during both pretraining and fine-tuning. Our goal is to test whether explicitly combining metadata with functional activity leads to better initialization and transfer of channel embeddings when layouts are unseen. More details of each are available in Figure 1 and Appendix B.

Pretraining. We follow CBraMod’s pipeline with hyperparameters and preprocessing of data (Wang et al., 2024) closely. The dataset used is the Temple University Hospital (TUH) EEG Corpus (Obeid & Picone, 2016), which consists of a primarily 19 channel 10-20 electrode configuration, across 14k subjects for a total of 27k hours of recording. Similar to CBraMod’s proposed pretraining, we perform data cleaning with a 100uV cut off to remove around 2/3 of the full dataset, leaving around 9k hours of pretraining data. This was shown to improve model downstream performance in the original work, so we adopt it for our exploration of different channel embedding techniques. Across variants, we keep the backbone, objective, and training setup fixed so that differences reflect the channel embedding design.

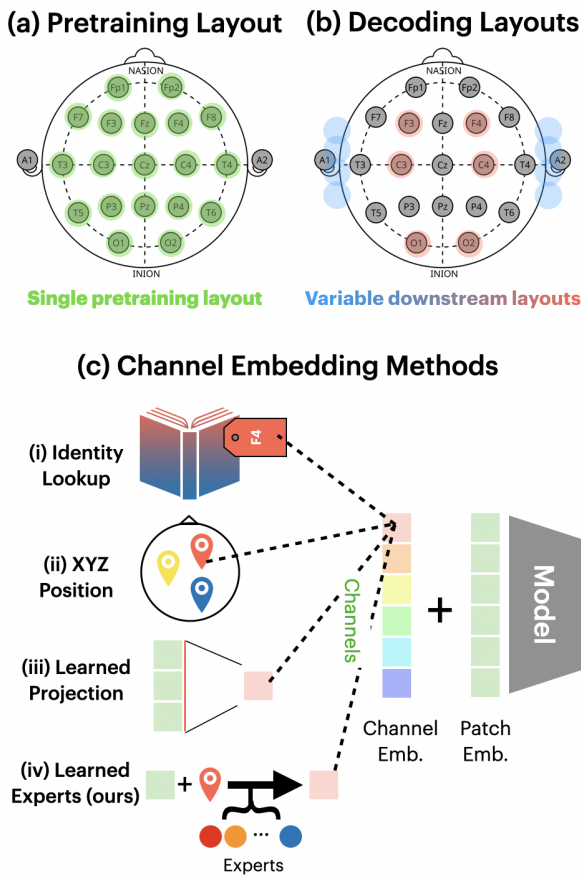


Figure 1. Layout Transfer Challenge + Channel Embedding Techniques. (a) Pretraining often occurs on a single pretraining layout. (b) Decoding needs to work on variable downstream layouts. (c) Channel Embedding Methods help identify the origin of functional activity, but many techniques do not learn transferable representations due to channel layout mismatch between pretraining and decoding. (c.i) Identity lookup channel embeddings can be learned (channel id) or fixed (APE). (c.ii) XYZ position based lookup (xyz) as used in (Chau et al., 2025). (c.iii) Learned projection (ACPE) as used in (Wang et al., 2024). (c.iv) Our method learns mixture models and expert embeddings during pretraining that can be leveraged during downstream decoding.

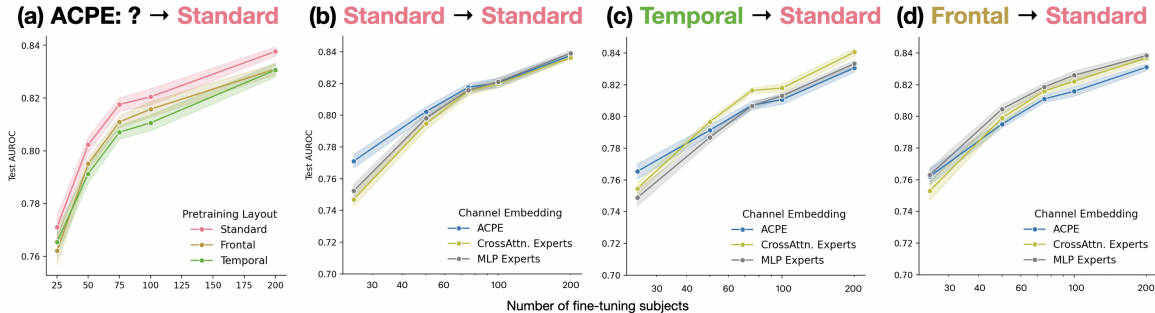


Figure 2. **Variable pretraining layouts.** Downstream layout decoding performance (y-axis) across number of fine-tuning subjects (x-axis), of different layout transfer problems (a) and of different channel embedding techniques across different layout transfer problems (b-d). (a) Pretraining with different layouts (colors) hurts ACPE’s ability to transfer performance to downstream layout (Standard). (b) Pretrain with Standard, fine-tune with Standard. (c) Pretrain with Temporal, fine-tune with Standard. (d) Pretrain with Frontal, fine-tune with Standard.

The main pretraining loss is a reconstruction loss on masked patches. For adapters (MLP and Cross Attention Experts) that require patch embedding as input, we use an average of the non-masked patches for each channel. For ACPE, the patch embeddings across time and space (no modification on masked vs non-masked) get multiplied by the learned CNN kernel to provide the positional embedding.

Downstream Decoding. We fine-tune using the hyperparameters reported in CBraMod (Wang et al., 2024). For sample-efficiency sweeps, we match the full fine-tuning compute budget by fixing the number of optimization steps based on CBraMod’s epoch schedule. We evaluate Temple University Abnormal EEG Detection (TUAB) (Obeid & Picone, 2016), an abnormal EEG classification task that closely matches the pretraining data family, and EESM17 (Mikkelsen et al., 2017), an ear-EEG sleep staging task with a disjoint downstream layout. For TUAB, the main metric shown in the toy transfer setting is test AUROC as a function of the number of fine-tuning subjects. For EESM17, the main figure shows relative Cohen’s kappa versus ACPE per held-out subject, while the appendix reports absolute performance across all channel embedding techniques. We use subject-disjoint train/val/test splits for TUAB and leave-one-subject-out (LOSO) evaluation for EESM17. Because ear-EEG channels exhibit different scale and variability, we apply instance normalization to EESM17 signals; otherwise we follow the CBraMod training setup.

3. Experiments

Variable pretraining layouts and sample efficiency. To understand how channel embedding techniques transfer between pretraining and downstream layouts, we construct a controlled toy transfer problem by pretraining each variant on 3 different electrode subset layouts from TUH: Standard, Frontal, Temporal. Each subset contains 6 electrodes for consistency in data quantity for the model. Precise electrode

selections are available in Appendix A. For each layout subset, we train 5 models from different random seeds to cover variability due to random initialization. Then, we evaluate each of these models on 5 random seeds of fine-tuning on the downstream dataset. The downstream dataset uses the Standard layout subset of electrodes, so 1 pretraining layout (Standard) perfectly matches the downstream dataset, while the other 2 (Frontal and Temporal) do not. We focus our sample efficiency sweep on using <1/10th (<200 out of 2k subjects) of the full TUAB training set. This setting tests layout transfer in a controlled, data constrained manner within the same dataset family between pretraining and downstream decoding.

Full layout pretraining. To investigate the real-world use case, we pretrain on the full TUH set of electrodes with 5 different random seeds. These models are then evaluated on a downstream dataset with widely different layouts, notably, EESM17 (Mikkelsen et al., 2017) which is a completely disjoint set of electrodes. Because EESM17 uses ear electrodes absent from TUH, this benchmark tests whether pretrained spatial structure can extrapolate beyond the scalp montage. We again use 5 fine-tuning seeds for downstream adaptation. This provides one of the most challenging layout transfer problems that one might face when trying to use pretrained models in the real-world, and serves as our main real-device transfer benchmark. Together, these settings let us study layout mismatch under controlled conditions and in a realistic cross-device transfer scenario.

4. Results

Variable pretraining layouts and sample efficiency. We first confirm that layout mismatch reduces performance: pretraining on Frontal or Temporal subsets and evaluating on Standard is worse than Standard-to-Standard transfer (Figure 2a). We also find that ACPE performs on par with or better than traditional channel embedding baselines (Ap-

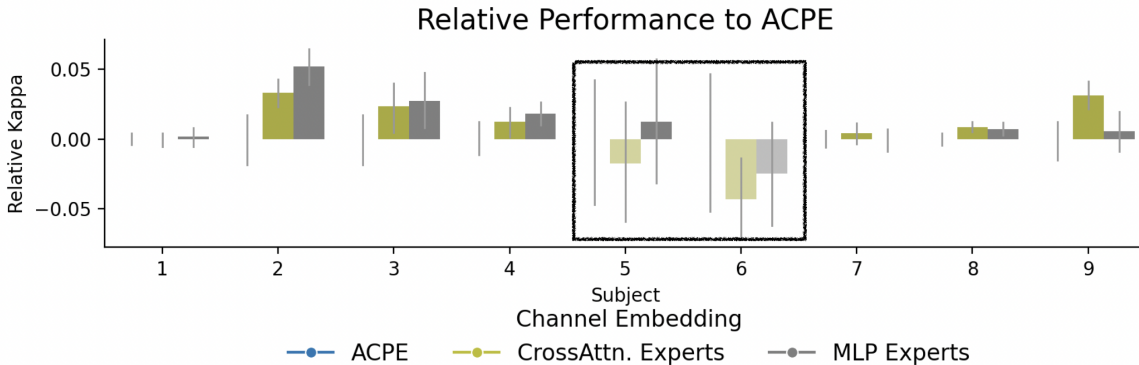


Figure 3. **Pretraining on Full TUH, fine-tune on ear EEG form factor for sleep staging classification.** We plot the delta performance between ACPE and our technique across subjects at N=6 training subjects as a focused comparison against the strongest prior learned baseline. We find that 7/9 subjects in leave-one-subject-out evaluations (x-axis) benefit from our technique (positive colored bars). The subjects that had worse performance with our method (box outline) were previously reported to have poor electrode contact (Mikkelsen et al., 2017) (subject note in Appendix D). Shown is mean delta of the labeled method against ACPE performance, with std across 5 pretraining x 5 fine-tuning seeds; absolute performance across all channel embedding techniques is provided in Appendix E.

pendix C), consistent with prior work. Importantly, our expert-based adapters are competitive with and often improve over ACPE under layout transfer (Figure 2bcd), indicating that additional spatial information learned during pretraining can be reused. This pattern is most visible when transferring to new layouts rather than when pretraining and fine-tuning use the same subset. In the controlled toy problem, Device Passport is most favorable in the mismatched-layout settings that motivate the paper, namely Temporal-to-Standard and Frontal-to-Standard transfer. As the number of fine-tuning subjects increases, performance improves across methods; ACPE tends to be quite strong at the very lowest-data regime, while Device Passport is stronger across most intermediate regimes, especially on layout transfer settings.

Full layout pretraining. In this realistic layout-transfer regime, Device Passport improves over the strongest prior learned baseline, ACPE, for most held-out subjects.

In transferring to completely new layouts (e.g. ear electrodes), Device Passport improves performance on EESM17 (Mikkelsen et al., 2017) relative to ACPE (Figure 3). The MLP expert variant improves over ACPE for 7/9 held-out subjects and ties on one additional subject, showing that the same qualitative advantage from the toy transfer setting persists in a realistic new-device transfer setting. In absolute terms, the MLP expert variant is best or second-best for all 9 EESM17 subjects in Appendix E, while the cross-attention variant is strongest on several subjects but less robust on the poor-contact cases. The subjects that reportedly had poor electrode contact (Subjects 5 (our 6) and 1 (our 5) in (Mikkelsen et al., 2017)) were hurt with our method, suggesting that robustness to noisy channels is a key limitation and opportunity for improvement. Taken together, these results suggest that metadata-guided expert embeddings can improve transfer to unseen layouts.

5. Discussion and Conclusion

We identify channel embeddings as an important bottleneck for reusing biosignal foundation models on unseen device layouts. In both controlled subset-transfer experiments and full-layout transfer to ear-EEG, Device Passport is most helpful when the downstream layout differs from pretraining, suggesting that pretrained spatial information can be reused more directly than with ordinary lookup or activity-only positional encodings. We focused on low fine-tuning sample regimes because they match the deployment setting for new biosignal devices: before large device-specific datasets exist, practitioners need pretrained models that adapt with limited labeled data.

These results support Device Passport in the layout-transfer settings we study, while leaving room for stronger robustness across all regimes. ACPE remains a competitive baseline, and Device Passport depends on useful channel metadata. Nevertheless, the pattern across controlled and real-device transfers suggests that metadata-guided expert embeddings provide a practical way to reuse pretrained spatial structure: downstream training can adapt the mixture rule without relearning channel representations from scratch. The clearest remaining limitation is sensitivity to noisy or poorly contacted channels, as seen in EESM17, which points to direct extensions such as synthetic channel corruption or explicit noise experts.

Overall, Device Passport offers a lightweight mechanism for adapting pretrained biosignal models to new health-sensing devices and montages with limited device-specific data. More broadly, these findings suggest that explicitly modeling how channel metadata and functional activity jointly define sensor identity can make biosignal foundation models more portable across real-world acquisition layouts.

References

- Abbaspourazad, S., Elachqar, O., Miller, A. C., Emrani, S., Nallasamy, U., and Shapiro, I. Large-scale training of foundation models for wearable biosignals. *arXiv preprint arXiv:2312.05409*, 2023.
- Azabou, M., Arora, V., Ganesh, V., Mao, X., Nachimuthu, S., Mendelson, M., Richards, B., Perich, M., Lajoie, G., and Dyer, E. A unified, scalable framework for neural population decoding. *Advances in Neural Information Processing Systems*, 36:44937–44956, 2023.
- Chau, G., Wang, C., Talukder, S., Subramaniam, V., Soedar-madji, S., Yue, Y., Katz, B., and Barbu, A. Population transformer: Learning population-level representations of neural activity. *ArXiv*, pp. arXiv–2406, 2025.
- Chu, X., Tian, Z., Zhang, B., Wang, X., and Shen, C. Conditional positional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021.
- Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C., Goj, R., Jas, M., Brooks, T., Parkkonen, L., et al. Meg and eeg data analysis with mne-python. *Frontiers in Neuroinformatics*, 7:267, 2013.
- Jiang, W.-B., Zhao, L.-M., and Lu, B.-L. Large brain model for learning generic representations with tremendous eeg data in bci. *arXiv preprint arXiv:2405.18765*, 2024.
- Kaifosh, P. and Reardon, T. R. A generic non-invasive neuromotor interface for human-computer interaction. *Nature*, 645(8081):702–711, 2025.
- Mikkelsen, K. B., Villadsen, D. B., Otto, M., and Kidmose, P. Automatic sleep staging using ear-eeg. *Biomedical engineering online*, 16(1):111, 2017.
- Obeid, I. and Picone, J. The temple university hospital eeg data corpus. *Frontiers in neuroscience*, 10:196, 2016.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Wang, J., Zhao, S., Luo, Z., Zhou, Y., Jiang, H., Li, S., Li, T., and Pan, G. Cbramod: A criss-cross brain foundation model for eeg decoding. *arXiv preprint arXiv:2412.07236*, 2024.

A. Decoding layouts

Electrode layouts. Figure 4 shows the exact channel subsets used for our pretraining–layout transfer experiments (Standard/Frontal/Temporal) and the downstream decoding layouts.

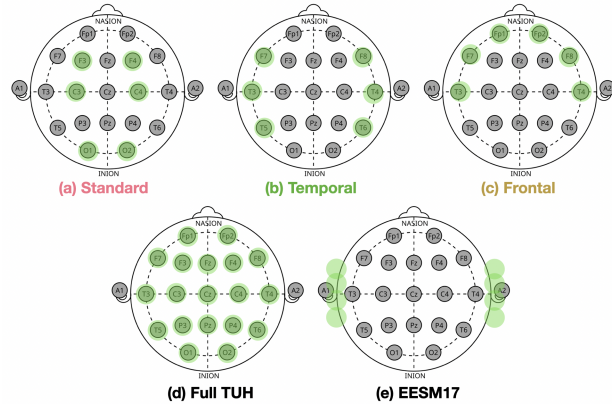


Figure 4. **Electrode layouts.** Electrode selections (green) used in our layout transfer experiments, including the TUH subset pretraining layouts (a-c), full pretraining (d), and downstream decoding ear-EEG (e) layouts.

B. Channel embedding details

APE. Absolute positional embeddings (APE) leverage sinusoidal positional encoding indexed by a channel’s ordinal number in the layout (Vaswani et al., 2017). This can misalign spatial information when the downstream layout reorders or replaces channels.

Channel ID. Learned channel embeddings are indexed by channel name (e.g., FP1) and are commonly used in neural foundation models (Azabou et al., 2023; Jiang et al., 2024). This works when the same 10–20 channel reappears at fine-tuning time, but new devices require randomly initialized embeddings.

XYZ. XYZ sinusoidal embeddings use electrode coordinates from MNE (meters) (Gramfort et al., 2013). We scale to millimeters, add 150 mm, round to the nearest millimeter, then index sinusoidal embeddings for each axis. Each axis contributes 66 dimensions (200/3), concatenated and zero-padded to 200 dimensions. This is effective with sufficient coordinate diversity (Chau et al., 2025) but may be limited under single-layout pretraining.

ACPE. Asymmetric conditional positional encoding (ACPE) is CBraMod’s channel-aware positional encoding scheme (Wang et al., 2024), adapted from conditional positional encoding (CPE) in vision (Chu et al., 2021). Instead of using fixed absolute embeddings, ACPE dynamically generates positional information from the patch embeddings

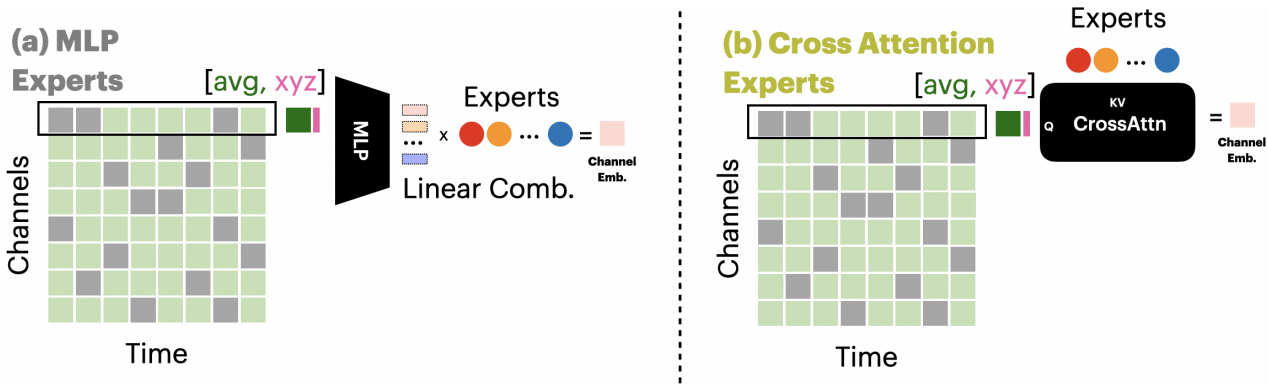


Figure 5. **Our method.** We learn (a) MLP or (b) cross-attention mechanisms that combine channel activity with xyz location to weight or attend to a set of learned experts.

themselves, making it more adaptable to varying channel layouts and time lengths. Concretely, ACPE applies a depth-wise 2D convolution over the spatial (channels) and temporal dimensions of the patch-embedding grid to produce positional encodings, which are then added back to the patch embeddings. The convolution kernel is asymmetric, using a larger spatial kernel than temporal kernel ($k_s > k_t$) to encode long-range spatial dependencies and shorter-range temporal dependencies, reflecting the structure of EEG signals. This dynamic, asymmetric design has been shown to outperform APE and symmetric CPE in CBraMod when transferring across EEG formats (Wang et al., 2024).

Expert adapters (MLP / Cross Attention). We learn a small set of expert embeddings during pretraining and freeze them for downstream use. A lightweight adapter uses channel activity plus xyz metadata to compute mixture weights over experts. The xyz inputs are kept in decimeter scale to match patch-embedding magnitudes. In the MLP variant, a softmax over experts produces the channel embedding. In the cross-attention variant, learned K/V embeddings are combined with a query projected from activity+metadata to produce the channel embedding. A visualization of our method can be found in Figure 5.

C. ACPE vs. traditional encodings

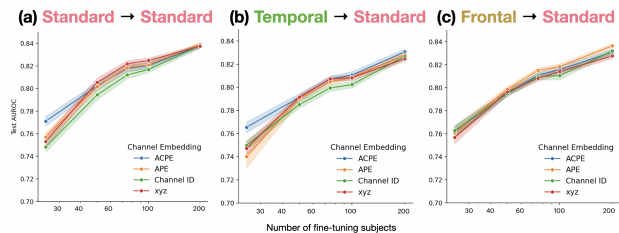


Figure 6. **ACPE compared with traditional positional encoding methods.** We compare ACPE against common channel embedding baselines (e.g., APE, Channel ID, and xyz-based encodings) in the variable-layout pretraining setting.

D. EESM17 Subject Note

The subject numbering used in our EESM17 plots follows our evaluation order rather than the numbering used in (Mikkelsen et al., 2017). In particular, the poor-contact subjects called out in the main text correspond to Subject 5 in (Mikkelsen et al., 2017) (our Subject 6) and Subject 1 in (Mikkelsen et al., 2017) (our Subject 5).

E. Sleep staging absolute performance

Table 1. Absolute Cohen’s kappa on sleep staging (EESM17) by subject for each channel embedding technique. Bold indicates best; underline indicates second best.

Model Name	Sub 1	Sub 2	Sub 3	Sub 4	Sub 5	Sub 6	Sub 7	Sub 8	Sub 9
ACPE	<u>0.742</u>	0.712	0.634	0.682	<u>0.325</u>	0.356	0.579	0.791	0.666
Cross-Attn. Experts	<u>0.742</u>	<u>0.746</u>	<u>0.657</u>	<u>0.694</u>	0.308	0.313	0.583	0.800	0.697
MLP Experts	0.743	0.764	0.661	0.700	0.338	<u>0.331</u>	0.579	0.798	0.671