

# VIBE: VISION TRANSFORMER BASED EXPERTS NETWORK FOR SSVEP DECODING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Steady-state visual evoked potential based brain-computer interfaces (SSVEP-BCIs) have attracted wide attention for their high information transfer rate (ITR) and non-invasiveness. However, existing deep learning methods for SSVEP-BCI decoding have reached a performance bottleneck, as they struggle to fully extract the complex neural signal features required for robust performance. Motivated by advances in vision and time series modeling, here we present a **VI**sion **TR**ansformer **B**ased **E**xpert network (VIBE), a multistage deep learning framework for SSVEP classification. VIBE integrates a Vision Transformer (ViT) module to generate rich spatiotemporal representations with data and network enhancement modules in a decoder for frequency recognition. We evaluate VIBE on two large benchmark datasets, including the Benchmark and the BETA dataset spanning 105 subjects. Notably, with just 0.4 seconds of stimulation, our VIBE achieves an ITR of 263.8 bits per minute (bpm) and 202.7 bpm on the Benchmark and BETA datasets, respectively. Experimental results demonstrate that VIBE consistently outperforms state-of-the-art baselines in offline experiments, highlighting its effectiveness as a high-performance decoding strategy for **individually calibrated** SSVEP-BCIs.

## 1 INTRODUCTION

Brain-computer interfaces (BCIs) provide a direct communication pathway between the brain and external devices, enabling interaction without relying on neuromuscular activity. BCIs have emerged as effective tools for augmentative communication and human-machine interaction, with broad potential applications ranging from neuroprosthesis (Willett et al., 2021; 2023) to the next-generation form of human-computer interaction (Gao et al., 2025). Among noninvasive paradigms, steady-state visual evoked potential based BCIs (SSVEP-BCIs) stand out for their non-invasiveness, high ITR, robustness, and scalability. SSVEPs are frequency-tagged neural responses that can be evoked by periodic visual stimulation, including flickering squares, reversing checkerboards, and moving gratings, and they are elicited over occipital cortex at the stimulation frequency and its harmonics. These frequency-tagged responses exhibit a high signal-to-noise ratio (SNR), enabling SSVEP-BCIs to implement high-speed spellers (Chen et al., 2015b), robotic control, and smart home systems. However, achieving high decoding accuracy under short time windows remains challenging in learning effective neural representations from noisy, data-constrained EEG recordings with complex spatiotemporal and spectral dynamics.

Advancements in SSVEP-BCI decoding have been driven by both traditional linear methods and emerging deep learning models. Early approaches such as canonical correlation analysis (CCA) (Bin et al., 2009) and its filter-bank extension (FBCCA) Chen et al. (2015a) established training-free plug-and-play frequency recognition, while subsequent training-based methods like task-related component analysis (TRCA) Nakanishi et al. (2017) and task-discriminant component analysis (TDCA) (Liu et al., 2021b) designed sophisticated spatial filters using individually calibrated data to significantly boost the decoding performance. However, these linear methods remain limited in capturing the nonlinear and hierarchical patterns of EEG signals. **To mitigate these issues, convolutional neural networks (CNNs) and, more recently, Transformer-based architectures have been explored for learning richer spatio-temporal representations from EEG, achieving improvements over traditional baselines (Li et al., 2020; Song et al., 2022). However, both families of models exhibit concrete gaps when applied to SSVEP decoding. CNNs depend on fixed, small receptive fields, which lim-**

its their ability to capture long-range temporal structure and harmonic relationships across wider time-scales features that are essential for differentiating densely spaced SSVEP frequencies. While Transformers are in principle capable of modeling global dependencies, existing works primarily deploy them in leave-one-subject-out (LOSO) evaluations, where broad cross-subject scenario is beneficial. In individually calibrated settings, however, this line of research often struggles to extract the fine-grained, local frequency-specific patterns required for high-precision decoding, and typically underperforms despite its expressive capacity.

To overcome these limitations, we introduce **VIBE** (Vision Transformer Based Experts Network), a multistage framework that unifies transformer-based sequence modeling with expert-driven specialization. VIBE employs a ViT module to perform temporal generation, expanding short input sequences into richer representations that preserve multi-scale temporal dependencies. It integrates a Mixture of Experts (MoE) decoder, where experts specialize in different subband-channel-temporal dynamics, and a load-balancing loss ensures diverse expert utilization for better generalization. On top of these architectural innovations, VIBE employs a staged training scheme that progressively pretrains and fine-tunes the ViT-based temporal generation and MoE-based decoding modules, adapting from population data to subject-specific dynamics. It further integrates data augmentation strategies, including temporal stitching, channel chunk shuffling, random temporal cropping, and decorrelation, to regularize training and enhance representation learning for EEG.

We evaluate VIBE on two large benchmark datasets, the Benchmark and the BETA datasets, spanning 105 subjects. Results show that VIBE consistently outperforms both classical and deep learning baselines, achieving **state-of-the-art** accuracy and ITR under short time windows. These findings establish VIBE as an effective decoding strategy for high-throughput SSVEP-BCIs.

In summary, our main contributions are threefold:

1. A novel hybrid framework that combines ViT-based temporal generation with MoE-based subband-channel-temporal decoding enhancement for SSVEP classification.
2. A staged training scheme that progressively adapts temporal generation and decoding modules from population to subject-specific data.
3. A suite of data augmentation methods designed for **SSVEP data**, improving robust representation learning from limited and noisy neural data.

## 2 RELATED WORK

**Traditional methods.** Early research focused on traditional frequency recognition methods, which can be broadly divided into training-free and training-based approaches. Canonical correlation analysis (CCA) (Bin et al., 2009) and its filter-bank extension (FBCCA) (Chen et al., 2015a) became widely used due to their plug-and-play training-free capability. Prior studies using individually calibrated data introduced multiset CCA (MsetCCA) (Zhang et al., 2014), L1-regularized multiway CCA (L1MCCA) (Zhang et al., 2013), and extended CCA (eCCA) (Nakanishi et al., 2014), which improved robustness by leveraging richer reference structures and regularization strategies. More sophisticated spatial filter based methods were developed to further boost accuracy. Task-related component analysis (TRCA) (Nakanishi et al., 2017) significantly improved SSVEP decoding by maximizing trial-to-trial reproducibility. To further address redundancy in TRCA’s ensemble design, task-discriminant component analysis (TDCA) (Liu et al., 2021b) eliminated the training of spatial filters class by class and leveraged spatio-temporal neural dynamics, making it a state-of-the-art method for enabling high-speed brain spellers. Despite these advances, traditional methods remain linear and limited in their ability to capture nonlinear, hierarchical representations of EEG.

**Deep learning methods.** Motivated by these limitations, recent research has turned toward learning complex representations from noisy signals in an end-to-end manner. Convolutional neural networks (CNNs) enabled data-driven feature extraction, analogous to filtering in EEG signal processing, and advances such as convolutional correlation analysis (ConvCA) (Li et al., 2020) and deep neural network classifiers (Guney et al., 2022) have surpassed linear baselines. Extensions incorporated fixed and dynamic template networks (Xiao et al., 2022), bidirectional Siamese correlation networks (Zhang et al., 2022), and multiscale CNNs with squeeze-and-excitation blocks (Jin et al., 2024). More recently, Transformer-based architectures leveraging self-attention to capture long-range temporal dependencies have been applied to EEG (Song et al., 2022; Wan et al., 2023),

including SSVEPformer (Chen et al., 2023), DG-Conformer (Liu et al., 2024), SSVEPPoolformer (Li et al., 2025a) and MTSNet (Lan et al., 2025) for cross-subject SSVEP classification. Hybrid approaches such as TRCA-Net (Deng et al., 2023) and discriminant compacted network (Li et al., 2025b) combine spatial filters with neural networks, while ConsenNet (Zhang et al., 2024) leverage a teacher-student framework to further improve performance. Most recently, Mamba-based models such as SUMamba (Dong et al., 2026) integrated multi-scale feature fusion to facilitate classification. However, CNNs remain limited in capturing global dependencies, and Transformers often neglect inductive biases specific to EEG, leaving the representation learning problem unresolved for high-throughput SSVEP-BCIs.

### 3 METHOD

We first define the notation used throughout this work. The multi-channel EEG signal is represented as  $X \in \mathbb{R}^{B \times C \times T}$ , where  $B$  denotes the number of filter banks,  $C$  the number of EEG electrodes (channels), and  $T$  the total number of sampled timestamps. In our experiments, we consider  $B = 3$  sub-bands extracted by band-pass filtering, and  $C = 9$  channels selected from classical montage for SSVEP classification (Chen et al., 2015b).

Fig. 1 is not so auto-explicative. This reviewer believe that there is too much text in it, without a clear indication about where to start looking at. Moreover, the caption only talks about points a) and e), but not b-d). I suggest redesigning the figure and its caption entirely.

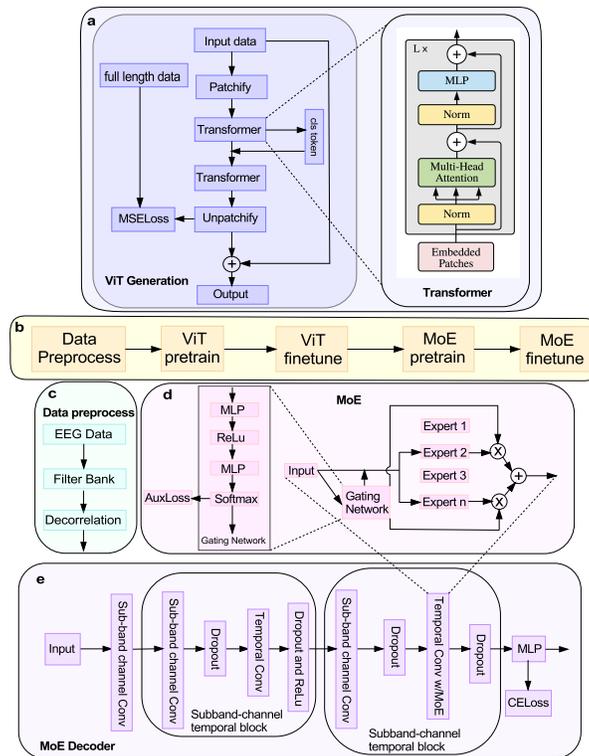


Figure 1: Framework of the proposed VIBE method. Panel (b) presents the overall pipeline, from short SSVEP EEG segments to final classification. Panel (c) depicts the data preparation strategy used during training and testing. Temporal sequence generation with the ViT is detailed in (a), identical for both pretrain and fine-tune stages. The MoE decoder is shown in (e), also identical for both pretrain and fine-tune stages. Panel (d) illustrates the final convolutional temporal layer integrated with the MoE, highlighting how expert activations contribute to the classification output.

### 3.1 ViT-BASED TEMPORAL GENERATION.

The first component of our framework is a ViT (Dosovitskiy et al., 2021) adapted for temporal sequence generation. We represent the multi-band EEG data as an image tensor  $X \in \mathbb{R}^{B \times C \times T_{\text{in}}}$ , where the three dimensions correspond to sub-bands ( $B = 3$ ), channels ( $C$ ), and time samples ( $T_{\text{in}}$ ). Following the ViT formulation, the signal is partitioned into non-overlapping patches of size  $(B, C, t)$ , where  $t$  denotes a small temporal window (e.g.,  $t = 10$  samples  $\approx 0.04$  s). Each patch is then flattened and linearly projected into a latent embedding, forming a sequence of tokens. Like standard ViTs, we include positional embeddings added to the patchified embedding. A learnable class token is attached to the sequence, which serves to expand the representation to match the target output length  $T_{\text{out}}$ . This latent sequence is then processed by a transformer-based decoder (also implemented as a ViT), and the resulting patches are unpatchified to reconstruct the output sequence  $\hat{X} \in \mathbb{R}^{B \times C \times T_{\text{out}}}$ , with  $T_{\text{out}} > T_{\text{in}}$ . The model was optimized using the mean squared error (MSE) loss. We clip the extended sequence of shape  $\mathbb{R}^{B \times C \times (T_{\text{out}} - T_{\text{in}})}$ , which represents the newly generated temporal samples. This generated segment was concatenated with the original input  $X$ , forming the final output representation.

### 3.2 MOE WITH CHANNEL-SUBBAND TEMPORAL DECODER

The second part of our framework is a decoder that jointly models subband, channel, and temporal dependencies using a MoE design.

#### 3.2.1 MOE

MoE (Shazeer et al., 2017) mechanism is designed to increase model capacity while keeping computational cost manageable through sparse activation. MoE has also been successfully applied to EEG decoding tasks (Yang et al., 2025). Instead of applying a single shared transformation to all inputs, an MoE layer maintains a set of  $E$  experts  $\{f_1, f_2, \dots, f_E\}$ , each parameterized as a learnable function (e.g., convolutional filters in our case). For each input token  $x$ , a gating network produces a probability distribution over experts, and only the top- $k$  experts are selected to process the input. Specifically, we take the mean over temporal domain as the input of the gating network, where the network consist of two layers of MLP and one ReLU as activation. The final output is then obtained as a weighted combination of the selected experts’ outputs, where the gating scores act as mixture coefficients. This strategy enables different experts to specialize on distinct temporal or spectral patterns in the EEG signal, enhancing both representation power and generalization.

**Auxiliary Load-Balancing Loss.** To encourage balanced utilization of experts and prevent overfitting, we introduce an auxiliary load-balancing loss. For each MoE layer, the gating network computes a probability distribution over  $E$  experts for each input token. Let  $G \in \mathbb{R}^{B \times E}$  denote the gate probabilities for a batch of  $B$  inputs, with  $G_{ij}$  representing the probability of assigning input  $i$  to expert  $j$ . The mean usage of each expert is then  $\bar{u}_j = \frac{1}{B} \sum_{i=1}^B G_{ij}$ . We define the auxiliary loss as the Kullback-Leibler (KL) divergence between the mean expert usage and a uniform distribution:

$$\mathcal{L}_{\text{aux}} = \lambda_{\text{aux}} \text{KL}(\log(\bar{u}) \| U)$$

where  $U$  is a uniform vector of length  $E$ , and  $\lambda_{\text{aux}}$  is a weighting coefficient. This loss encourages all experts to be used approximately equally, preventing collapse onto a small subset of experts and improving generalization.

#### 3.2.2 CHANNEL-SUBBAND FUSED TEMPORAL DECODER

The decoder first integrates information across sub-bands and channels by treating the data as a combined subband-channel dimension of shape  $(B * C, 1, T)$ . A convolutional layer with kernel  $(1, 1)$  is applied to extract higher-level spectral-channel features. Next, a subband-channel-temporal block is applied that consists of two convolutional layers: one over the subband-channel dimension and one over the temporal dimension. This block pools information separately from the spectral-channel and temporal dimensions. The block is repeated twice to progressively capture richer patterns across

channels, sub-bands, and time. Nonlinearities and dropout are applied between layers for regularization. Finally, the extracted features are flattened and passed through a fully connected layer to produce the class logits. In our model, we replace the last temporal convolution layer with a MoE layer, where experts in the networks employ the original temporal conv.

### 3.3 DATA AUGMENTATION

In the following section, the dataset is denoted as  $X \in \mathbb{R}^{S \times B \times N \times M \times T \times C}$ , where  $S$  is the number of subjects,  $B$  the number of subbands,  $N$  the number of trials,  $M$  the number of targets,  $T$  the temporal length, and  $C$  the number of EEG channels.

**Cross-Subject Temporal Stitching** As a generalization of (Lotte, 2015), the EEG signals are first divided into short temporal segments, or *chunks*, across the time dimension. For each chunk, we randomly select a segment of the same duration from any subject and trial in the dataset, preserving the original subband, channel, and target labels. The selected segments are combined across time to form a new synthetic trial, maintaining the original subband, channel, and target structure. This approach allows the creation of entirely new temporal patterns by sampling from different subjects and trials, rather than modifying the original trial. Namely,  $\oplus_i X_{s_i, :, n_i, :, i\tau : \tau(i+1), :} \in \mathbb{R}^{1 \times B \times 1 \times M \times T \times C}$  is a generated piece of data of, where  $\{0, \tau, 2\tau, \dots\}$  is the time chunk sequence and each  $s_i, n_i$  is randomly selected among  $S, N$ ,

**Channel Chunk Shuffle** Given a random subject, input data is first divided into consecutive chunks along the time dimension. For each chunk, with a certain probability, two channels are randomly selected and swapped, while all other dimensions—including subbands, trial, target labels remain unchanged. Explicitly, given a subject’s trial and target, if time chunks  $(t_i\tau, t_j\tau)$  are selected to shuffle with  $(c_{i_0}, c_{i_1}), (c_{j_0}, c_{j_1}) \in S_c$  are transpositions corresponding to  $t_i, t_j$  as channel swap,  $X_{s, :, n, m, t_i\tau : t_{i+1}\tau, c_{i_0}}$  replaces the channel  $c_{i_1}$  and similarly for  $t_j, (c_{j_0}, c_{j_1})$ .

**Random Temporal Crop** *Gaze-shift and stimulus-locking latencies can vary across individuals in practical systems. Such variability has been extensively documented in the SSVEP literature (Liu et al., 2021b; Wang et al., 2016; Pan et al., 2011; Lemm et al., 2005; Dornhege et al., 2006; Wu et al., 2008; Qi et al., 2015)* Therefore, inspired by (Liu et al., 2021b), a *Random Temporal Crop (RTC)* augmentation is utilized to increase temporal diversity in the training data. For any chosen subject, we preserve the original trial, target label, subband, and channel structure. With a given probability, we randomly select a short segment of time (e.g. 0.03s) from the data, keeping only the latter portion and discarding the former. The cropped segment is zero-padded at the end to restore the trial to its original temporal length.

**Channel Decorrelation** We adopt a covariance-based whitening procedure across channels, conditioned on each subject and subband. For each subject and subband, we first compute the mean trial across all training trials to obtain a representation. This mean trial is used to estimate the channel covariance, from which a whitening matrix is derived (He & Wu, 2019; Liu et al., 2021a). The whitening matrix is applied to both training and test data, effectively reducing subject, trial-level variability while preserving the temporal and target-related structure of the signals. The decorrelation procedure emphasizes stable patterns across different subject and subband.

### 3.4 TRANSFER LEARNING.

Following the transfer learning strategy of (Guney et al., 2022), we adopt a staged training procedure to strengthen representation ability. Our model comprises two main components: a ViT encoder for temporal length generation and a MoE-based decoder for subband-channel and temporal integration. Thus, the transfer learning process is added to these component.

In the first stage, the ViT encoder is trained in a generative manner, reconstructing the temporal sequence from shorter inputs. In the second stage, this encoder is fine-tuned separately for each subject, where the global model parameters are re-initialized and adapted using only subject-specific data. The decoder is trained in a similar two-step fashion: first, a global MoE decoder is optimized using the pooled training data across all subjects, and subsequently, a subject-specific fine-tuning step is applied to adapt the decoder to individual variability.

## 4 EXPERIMENTS

### 4.1 DATASET

The experiments were carried out on two public 40-target SSVEP datasets: the Benchmark dataset (Wang et al., 2016) and the BETA dataset (Liu et al., 2020). Both datasets employed the joint frequency and phase modulation (JFPM) method to encode target stimuli. The data acquisition equipment for the Benchmark and BETA datasets is identical; however, the Benchmark dataset was collected in a controlled laboratory environment within an electromagnetic shielding room, whereas the BETA dataset was recorded in a more naturalistic setting, reflecting real-world conditions. All experiments, including comparisons with state-of-the-art methods, were performed on these two datasets. This allows us to evaluate the performance of our decoding approach under both controlled and realistic acquisition conditions.

### 4.2 PREPROCESSING

The same preprocessing pipeline was applied to both datasets. Nine electrodes (Pz, PO5, PO3, POz, PO4, PO6, O1, Oz, and O2) were selected for analysis. The EEG signals were downsampled to 250 Hz. To account for visual response latency, we considered delays of 0.14 s for Benchmark and 0.13s for BETA, consistent with previous studies (Chen et al., 2015b). For each trial, data segments of length  $t$  seconds were extracted in the time windows  $[0.14, 0.14 + t]$  s and  $[0.13, 0.13 + t]$  s after stimulus onset for Benchmark and BETA, respectively.

We apply a filter-bank approach as a preprocessing step to enhance SSVEP signals (Chen et al., 2015a). Data passes through three band-pass filters with frequency ranges  $(8N, 90)$  Hz, where  $N = 1, 2, 3$ , and filtered signals are concatenated along the sub-band dimension. This procedure captures multiple harmonics and improves the signal representation for subsequent decoding.

### 4.3 BASELINE MODELS

**Deep Learning Models.** **DNN** (Guney et al., 2022) is a dense convolutional neural network that processes time-series data and incorporates a fine-tuning stage to boost performance. **SSVEPformer** (Chen et al., 2023) is a transformer-based neural network that takes complex spectra as input, leveraging a transformer encoder and fully connected layer to extract phase and frequency features. **TR-CANet** (Deng et al., 2023) applies TRCA-based spatial filtering to the input data, followed by a DNN for feature learning.

**Traditional Models.** **TDCA** (Liu et al., 2021b) addresses the redundancy of stimulus-specific spatial filters in TRCA and the underutilization of temporal information. It enhances the performance of individually calibrated SSVEP-BCIs by learning task-discriminative spatiotemporal components. **TRCA** (Nakanishi et al., 2017) derives spatial filters by maximizing SSVEP reproducibility across trials, while **eTRCA** extends this by ensembling filters across all frequencies. **eCCA** (Nakanishi et al., 2014) introduces a combination of spatial filters derived from canonical correlation analysis (CCA) and employs a user-specific target identification algorithm based on individual calibration data. **msTRCA** (Wong et al., 2020) extends TRCA with a multi-stimulus learning scheme that leverages data from both target and non-target stimuli.

### 4.4 EXPERIMENTAL SETUP

We employed  $k$ -fold cross-validation, with  $k = 6$  for Benchmark and  $k = 4$  for BETA. For each subject, one block of EEG data was designated as the test set, while the remaining blocks were used for training within that fold. All training follows the four-stage procedure: ViT generative pretraining, ViT subject-specific fine-tuning, MoE decoder pretraining, and MoE subject-specific fine-tuning. Further implementation details are provided in the relevant subsectionA of the Appendix.

## 5 RESULT

To evaluate the performance of algorithms among different data lengths, we report both classification accuracy and ITR. The ITR, measured in **bits per minute (bpm)**, is defined as (Wolpaw et al., 2002):

$$ITR(P, T, M) = \left( \log_2 M + P \log_2 P + (1 - P) \log_2 \frac{1 - P}{M - 1} \right) \frac{60}{T} \quad (1)$$

Here,  $M$  denotes the number of target classes,  $P$  denotes the classification accuracy, and  $T$  (in seconds) represents the total selection duration, including gaze time and a fixed gaze shift of 0.5 s.

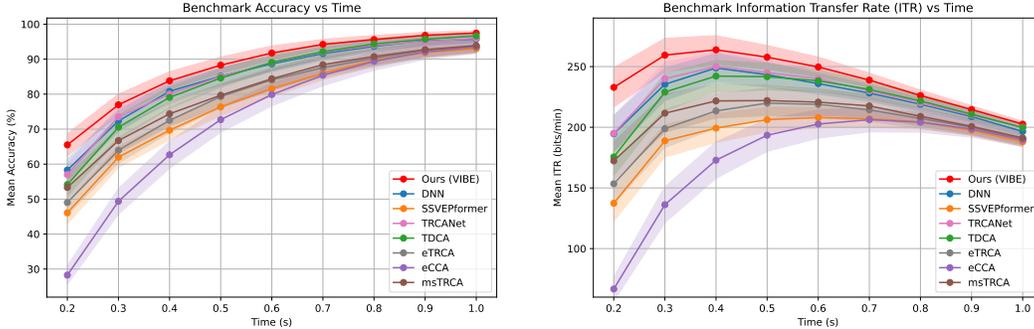


Figure 2: The left panel shows the mean classification accuracy, and the right panel shows the mean information transfer rate (ITR) across all 35 subjects in the Benchmark dataset. Shaded regions indicate the standard errors for subjects.

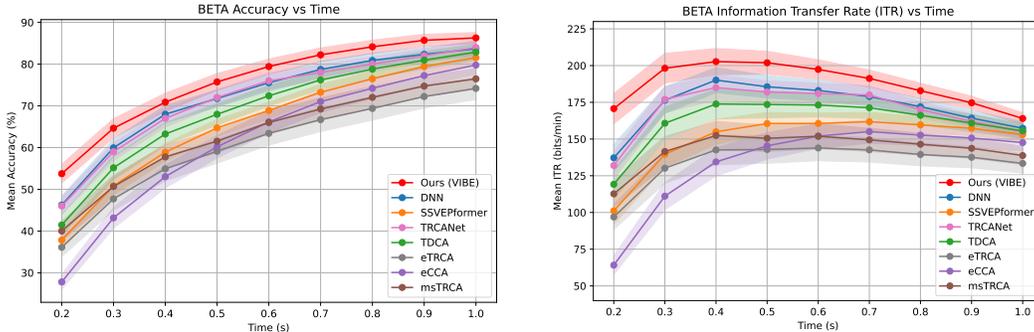


Figure 3: The left panel shows the mean classification accuracy, and the right panel shows the mean ITR across all 70 subjects in the BETA dataset. Shaded regions indicate the standard errors for subjects.

Figures 2 and 3 present the average classification accuracy and ITR of the proposed VIBE network evaluated on Benchmark and BETA across different data lengths. At the shortest data length (0.2 s), VIBE achieved the largest accuracy advantage over all other methods, highlighting its superior capability for rapid SSVEP decoding (Benchmark: 65.5% vs. 58.8%; BETA: 54.1% vs. 46.2%). The maximum ITR for VIBE was observed at 0.4 s, reaching  $263.8 \pm 11.7$  bpm for Benchmark and  $202.7 \pm 8.9$  bpm for BETA, exceeding the corresponding values of the DNN baseline ( $248.8 \pm 11.8$  bpm and  $190.1 \pm 8.4$  bpm, respectively). For a 1 s data length, VIBE maintained superior classification performance compared with DNN (Benchmark:  $97.4 \pm 0.7\%$  vs.  $95.7 \pm 1.1\%$ ; BETA:  $86.3 \pm 1.3\%$  vs.  $83.7 \pm 1.6\%$ ). Collectively, these results demonstrated that VIBE effectively decodes SSVEP responses across a range of time windows, with particularly pronounced benefits under short observation periods.

The performance of each method was evaluated in terms of decoding accuracy and ITR across data lengths. A two-way (method  $\times$  data length) repeated-measures ANOVA (Greenhouse–Geisser

corrected) revealed a statistically significant interaction between method and data length for both datasets (Benchmark:  $F(56, 1904)$ ,  $p < 0.001$ ; BETA:  $F(56, 3864)$ ,  $p < 0.001$ ). The eight methods included in this analysis correspond to those described in the Experiments Section. The detailed results are provided in Appendix Table 7. These findings indicate that the effect of data length on decoding performance depended on the method used, and vice versa, highlighting significant differences in performance trends across methods and data lengths.

For both Benchmark and BETA, paired t-tests revealed that our proposed VIBE method achieved significantly higher decoding accuracies than the deep learning baseline (DNN) and the traditional method (TDCA) across all evaluated data lengths (all:  $p < 0.05$ ). The details of these results are summarized in Table 8, in Appendix. The advantage of VIBE was especially pronounced at short data lengths (e.g., 0.2 s, Benchmark: VIBE vs. DNN:  $p = 1.0 \times 10^{-12}$ ; VIBE vs. TDCA:  $p = 2.2 \times 10^{-14}$ ; BETA: VIBE vs. DNN:  $p = 4.8 \times 10^{-18}$ ; VIBE vs. TDCA:  $p = 3.3 \times 10^{-25}$ ), demonstrating that our method was more robust under very short EEG segments. As data length increases, all methods converged towards similar performance, but VIBE consistently maintained a significant edge, indicating its effectiveness in both short- and long-window SSVEP decoding. **After applying the Holm-Bonferroni correction to control for multiple comparisons, all paired t-test results remained statistically significant ( $p < 0.05$ ), indicating that VIBE consistently outperforms the compared methods across data lengths and datasets.**

## 5.1 ABLATION STUDY

To better understand the contribution of each component in our framework, we conducted an ablation study on the Benchmark dataset with a 0.2 s data length. A brief summary of the ablation is presented in Table 10. The original model achieved an accuracy of **65.5%**.

**MoE** Removing the MoE module resulted in an accuracy of **64.2%**, highlighting its importance. Further analysis of MoE placement across different layers is provided in Appendix D.1.

**ViT regeneration.** Removing the ViT regeneration module led to a performance drop to **61.8%**, highlighting its essential role in feature representation. Further exploration of the effect of varying ViT generation time length is provided in Appendix D.3, which indicated that the optimal generation time depended on the input trial length: longer generation times benefited short trials, while shorter generation times were preferable for longer trials.

**Data augmentation.** Two augmentation strategies were employed: a decorrelation-based augmentation and an additional data generation module (the three methods described in Section 3.3) for the MoE decoder. When only the decoder-specific augmentation was removed, the accuracy decreased to **63.5%**; when both strategies were removed, the accuracy further dropped to **62.1%**. Further analysis of the effect of removing each augmentation is provided in Appendix D.2.

These results confirm that each module contributes positively, with all these three blocks play critical roles, and that the optimal ViT generation time length is data-length dependent.

## 5.2 FEATURE VISUALIZATION VIA T-SNE

To explore the reasons behind the superior performance of our model, we applied t-distributed Stochastic Neighbor Embedding (t-SNE) (Maaten & Hinton, 2008) to visualize the features learned from the final fully connected layer. We examined only our model, comparing the full version with an ablated version in which the ViT, data augmentation, and MoE were removed, using a data length of 0.6 s. In Appendix Figure 4, dots of the same color in the full model (left) formed more compact and dense clusters than in the ablated model (right). The circled clusters highlighted representative examples. This increased density indicated that the three key modules contributed to generating more discriminative and tightly grouped feature representations.

## 5.3 SUBJECT-WISE ITR VISUALIZATION

In Appendix Figure 5, we visualized the ITR for each subject using radar plots to compare our model (VIBE) with baseline methods (DNN and TDCA) at a data length of 0.2 s. Four radar plots were presented, corresponding to the Benchmark and BETA datasets, and comparing VIBE with DNN and TDCA, respectively. Each spoke in the radar plot represented an individual subject, and the distance

432 from the center indicated the ITR value. Across both datasets, VIBE consistently achieved higher  
433 ITR values for most subjects compared to the baseline models, illustrating its superior performance  
434 and robustness in short-duration SSVEP decoding.

## 436 6 DISCUSSION

### 437 6.1 NEURAL UNDERPINNINGS OF THE PROPOSED MODULES

438 In our study, the effectiveness of ViT-based regeneration can be attributed to the temporal nature  
439 of the SSVEP signals. Thus, the regenerated segments effectively extended the temporal window  
440 available to the decoder, which provided richer frequency-level information (target frequencies lie  
441 within the 8-15.8 Hz range). The regeneration step ensured that the decoder could access more  
442 complete frequency cycles, especially when the original data length was short. Although each ViT  
443 patch embedding only encoded a tiny fraction of data length (e.g. 0.04 s), it preserved additional  
444 temporal information for the decoder to facilitate classification.

445 Data augmentation plays a crucial role in enhancing the robustness of the model by introducing  
446 variability and simulating real-world scenarios. Several techniques have been implemented in this  
447 study, each inspired by physiological and contextual considerations related to EEG signals. First,  
448 the random temporal crop augmentation addresses inter-subject and task-dependent variability in  
449 SSVEP latency. While the standard latency is incorporated in the data preprocessing pipeline, the  
450 actual latency for each individual can differ, so this augmentation randomly samples temporal seg-  
451 ments within each trial to learn latency-tolerant features rather than overfitting to a fixed window,  
452 improving generalization. Second, the channel chunk shuffle augmentation is motivated by the  
453 dipole-source origin of EEG and distortions from volume conduction and other artifacts, and it ran-  
454 domly shuffles chunks of channels to simulate varied electrode placements and signal quality. This  
455 promotes invariance to sensor positioning and improves generalization across hardware setups and  
456 individuals. Third, cross-subject temporal stitching encourages the decoder to focus on frequency-  
457 level information rather than subject-specific features by stitching trials across subjects, exposing  
458 it to diverse temporal patterns and yielding generalized frequency responses that reflect underlying  
459 physiology. **This increased diversity of neural signals mitigates overfitting to individual trials and  
460 promotes more robust class-specific representations, ultimately improving the model’s generaliza-  
461 tion performance.**

462 The MoE mechanism is particularly valuable in the final temporal convolution layer of the model,  
463 where different experts can specialize in learning distinct temporal patterns relevant to specific task  
464 targets. Some experts may focus on shorter, rapid temporal responses, while others may specialize  
465 in longer, more sustained patterns, enabling the model to better capture the full range of temporal  
466 dynamics. This adaptability allows the model to allocate different experts to process different parts  
467 of the temporal signal.

### 468 6.2 TRAINING AND TESTING TIME ANALYSIS

469 For VIBE, the two pretraining stages were performed using data from all subjects (excluding test  
470 data), while the fine-tuning stages employed data from a single subject. Table 9 in the Appendix  
471 summarizes the training times for each stage and the testing time for a single 0.4 s trial, with all  
472 experiments conducted on an NVIDIA RTX 4090 GPU. The pretraining stages accounted for the  
473 majority of the training time, whereas fine-tuning for a specific subject could be completed in ap-  
474 proximately 17 seconds for BETA and 1 minute for Benchmark. The difference in training time  
475 between the two datasets is due to the differing number of epochs in each stage. Testing a single  
476 trial required less than 1 ms, which is negligible compared to the data duration. These findings indi-  
477 cate that VIBE provides a practical and efficient solution for SSVEP decoding in BCI applications.

### 480 6.3 LIMITATION AND FUTURE DIRECTION

481 One limitation of this work is that certain subjects exhibit performance that deviates markedly from  
482 the overall distribution, underscoring the need for more generalized approaches capable of handling  
483 inter-subject variability. Future investigations could therefore benefit from conducting experiments  
484 in alternative evaluation settings, such as performing cross-validation across subjects rather than  
485

486 **individually calibrated scenario**, or evaluate on other EEG decoding tasks (Song et al., 2024; Jiang  
487 et al., 2024; Wang et al., 2023), to provide a more rigorous assessment of generalization. Finally,  
488 an important direction for future research is the implementation of online experiments, wherein  
489 new patients are directly evaluated, to provide a realistic assessment of the model’s effectiveness in  
490 practical BCI applications.

## 491 492 7 CONCLUSION

493  
494 **In this study, we proposed the Vision Transformer Based Expert network (VIBE), a multistage deep**  
495 **learning framework that integrates a ViT-MoE architecture and novel data enhancement approaches**  
496 **tailored for SSVEP data. By leveraging information from short-duration EEG recordings, VIBE**  
497 **learns effective and discriminative neurophysiological representations for individually calibrated**  
498 **SSVEP decoding.** Evaluations on two benchmark datasets demonstrate that VIBE significantly im-  
499 proves both decoding accuracy and ITR. Overall, these results establish VIBE as a strong candidate  
500 for SSVEP decoding and support continued progress in **high-speed** BCI research.

## 501 502 8 REPRODUCIBILITY STATEMENT

503  
504 All datasets used in this work are publicly available and open-sourced. To facilitate reproducibility,  
505 we provide the complete code for our models and experiments alongside the submission **in the**  
506 **supplementary material.** Detailed descriptions of model architectures, training procedures, and data  
507 preprocessing steps are included in the main text, Appendix, ensuring that independent researchers  
508 can replicate our results.

## 509 510 9 ETHICS STATEMENT

511  
512 This work adheres to the ICLR Code of Ethics. All datasets used are publicly available and open-  
513 sourced. Specifically, the Benchmark and BETA datasets were collected under protocols approved  
514 by the respective institutions; For example, the BETA dataset protocol was approved by the Ethics  
515 Committee of Tsinghua University (No. 20190002) as reported in the original publication. No  
516 additional human subjects were involved in this study. The study focuses on computational modeling  
517 and analysis, without potential for harmful applications. All authors have read and complied with  
518 the ICLR Code of Ethics.

## 519 520 REFERENCES

- 521  
522 Guangyu Bin, Xiaorong Gao, Zheng Yan, Bo Hong, and Shangkai Gao. An online multi-channel  
523 ssvep-based brain–computer interface using a canonical correlation analysis method. *Journal of*  
524 *neural engineering*, 6(4):046002, 2009.
- 525 Jianbo Chen, Yangsong Zhang, Yudong Pan, Peng Xu, and Cuntai Guan. A transformer-based deep  
526 neural network model for ssvep classification. *Neural Networks*, 164:521–534, 2023.
- 527 Xiaogang Chen, Yijun Wang, Shangkai Gao, Tzyy-Ping Jung, and Xiaorong Gao. Filter bank canon-  
528 ical correlation analysis for implementing a high-speed ssvep-based brain–computer interface.  
529 *Journal of neural engineering*, 12(4):046008, 2015a.
- 530 Xiaogang Chen, Yijun Wang, Masaki Nakanishi, Xiaorong Gao, Tzyyping Jung, and Shangkai Gao.  
531 High-speed spelling with a noninvasive brain-computer interface. *Proceedings of the National*  
532 *Academy of Sciences of the United States of America*, 112(44):201508080, 2015b.
- 533 Yang Deng, Qingyu Sun, Ce Wang, Yijun Wang, and S Kevin Zhou. Trca-net: using trca filters to  
534 boost the ssvep classification with convolutional neural network. *Journal of Neural Engineering*,  
535 20(4):046005, 2023.
- 536 Liuyuan Dong, Chengzhi Xu, Xuyang Wang, Ruizhen Xie, Guangbo Lei, Yimeng Li, and Wanli  
537 Yang. Sumamba: A mamba-based deep learning model with multi-scale feature fusion for ssvep  
538 classification. *Biomedical Signal Processing and Control*, 112:108376, 2026.

- 540 Guido Dornhege, Benjamin Blankertz, Matthias Krauledat, Florian Losch, Gabriel Curio, and K-  
541 R Muller. Combined optimization of spatial and temporal filters for improving brain-computer  
542 interfacing. *IEEE transactions on biomedical engineering*, 53(11):2274–2281, 2006.
- 543  
544 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
545 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An im-  
546 age is worth 16x16 words: Transformers for image recognition at scale. In *International Confer-*  
547 *ence on Learning Representations*, 2021.
- 548 Xiaorong Gao, Yijun Wang, Xiaogang Chen, Bingchuan Liu, and Shangkai Gao. Brain–computer  
549 interface—a brain-in-the-loop communication system. *Proceedings of the IEEE*, 113(5):478–511,  
550 2025. doi: 10.1109/JPROC.2025.3600389.
- 551 Osman Berke Guney, Muhtasham Oblokulov, and Huseyin Ozkan. A deep neural network for ssvep-  
552 based brain-computer interfaces. *IEEE transactions on biomedical engineering*, 69(2):932–944,  
553 2022.
- 554  
555 He He and Dongrui Wu. Transfer learning for brain–computer interfaces: A euclidean space data  
556 alignment approach. *IEEE Transactions on Biomedical Engineering*, 67(2):399–410, 2019.
- 557  
558 Wei-Bang Jiang, Li-Ming Zhao, and Bao-Liang Lu. Large brain model for learning generic represen-  
559 tations with tremendous eeg data in bci. In *International Conference on Learning Representations*,  
560 2024.
- 561 Jing Jin, Xiao Wu, Ian Daly, Weijie Chen, Xinjie He, Xingyu Wang, and Andrzej Cichocki. Squeeze  
562 and excitation-based multiscale cnn for classification of steady-state visual evoked potentials.  
563 *IEEE Internet of Things Journal*, 2024.
- 564  
565 Zhen Lan, Zixing Li, Chao Yan, Xiaojia Xiang, Dengqing Tang, Min Wu, and Zhenghua Chen.  
566 Mtsnet: Convolution-based transformer network with multi-scale temporal-spectral feature fusion  
567 for ssvep signal decoding. *IEEE Journal of Biomedical and Health Informatics*, 2025.
- 568  
569 Steven Lemm, Benjamin Blankertz, Gabriel Curio, and K-R Muller. Spatio-spectral filters for im-  
570 proving the classification of single trial eeg. *IEEE transactions on biomedical engineering*, 52(9):  
1541–1548, 2005.
- 571  
572 Chunquan Li, Zhiyuan Liao, Yuxin Cheng, Zitao Wang, Junyun Wu, Ruijun Liu, and Peter X  
573 Liu. Ssveppoolformer: An improved poolformer model with the adaptive denoising algorithm  
574 for ssvep-eeg signal classification. *IEEE Transactions on Consumer Electronics*, 2025a.
- 575  
576 Dian Li, Yongzhi Huang, Ruixin Luo, Lingjie Zhao, Xiaolin Xiao, Kun Wang, Weibo Yi, Minpeng  
577 Xu, and Dong Ming. Enhancing detection of ssveps using discriminant compacted network.  
*Journal of Neural Engineering*, 22(1):016043, 2025b.
- 578  
579 Yao Li, Jiayi Xiang, and Thenkurussi Kesavadas. Convolutional correlation analysis for enhancing  
580 the performance of ssvep-based brain-computer interface. *IEEE Transactions on Neural Systems*  
*and Rehabilitation Engineering*, 28(12):2681–2690, 2020.
- 581  
582 Bingchuan Liu, Xiaoshan Huang, Yijun Wang, Xiaogang Chen, and Xiaorong Gao. Beta: A large  
583 benchmark database toward ssvep-bci application. *Frontiers in neuroscience*, 14:627, 2020.
- 584  
585 Bingchuan Liu, Xiaogang Chen, Xiang Li, Yijun Wang, Xiaorong Gao, and Shangkai Gao. Align  
586 and pool for eeg headset domain adaptation (alpha) to facilitate dry electrode based ssvep-bci.  
*IEEE Transactions on Biomedical Engineering*, 69(2):795–806, 2021a.
- 587  
588 Bingchuan Liu, Xiaogang Chen, Nanlin Shi, Yijun Wang, Shangkai Gao, and Xiaorong Gao. Im-  
589 proving the performance of individually calibrated ssvep-bci by task-discriminant component  
590 analysis. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 29:1998–2007,  
591 2021b.
- 592  
593 Jiawei Liu, Ruimin Wang, Yuankui Yang, Yuan Zong, Yue Leng, Wenming Zheng, and Sheng Ge.  
Convolutional transformer-based cross subject model for ssvep-based bci classification. *IEEE*  
*Journal of Biomedical and Health Informatics*, 2024.

- 594 Fabien Lotte. Signal processing approaches to minimize or suppress calibration time in oscillatory  
595 activity-based brain–computer interfaces. *Proceedings of the IEEE*, 103(6):871–890, 2015.  
596
- 597 Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*  
598 *learning research*, 9(Nov):2579–2605, 2008.
- 599 Masaki Nakanishi, Yijun Wang, Yu Te Wang, Yasue Mitsukura, and Tzyyung Jung. A high-speed  
600 brain speller using steady-state visual evoked potentials. *International Journal of Neural Systems*,  
601 24(6):1450019, 2014.  
602
- 603 Masaki Nakanishi, Yijun Wang, Xiaogang Chen, Yu-Te Wang, Xiaorong Gao, and Tzyy-Ping Jung.  
604 Enhancing detection of ssveps for a high-speed brain speller using task-related component analy-  
605 sis. *IEEE Transactions on Biomedical Engineering*, 65(1):104–112, 2017.  
606
- 607 Jie Pan, Xiaorong Gao, Fang Duan, Zheng Yan, and Shangkai Gao. Enhancing the classification ac-  
608 curacy of steady-state visual evoked potential-based brain–computer interfaces using phase con-  
609 strained canonical correlation analysis. *Journal of neural engineering*, 8(3):036027, 2011.
- 610 Feifei Qi, Yuanqing Li, and Wei Wu. Rstfc: A novel algorithm for spatio-temporal filtering and  
611 classification of single-trial eeg. *IEEE transactions on neural networks and learning systems*, 26  
612 (12):3070–3082, 2015.
- 613 Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh,  
614 and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local-  
615 ization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626,  
616 2017.  
617
- 618 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and  
619 Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In  
620 *International Conference on Learning Representations*, 2017.
- 621 Yonghao Song, Qingqing Zheng, Bingchuan Liu, and Xiaorong Gao. Eeg conformer: Convolu-  
622 tional transformer for eeg decoding and visualization. *IEEE Transactions on Neural Systems and*  
623 *Rehabilitation Engineering*, 31:710–719, 2022.  
624
- 625 Yonghao Song, Bingchuan Liu, Xiang Li, Nanlin Shi, Yijun Wang, and Xiaorong Gao. Decoding  
626 Natural Images from EEG for Object Recognition. In *International Conference on Learning*  
627 *Representations*, 2024.
- 628 Zhijiang Wan, Manyu Li, Shichang Liu, Jiabin Huang, Hai Tan, and Wenfeng Duan. Eegformer: A  
629 transformer–based brain activity classification method using eeg signal. *Frontiers in neuroscience*,  
630 17:1148855, 2023.  
631
- 632 Christopher Wang, Vighnesh Subramaniam, Adam Uri Yaari, Gabriel Kreiman, Boris Katz, Igna-  
633 cio Cases, and Andrei Barbu. Brainbert: Self-supervised representation learning for intracranial  
634 recordings. In *International Conference on Learning Representations*, 2023.
- 635 Yijun Wang, Xiaogang Chen, Xiaorong Gao, and Shangkai Gao. A benchmark dataset for ssvep-  
636 based brain–computer interfaces. *IEEE Transactions on Neural Systems and Rehabilitation En-*  
637 *gineering*, 25(10):1746–1752, 2016.  
638
- 639 Francis R Willett, Donald T Avansino, Leigh R Hochberg, Jaimie M Henderson, and Krishna V  
640 Shenoy. High-performance brain-to-text communication via handwriting. *Nature*, 593(7858):  
641 249–254, 2021.
- 642 Francis R Willett, Erin M Kunz, Chaofei Fan, Donald T Avansino, Guy H Wilson, Eun Young  
643 Choi, Foram Kamdar, Matthew F Glasser, Leigh R Hochberg, Shaul Druckmann, et al. A high-  
644 performance speech neuroprosthesis. *Nature*, 620(7976):1031–1036, 2023.  
645
- 646 Jonathan R Wolpaw, Niels Birbaumer, Dennis J McFarland, Gert Pfurtscheller, and Theresa M  
647 Vaughan. Brain–computer interfaces for communication and control. *Clinical neurophysiology*,  
113(6):767–791, 2002.

- 648 Chi Man Wong, Feng Wan, Boyu Wang, Ze Wang, Wenya Nan, Ka Fai Lao, Peng Un Mak, Mang I  
649 Vai, and Agostinho Rosa. Learning across multi-stimulus enhances target recognition methods in  
650 ssvep-based bcis. *Journal of neural engineering*, 17(1):016026, 2020.
- 651 Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt:  
652 Introducing convolutions to vision transformers. In *Proceedings of the IEEE/CVF international  
653 conference on computer vision*, pp. 22–31, 2021.
- 654 Wei Wu, Xiaorong Gao, Bo Hong, and Shangkai Gao. Classifying single-trial eeg during motor  
655 imagery by iterative spatio-spectral patterns learning (isspl). *IEEE Transactions on Biomedical  
656 Engineering*, 55(6):1733–1743, 2008.
- 657 Xiaolin Xiao, Lichao Xu, Jin Yue, Baizhou Pan, Minpeng Xu, and Dong Ming. Fixed template  
658 network and dynamic template network: novel network designs for decoding steady-state visual  
659 evoked potentials. *Journal of Neural Engineering*, 19(5):056049, 2022.
- 660 Xiaoli Yang, Yurui Li, Jianyu Zhang, Huiyuan Tian, Shijian Li, and Gang Pan. Evomoe: Evolution-  
661 ary mixture-of-experts for ssvep-eeg classification with user-independent training. *IEEE Journal  
662 of Biomedical and Health Informatics*, 2025.
- 663 Xinyi Zhang, Shuang Qiu, Yukun Zhang, Kangning Wang, Yijun Wang, and Huiguang He. Bidi-  
664 rectional siamese correlation analysis method for enhancing the detection of ssveps. *Journal of  
665 Neural Engineering*, 19(4):046027, 2022.
- 666 Xinyi Zhang, Wei Wei, Shuang Qiu, Xujin Li, Yijun Wang, and Huiguang He. Enhancing ssvep-  
667 based bci performance via consensus information transfer among subjects. *IEEE Transactions on  
668 Neural Networks and Learning Systems*, 2024.
- 669 Yu Zhang, Guoxu Zhou, Jing Jin, Minjue Wang, Xingyu Wang, and Andrzej Cichocki. L1-  
670 regularized multiway canonical correlation analysis for ssvep-based bci. *IEEE Transactions on  
671 Neural Systems and Rehabilitation Engineering*, 21(6):887–896, 2013.
- 672 Yu Zhang, Guoxu Zhou, Jing Jin, Xingyu Wang, and Andrzej Cichocki. Frequency recognition  
673 in ssvep-based bci using multiset canonical correlation analysis. *International journal of neural  
674 systems*, 24(04):1450013, 2014.

## 678 A IMPLEMENTATION DETAILS

681 For both ViT stages, the patch temporal length is set to 10 and the hidden dimension to 48, with  
682 a dropout rate of 0.8. For all trial durations except 0.2 s, the generative time length is 0.04 s,  
683 while for trials of 0.2 s, the generative time length equals 0.2 s. In the MoE decoder stage, the  
684 MoE is applied only to the second temporal layer, using 4 experts. The gating network consists  
685 of two MLPs, each with an intermediate dimension of 100, and only the top expert is selected.  
686 For implementation convenience, the input data is reshaped to  $(B \times C, 1, T)$ , where  $B$  represents  
687 the subband dimension. Consequently, all subband-channel layers use a kernel size of  $(1, 1)$ , and  
688 the two temporal convolution layers use a kernel size of 10. The output channels are 200 for all  
689 subband-channel layers and 120 for temporal layers. All dropout layers have a probability of 0.1,  
690 except for the layer before the final flattening and MLP, which uses 0.95. During fine-tuning, the  
691 dropout probability of all intermediate layers is reset to 0.5.

692 For data augmentation, cross-subject temporal stitching is performed with a time chunk of 30. For  
693 channel chunk shuffling, the chunk size is 20 with a swap probability of 0.3. Random temporal  
694 cropping is applied with an activation probability of 0.4, selecting a short segment of 0.02 s to 0.06  
695 s. For each augmentation method, additional data corresponding to 20% of the original dataset size  
696 is generated.

697 The ViT learning rate is set to 0.0001 during general pretraining and 0.00001 during subject-specific  
698 fine-tuning, while the decoder learning rate is fixed at 0.0001. The Adam optimizer is used with a  
699 weight decay of 0.0001, and an L2 regularization penalty of 0.001 is applied to the decoder. The  
700 batch size for both datasets is 32.

701 The number of training epochs for each stage differs between the Benchmark and BETA datasets,  
as summarized in Table 1.

Table 1: Stage-wise training epochs for Benchmark and BETA datasets.

Dataset	ViT Pretrain	ViT Transfer	Decoder Pretrain	Decoder Transfer
Benchmark	300	1000	1500	1000
BETA	300	500	500	700

## B MODEL SIZE COMPARISON

We report both the total number of trainable parameters and the corresponding memory usage in FP32 precision. As summarized in Table 2, all models remain lightweight, with memory requirements well within the range suitable for real time or embedded deployment. For this analysis, we only consider the deep learning-based models. Traditional methods are excluded because they do not maintain persistent trainable parameters and their memory footprint is dominated by temporary data buffers rather than model weights, making a direct comparison of “model size” with neural architectures not meaningful. Note that TRCANet and the DNN share the same decoder architecture, resulting in exactly the same number of parameters. This comparison highlights the balance between model capacity and efficiency across the evaluated architectures.

Table 2: Model size comparison: parameters and approximate memory (FP32).

Method	Parameters	Approx. Size (MB, FP32)
DNN	0.18M	0.7
TRCANet	0.18M	0.7
SSVEPformer	9.26 M	37.1
Ours	1.18M	4.7

## C DECORRELATION DETAILS

For each subject  $s$  and subband  $b$ , we first compute the mean across training trials:

$$\mu^{(s,b)} = \frac{1}{N} \sum_{n \in \text{train trials}} X_{n, :, :, :}^{(s,b)} \in \mathbb{R}^{M \times T \times C}. \quad (2)$$

The aggregated mean activity is reshaped into  $\mu^{(s,b)} \in \mathbb{R}^{C \times (M \cdot T)}$ , and used to compute the channel covariance matrix:

$$\text{Cov}^{(s,b)} = \frac{1}{M \cdot T} \mu^{(s,b)} \left( \mu^{(s,b)} \right)^\top \in \mathbb{R}^{C \times C}. \quad (3)$$

The whitening matrix is defined as

$$W^{(s,b)} = \left( \text{Cov}^{(s,b)} \right)^{-\frac{1}{2}}, \quad (4)$$

and decorrelation is applied to both training and test data as

$$\tilde{X}_{n, :, :, :}^{(s,b)} = W^{(s,b)} X_{n, :, :, :}^{(s,b)} \quad \forall n. \quad (5)$$

By using the trial-averaged activity to construct the covariance, this procedure reduces trial-level variability while preserving target and temporal structure, and ensures that whitening is guided by stable patterns rather than noisy single-trial fluctuations.

## D FURTHER ABLATION STUDY

### D.1 MoE

The results of applying the MoE module at different layers are summarized in Table 3. When MoE was applied to the first subband-channel layer, the accuracy was **65.0%**, and applying it to the second subband-channel layer yielded **65.2%**. In contrast, applying MoE to the first temporal layer

resulted in a lower accuracy of **63.9%**, while placing it on the second temporal layer achieved an accuracy of **64.8%**. These results suggest that the second temporal layer and the subband layer were particularly important for MoE, as they contributed more significantly to improving performance compared to other layers. This highlights the importance of capturing frequency and temporal dynamics at these stages of the model.

Table 3: MoE Ablation Study: Different MoE Configurations. The Subbands listed refer to the Sub-band channel Conv layer in Subband-channel temporal blocks. Tested on Benchmark for 0.2s.

MoE Type	Subband 1	Subband 2	Both Temporal	Temporal 1
Accuracy (%)	65.0	65.2	64.8	63.9

## D.2 DATA AUGMENTATION

The impact of different data augmentation strategies is summarized in Table 4. Removing temporal stitching, channel shuffle, or temporal crop resulted in minor decreases in accuracy of around 1%, while omitting decorrelation caused the largest drop to **63.3%**. These results indicate that all augmentation components contributed to model performance, with decorrelation having the most significant effect. Notably, removing all three data generation methods resulted in a 2% decrease, suggesting that each method provided complementary benefits along different dimensions. **Table 5 further shows the effect of progressively adding augmentation components, with accuracy steadily increasing from using only decorrelation to including stitching, crop, and shuffle.**

Table 4: Data Augmentation Ablation Study: Different Data Augmentation Configurations. Tested on Benchmark for 0.2s.

Augmentation	No Stitching	No Channel Shuffle	No Temp Crop	No Decorrelation
Accuracy (%)	64.5	64.9	64.7	63.3

Table 5: **Data Augmentation Ablation Study: Different Data Augmentation Configurations, showing the effect of progressively adding augmentation components. Tested on Benchmark for 0.2s.**

Augmentation	Decorrelation	Decorrelation + Stitching	Decorrelation + Stitching + Crop	Decorrelation + Stitching + Crop + Shuffle
Accuracy (%)	64.1	64.7	64.9	65.5

## D.3 EFFECT OF ViT GENERATION TIME LENGTH

Table 6 shows the impact of varying the ViT generation time length on classification accuracy. On the Benchmark dataset at 0.2 s, performance improved steadily from **64.5%** (0.04s) to **65.5%** (0.2 s). On the BETA dataset, a similar trend was observed, with accuracy increasing from **52.47%** (0.04 s) to **54.13%** (0.2 s). We note that for other data lengths (0.3 s to 1.0 s), the generated data augmentation achieving the best result was fixed at 0.04 s. To illustrate the effect of longer generation times, we performed the same experiment on the 0.4 s data length. The shortest generation time of 0.04 s achieved the highest accuracy (Benchmark: **83.85%**, BETA: **70.85%**), while increasing the generation time gradually decreased performance across other settings by up to 1.5%.

## E SUPPLEMENTARY TABLES AND FIGURES

Table 6: Effect of ViT generation time length on classification accuracy (%).

Dataset	0.04 s	0.08 s	0.12 s	0.16 s	0.20 s
Benchmark (0.2 s)	64.50	64.91	65.28	65.01	<b>65.50</b>
BETA (0.2s)	52.47	52.64	52.72	53.26	<b>54.13</b>
Benchmark (0.4 s)	<b>83.85</b>	83.14	82.60	82.64	82.71
BETA (0.4s)	<b>70.85</b>	70.50	69.96	69.17	69.50

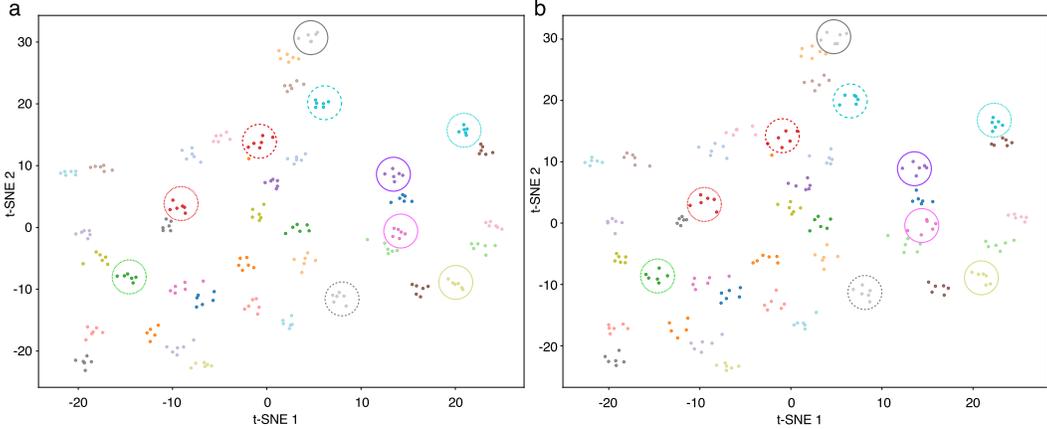


Figure 4: t-SNE visualization of learned features for a representative subject from the Benchmark dataset, using a data length of 0.6 s. Left: full model; Right: ablated model (without ViT, data augmentation, or MoE).

Table 7: Greenhouse–Geisser corrected two-way repeated-measures ANOVA results for the interaction effect between data length and method.

Effect	Benchmark (Accuracy)	Benchmark (ITR)	BETA (Accuracy)	BETA (ITR)
$F$	22.598	22.336	25.256	32.494
$p_{GG}$	$3.45 \times 10^{-7}$	$6.17 \times 10^{-7}$	$7.05 \times 10^{-8}$	$4.86 \times 10^{-10}$

Table 8: Paired t-test p-values comparing VIBE with DNN and TDCA for Benchmark and BETA across data lengths (0.2–1.0 s).

Data length (s)	Benchmark		BETA	
	VIBE vs DNN	VIBE vs TDCA	VIBE vs DNN	VIBE vs TDCA
0.2	$1.044 \times 10^{-12}$	$2.211 \times 10^{-14}$	$4.821 \times 10^{-18}$	$3.305 \times 10^{-25}$
0.3	$1.713 \times 10^{-11}$	$1.947 \times 10^{-7}$	$7.752 \times 10^{-13}$	$1.880 \times 10^{-17}$
0.4	$2.147 \times 10^{-7}$	$5.304 \times 10^{-7}$	$1.144 \times 10^{-7}$	$8.812 \times 10^{-13}$
0.5	$1.269 \times 10^{-6}$	$1.073 \times 10^{-5}$	$2.888 \times 10^{-9}$	$1.010 \times 10^{-14}$
0.6	$2.425 \times 10^{-7}$	$4.578 \times 10^{-5}$	$1.933 \times 10^{-10}$	$1.638 \times 10^{-14}$
0.7	$8.017 \times 10^{-5}$	$5.435 \times 10^{-3}$	$1.100 \times 10^{-8}$	$7.185 \times 10^{-11}$
0.8	$1.284 \times 10^{-3}$	$8.630 \times 10^{-3}$	$1.463 \times 10^{-9}$	$1.445 \times 10^{-11}$
0.9	$1.398 \times 10^{-3}$	$2.856 \times 10^{-2}$	$6.985 \times 10^{-10}$	$7.206 \times 10^{-9}$
1.0	$2.875 \times 10^{-3}$	$3.825 \times 10^{-2}$	$1.275 \times 10^{-6}$	$4.638 \times 10^{-7}$

Table 9: Training and testing times for VIBE on two datasets for time data length 0.4s. Times are in seconds, except for the test stages, which are in milliseconds.

Dataset	ViT			MoE Decoder		
	Train (s)	Finetune (s)	Test (ms)	Train (s)	Finetune (s)	Test (ms)
Benchmark	270.8	19.4	0.7	4180.3	43.2	0.09
BETA	328.9	5.7	0.7	1681.1	11.8	0.09

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

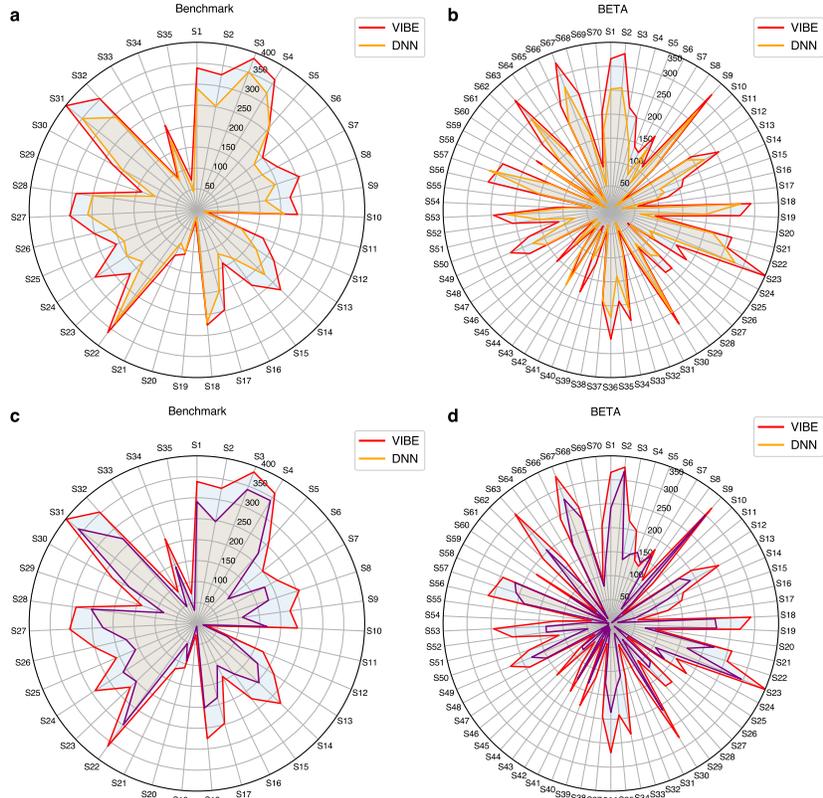


Figure 5: Subject-wise ITR plot. The methods were evaluated at a data length of 0.2s.

Table 10: Ablation General Results: No ViT, No MoE, and No Data Augmentation. Tested on Benchmark for 0.2s.

Model	No ViT	No MoE	No Data Augmentation
Accuracy(%)	61.8	64.2	62.1

## F LLM USAGE

Large language models (LLMs) were used solely for proofreading this manuscript to improve language clarity and readability. They were not involved in generating ideas, designing experiments, implementing methods, or analyzing results. All scientific contributions are entirely the work of the authors.

## G ANALYZING MULTI-STAGE CONTRIBUTIONS WITH MOE

In this analysis, we investigate how different stages of VIBE influence the model’s frequency domain behavior under a controlled experimental setting. We use slightly modified parameters relative to the main experiments: each input consists of a 1s EEG segment, and the Stage 1 ViT generator produces an additional 0.2s extension. The decoder is implemented with eight temporal experts in the MoE layer. Both the Stage 3 (non-finetuned) and Stage 4 (finetuned) decoders are trained using the ViT augmented data, and all visualizations are performed using the same training set to remove cross-split variability. For each subject and trial, we extract the activation from the final temporal Conv2D layer and average across channels as well as subjects to obtain a single temporal sequence per target. This sequence is zero-padded to 5s to achieve sufficient spectral resolution for targets ranging from 8 Hz to 15.8 Hz, after which we apply an FFT. We compute the signal-to-noise ratio (SNR) by taking the peak amplitude at the target frequency and its second harmonic, normalizing each by the mean amplitude of neighboring frequency bins, and averaging the two values. We then perform paired t-tests to assess statistical significance. The analysis compares four conditions: (1) Stage 3 versus Stage 4 decoders, and (2) ViT augmented input (original+generated) versus original only input with zero-padding. The corresponding FFT visualizations and t-test statistics are summarized in the accompanying Figure 6 and Table 11. The FFT visualizations and t-test results consistently show that both the inclusion of ViT-generated data and decoder fine-tuning improve spectral responses. Specifically, the amplitude at the target frequency and its second harmonic is highest for FT with full input, followed by FT with original input, then BASE full, and finally BASE original. All comparisons are statistically significant, highlighting that ViT augmentation and Stage 4 fine-tuning are important for enhancing frequency-specific features in the model.

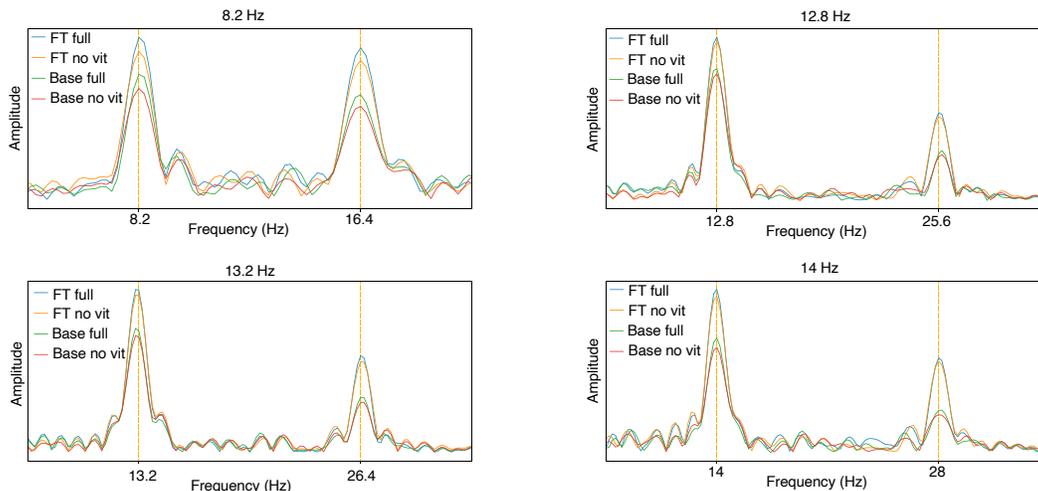


Figure 6: FFT visualizations of activation sequences from the final temporal layer. FT refers to Stage 4 and BASE refers to Stage 3. Full uses the complete input sequence (original data + Stage 2 ViT-generated data), while Original keeps only the original data and zero-pads the rest. Yellow dashed vertical lines indicate the target frequency and its second harmonic. The amplitude at the target frequency and its second harmonic is highest for FT with full input, followed by FT with original input, then BASE full, and finally BASE original. These results demonstrate that both incorporating ViT-generated data and fine-tuning the decoder strengthen the neuro physiological spectral responses.

Table 11: Paired t-test results comparing SNR differences from the final temporal layer activations. FT refers to Stage 4 and BASE refers to Stage 3. Full uses the complete input sequence (original data + Stage 2 ViT-generated data), while Original keeps only the original data and zero-pads the rest.

Comparison	t-stat	p-value	n
FT (full – original)	3.3941	$1.59 \times 10^{-3}$	40
BASE (full – original)	3.0741	$3.84 \times 10^{-3}$	40
FT full – BASE full	9.2720	$2.08 \times 10^{-11}$	40
FT original – BASE original	7.6147	$3.12 \times 10^{-9}$	40

## H VISUALIZATION ANALYSIS FOR THE MOE TEMPORAL LAYER

We perform a qualitative analysis of the MoE temporal layer using the fine-tuned decoder for a randomly selected subject. Grad-CAM (Selvaraju et al., 2017) is applied to the final MoE temporal convolution layer to obtain activation importance over time. In Fig 8, the resulting temporal map is averaged across channels to produce a single 1D sequence. As in the previous analyses, this sequence is zero-padded to 5s to ensure sufficient spectral resolution before applying the FFT. In addition to the frequency-domain visualization, we record the MoE expert selection count Fig 7 for each target frequency to examine how the mixture-of-experts distributes attention across different spectral components. The MoE expert selection patterns reveal that the experts do not collapse into a single dominant expert; instead, different target frequencies elicit distinct expert activation distributions. This diversity indicates that individual experts specialize in different temporal-spectral patterns rather than redundantly modeling the same structure. Complementing this, the Grad-CAM analysis shows clear spectral peaks at the target frequency and its harmonic after FFT, demonstrating that the MoE temporal layer effectively pools and amplifies frequency-specific structure. Together, these results confirm that the MoE architecture meaningfully decomposes the temporal dynamics and contributes specialized processing across targets.

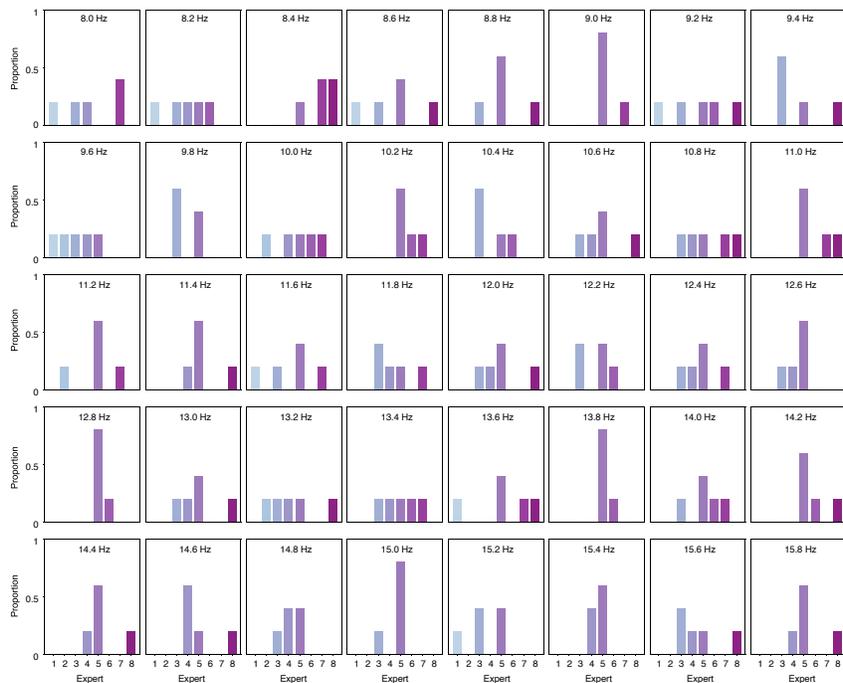


Figure 7: MoE expert selection counts for all 40 target labels in a randomly selected subject. The x-axis represents the different experts (1–8), and the y-axis shows the proportion of times each expert is selected. MoE expert activations vary systematically across target frequencies, reflecting frequency-specific patterns

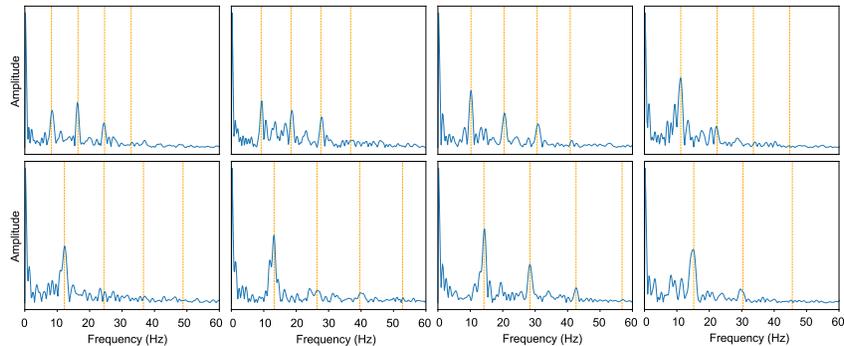


Figure 8: Grad-CAM visualization of the final MoE temporal layer, averaged across channels, for a randomly selected subject. FFTs are shown for selected target frequencies: 8, 9, 10, 11 Hz (top row) and 12, 13, 14, 15 Hz (bottom row) to illustrate temporal importance patterns. Yellow dashed vertical lines indicate the target frequency and its harmonics. Across all target frequencies, the Grad-CAM results consistently show peaks at the target frequency and its second harmonic, demonstrating that the MoE layer captures frequency-specific temporal activations.

## I SPECTRAL VALIDATION OF DATA AUGMENTATION

We provide spectral analysis for the two data augmentation methods used in our framework, illustrating their preservation on SSVEP frequency structure. Because stitch-augmented data and channel-chunk-shuffle augmentation are less intuitive than our other augmentation strategies, we include explicit spectral analyses to clarify how they operate.

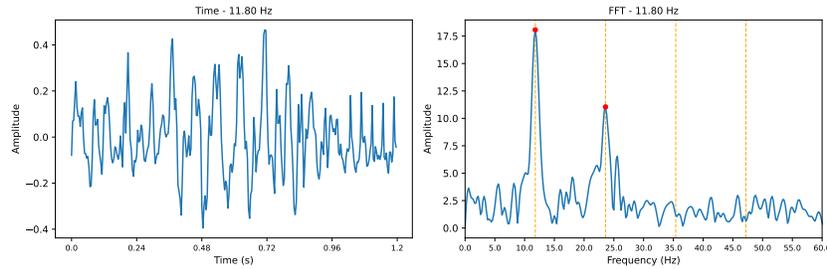
### I.1 SPECTRAL VALIDATION OF STITCH-AUGMENTED DATA

To assess whether the Cross-Subject Temporal Stitching augmentation preserves useful SSVEP structure, we visualize the temporal and spectral profiles of the augmented samples. We use 1 s temporal length original data as input from Benchmark. After Stage 2 ViT generation (add 0.2s generate data) and Stitch recomposition, we select a single subband, average across channels and trials dimensions, zero-pad the resulting sequence to 5s, and apply an FFT. The corresponding temporal traces and spectra consistently exhibit clear peaks at the target frequency and its harmonic, indicating that the stitched signals retain the essential SSVEP signatures. Although Stitch mixes temporal chunks across channels and subjects, this variability diversifies the training distribution for the decoder while maintaining the frequency-specific structure required for accurate decoding. The accompanying plot (Fig. 9) contains two subpanels: the left panel shows the stitched time-series signal, and the right panel shows the FFT magnitude spectrum. In the spectral panel, clear peaks appear at the target frequency and its harmonic ( $2 \times \text{freq}$ ), confirming that the stitched samples preserve the characteristic SSVEP structure.

### I.2 SPECTRAL VALIDATION OF CHANNEL CHUNK SHUFFLE AUGMENTED DATA

To verify that the Channel Chunk Shuffle augmentation preserves essential SSVEP structure, we visualized the spectral profiles of augmented samples. We use 1 s temporal length original data as input from Benchmark. Following Stage 2 ViT generation (add 0.5s data) and shuffle recomposition, we select a single subband and channel, average across trials, zero-pad the resulting sequence to 5 s, and compute the FFT. By shuffling within defined temporal chunks, we expose the network to different channel arrangements, enhancing generalization while preserving physiologically meaningful spectral features. The accompanying plot (Fig. 10) shows that prominent peaks at the target frequency and its harmonic, indicating that the shuffled signals maintain the frequency-specific structure required for accurate decoding.

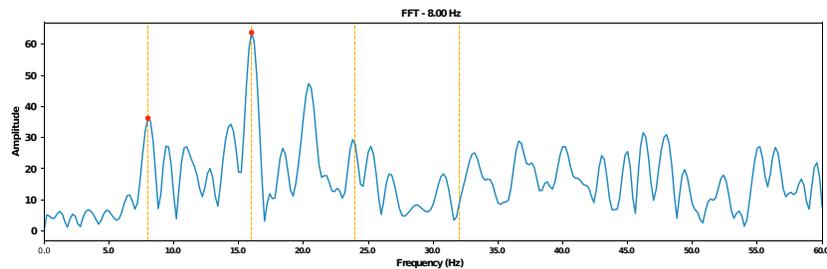
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094



1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113

Figure 9: Time series (left) and corresponding FFT magnitude spectrum (right) of a Cross-Subject Temporal Stitching sample for a representative target frequency. We use 1 s temporal length original data as input from Benchmark and have ViT generate additional 0.2 s. In the spectral panel, the red dots on the curves indicate peaks at the target frequency and its second harmonic ( $2 \times \text{freq}$ ), where target frequency and harmonics are highlighted by yellow dashed vertical lines. These clear spectral peaks confirm that the stitched samples preserve the characteristic SSVEP spectral structure.

1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122



1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

Figure 10: FFT magnitude spectrum (right) of a Channel Chunk Shuffle sample for a representative target frequency. We use 1 s temporal length original data as input from Benchmark and have ViT generate additional 0.5 s. The red dots on the curves indicate peaks at the target frequency and its second harmonic ( $2 \times \text{freq}$ ), where target frequency and harmonics are highlighted by yellow dashed vertical lines. These clear spectral peaks confirm that the re-composited sample preserve the characteristic SSVEP spectral structure.

## J COMPREHENSIVE RESULT TABLES AND ADDITIONAL BASELINES

Tables 12,13,14,15 summarize the performance of our model compared to baseline methods under different time lengths (0.2 to 1.0 s) on both the Benchmark and beta datasets. Tables 1 and 3 report classification accuracy (%), while Tables 2 and 4 present ITR. Across all conditions, our model consistently outperforms all baseline methods, demonstrating superior accuracy and efficiency. Additionally, we include comparisons with three recently proposed methods, Dis-ComNet (Li et al., 2025b), SESCNN (Jin et al., 2024), and ConsenNet (Zhang et al., 2024), which were not included in the result Fig 2,3. Our model achieves the highest accuracy and competitive iteration performance across all settings, highlighting its effectiveness and robustness.

**Table 12: Classification accuracy (%) of various methods at different time lengths (0.2 - 1.0 s) on the Benchmark dataset. Values are reported as mean  $\pm$  standard deviation.**

Method	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
DNN	58.3 $\pm$ 3.2	72.3 $\pm$ 3.0	80.8 $\pm$ 2.7	85.1 $\pm$ 2.5	88.6 $\pm$ 2.2	91.6 $\pm$ 1.8	93.7 $\pm$ 1.6	95.2 $\pm$ 1.3	95.8 $\pm$ 1.1
TDCA	54.1 $\pm$ 3.4	70.5 $\pm$ 3.4	79.1 $\pm$ 3.1	84.6 $\pm$ 2.7	89.1 $\pm$ 2.3	92.2 $\pm$ 2.0	94.3 $\pm$ 1.6	95.7 $\pm$ 1.2	96.6 $\pm$ 1.0
TRCA	49.0 $\pm$ 3.6	64.0 $\pm$ 3.6	72.5 $\pm$ 3.5	79.3 $\pm$ 3.2	84.1 $\pm$ 3.0	87.6 $\pm$ 2.8	90.2 $\pm$ 2.5	92.5 $\pm$ 2.1	93.6 $\pm$ 1.9
eCCA	28.3 $\pm$ 2.7	49.4 $\pm$ 3.7	62.7 $\pm$ 3.9	72.7 $\pm$ 3.6	79.9 $\pm$ 3.3	85.4 $\pm$ 3.0	89.4 $\pm$ 2.6	92.1 $\pm$ 2.1	93.5 $\pm$ 1.8
msTRCA	53.3 $\pm$ 3.4	66.7 $\pm$ 3.6	74.3 $\pm$ 3.4	79.8 $\pm$ 3.2	84.4 $\pm$ 3.0	88.4 $\pm$ 2.7	90.7 $\pm$ 2.4	92.7 $\pm$ 2.0	93.9 $\pm$ 1.8
SSVEPformer	46.1 $\pm$ 3.1	62.0 $\pm$ 3.5	69.7 $\pm$ 3.3	76.3 $\pm$ 3.1	81.5 $\pm$ 3.0	86.0 $\pm$ 2.7	89.7 $\pm$ 2.4	91.9 $\pm$ 2.0	93.0 $\pm$ 1.8
TRCANet	58.4 $\pm$ 3.1	72.3 $\pm$ 3.0	80.8 $\pm$ 2.7	85.2 $\pm$ 2.5	88.8 $\pm$ 2.2	91.5 $\pm$ 1.8	93.7 $\pm$ 1.6	95.2 $\pm$ 1.3	95.5 $\pm$ 1.1
Dis-ComNet	54.2 $\pm$ 2.9	68.4 $\pm$ 3.6	76.0 $\pm$ 3.1	82.6 $\pm$ 3.0	86.3 $\pm$ 2.7	89.1 $\pm$ 2.8	92.2 $\pm$ 2.2	93.1 $\pm$ 1.9	95.6 $\pm$ 1.8
SESCNN	60.1 $\pm$ 3.2	73.5 $\pm$ 2.9	81.7 $\pm$ 3.3	85.6 $\pm$ 2.9	88.7 $\pm$ 2.6	91.2 $\pm$ 2.6	93.6 $\pm$ 2.7	94.9 $\pm$ 2.5	95.6 $\pm$ 2.7
Ours	65.5 $\pm$ 3.2	77.0 $\pm$ 2.9	83.8 $\pm$ 2.5	88.3 $\pm$ 2.2	91.7 $\pm$ 2.1	94.2 $\pm$ 1.9	95.6 $\pm$ 1.6	96.8 $\pm$ 1.1	97.4 $\pm$ 0.7

**Table 13: ITR (bits/min) of various methods at different time lengths (0.2 - 1.0 s) on the Benchmark dataset. Values are reported as mean  $\pm$  standard deviation.**

Method	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
DNN	194.6 $\pm$ 14.8	235.4 $\pm$ 13.5	248.8 $\pm$ 11.8	243.3 $\pm$ 10.2	236.0 $\pm$ 8.5	228.1 $\pm$ 6.9	218.9 $\pm$ 5.8	208.7 $\pm$ 4.4	196.5 $\pm$ 3.8
TDCA	173.9 $\pm$ 15.5	227.5 $\pm$ 15.3	240.6 $\pm$ 13.0	240.5 $\pm$ 11.0	237.3 $\pm$ 8.8	230.3 $\pm$ 7.4	220.8 $\pm$ 5.8	210.0 $\pm$ 4.4	199.1 $\pm$ 3.4
msTRCA	151.7 $\pm$ 16.1	197.3 $\pm$ 15.6	212.0 $\pm$ 14.2	218.5 $\pm$ 12.3	217.7 $\pm$ 10.7	213.4 $\pm$ 9.5	206.3 $\pm$ 8.0	199.2 $\pm$ 6.6	189.5 $\pm$ 5.7
eCCA	65.2 $\pm$ 10.0	134.8 $\pm$ 14.9	171.4 $\pm$ 15.0	192.2 $\pm$ 13.4	201.6 $\pm$ 11.7	205.2 $\pm$ 10.1	203.3 $\pm$ 8.3	197.8 $\pm$ 6.7	188.9 $\pm$ 5.6
msTRCA	170.5 $\pm$ 16.0	210.1 $\pm$ 15.8	220.1 $\pm$ 14.0	220.6 $\pm$ 12.5	219.5 $\pm$ 10.9	216.7 $\pm$ 9.4	208.2 $\pm$ 7.9	199.9 $\pm$ 6.5	190.6 $\pm$ 5.6
SSVEPformer	137.4 $\pm$ 13.9	188.8 $\pm$ 15.3	199.3 $\pm$ 13.5	206.2 $\pm$ 11.9	208.0 $\pm$ 10.7	207.0 $\pm$ 9.2	204.4 $\pm$ 7.8	197.5 $\pm$ 6.4	187.9 $\pm$ 5.6
TRCANet	196.3 $\pm$ 14.8	235.5 $\pm$ 13.5	248.6 $\pm$ 11.8	243.2 $\pm$ 10.2	236.2 $\pm$ 8.5	226.9 $\pm$ 6.8	217.8 $\pm$ 5.8	207.6 $\pm$ 4.4	194.8 $\pm$ 3.7
Dis-ComNet	185.2 $\pm$ 15.1	225.6 $\pm$ 13.8	236.0 $\pm$ 13.5	234.3 $\pm$ 10.3	230.2 $\pm$ 9.8	220.9 $\pm$ 8.3	210.3 $\pm$ 7.4	205.7 $\pm$ 6.9	198.1 $\pm$ 5.4
SESCNN	204.5 $\pm$ 14.5	237.8 $\pm$ 15.9	253.0 $\pm$ 12.8	245.4 $\pm$ 11.4	236.6 $\pm$ 10.7	229.3 $\pm$ 9.5	212.2 $\pm$ 7.1	204.1 $\pm$ 7.4	189.9 $\pm$ 6.1
Ours	232.8 $\pm$ 15.7	259.4 $\pm$ 14.1	263.8 $\pm$ 11.7	257.7 $\pm$ 9.4	249.6 $\pm$ 8.9	238.8 $\pm$ 7.2	226.1 $\pm$ 6.2	214.5 $\pm$ 5.1	202.5 $\pm$ 4.2

**Table 14: Classification accuracy (%) of various methods at different time lengths (0.2 - 1.0 s) on the BETA dataset. Values are reported as mean  $\pm$  standard deviation. ConsenNet was evaluated under a slightly different setting: all but one subject were used as their training set. In their fine-tuning stage, they used the first three calibration blocks from the new subject as training and the remaining blocks as testing. Thus, for each test trial in the BETA dataset, all other trials from that subject, as well as all data from other subjects, were seen during training. They used more data for training compared to our setup, incorporating all other trials from the test subject as well as data from all remaining subjects.**

Method	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
DNN	46.2 $\pm$ 2.1	59.9 $\pm$ 2.1	68.0 $\pm$ 2.1	71.8 $\pm$ 2.0	75.5 $\pm$ 2.0	78.7 $\pm$ 1.9	80.9 $\pm$ 1.7	82.3 $\pm$ 1.7	83.7 $\pm$ 1.6
TDCA	41.4 $\pm$ 2.4	55.2 $\pm$ 2.8	63.3 $\pm$ 2.7	68.0 $\pm$ 2.6	72.4 $\pm$ 2.4	76.2 $\pm$ 2.3	78.8 $\pm$ 2.1	80.9 $\pm$ 2.0	82.9 $\pm$ 1.9
msTRCA	36.1 $\pm$ 2.2	47.7 $\pm$ 2.7	54.9 $\pm$ 2.9	59.1 $\pm$ 2.9	63.4 $\pm$ 2.8	66.7 $\pm$ 2.9	69.4 $\pm$ 2.8	72.2 $\pm$ 2.8	74.2 $\pm$ 2.7
eCCA	27.8 $\pm$ 1.8	43.2 $\pm$ 2.5	53.0 $\pm$ 2.7	60.2 $\pm$ 2.7	66.2 $\pm$ 2.6	71.0 $\pm$ 2.6	74.2 $\pm$ 2.5	77.3 $\pm$ 2.4	79.8 $\pm$ 2.2
msTRCA	40.0 $\pm$ 2.3	50.7 $\pm$ 2.7	57.8 $\pm$ 2.7	61.5 $\pm$ 2.7	66.0 $\pm$ 2.7	69.2 $\pm$ 2.7	72.0 $\pm$ 2.6	74.7 $\pm$ 2.5	76.5 $\pm$ 2.5
SSVEPformer	37.8 $\pm$ 1.9	50.8 $\pm$ 2.3	58.9 $\pm$ 2.4	64.7 $\pm$ 2.4	68.9 $\pm$ 2.4	73.3 $\pm$ 2.3	76.5 $\pm$ 2.3	79.4 $\pm$ 2.2	81.5 $\pm$ 2.1
TRCANet	46.1 $\pm$ 2.1	60.0 $\pm$ 2.1	67.9 $\pm$ 2.1	71.7 $\pm$ 2.0	75.6 $\pm$ 2.0	78.9 $\pm$ 1.9	80.9 $\pm$ 1.7	82.2 $\pm$ 1.7	83.8 $\pm$ 1.7
SESCNN	46.9 $\pm$ 3.2	58.9 $\pm$ 3.8	68.1 $\pm$ 3.5	71.2 $\pm$ 3.9	75.9 $\pm$ 4.1	79.1 $\pm$ 3.4	81.1 $\pm$ 4.3	82.8 $\pm$ 3.9	84.1 $\pm$ 4.7
ConsenNet	-	-	67.9 $\pm$ 2.1	-	77.7 $\pm$ 2.2	-	83.6 $\pm$ 1.9	-	-
Ours	54.1 $\pm$ 2.3	64.6 $\pm$ 2.0	70.8 $\pm$ 1.9	75.7 $\pm$ 1.9	79.4 $\pm$ 1.8	82.2 $\pm$ 1.8	84.2 $\pm$ 1.7	85.7 $\pm$ 1.6	86.3 $\pm$ 1.4

## K COMPARISON OF GENERATORS FOR TEMPORAL SEQUENCE EXTENSION

To rigorously evaluate the design choice of using a ViT for temporal sequence generation, we conducted a comparative study against an alternative transformer-based architecture, the Convolutional

Table 15: ITR (bits/min) of various methods at different time lengths (0.2 - 1.0 s) on the BETA dataset. Values are reported as mean  $\pm$  standard deviation. ConsenNet was evaluated under a slightly different setting: all but one subject were used as their training set. In their fine-tuning stage, they used the first three calibration blocks from the new subject as training and the remaining blocks as testing. Thus, for each test trial in the BETA dataset, all other trials from that subject, as well as all data from other subjects, were seen during training. They used more data for training compared to our setup, incorporating all other trials from the test subject as well as data from all remaining subjects.

Method	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
DNN	137.2 $\pm$ 9.1	176.4 $\pm$ 9.0	190.1 $\pm$ 8.5	185.6 $\pm$ 7.6	183.0 $\pm$ 7.0	178.9 $\pm$ 6.3	172.0 $\pm$ 5.5	164.2 $\pm$ 5.0	157.4 $\pm$ 4.6
TDCA	118.3 $\pm$ 9.7	159.9 $\pm$ 11.1	173.0 $\pm$ 10.1	172.8 $\pm$ 9.1	172.4 $\pm$ 8.1	170.6 $\pm$ 7.3	165.5 $\pm$ 6.5	160.2 $\pm$ 5.9	154.8 $\pm$ 5.2
msTRCA	95.8 $\pm$ 8.5	129.4 $\pm$ 10.2	141.8 $\pm$ 10.1	142.1 $\pm$ 9.3	143.2 $\pm$ 8.7	141.8 $\pm$ 8.3	138.8 $\pm$ 7.7	136.9 $\pm$ 7.2	132.8 $\pm$ 6.6
eCCA	63.1 $\pm$ 6.4	110.2 $\pm$ 9.1	133.6 $\pm$ 9.5	144.5 $\pm$ 9.0	151.2 $\pm$ 8.3	154.3 $\pm$ 7.8	151.9 $\pm$ 7.0	150.0 $\pm$ 6.5	146.8 $\pm$ 5.8
msTRCA	111.7 $\pm$ 9.0	140.5 $\pm$ 10.2	151.6 $\pm$ 9.7	149.7 $\pm$ 8.9	151.2 $\pm$ 8.5	148.8 $\pm$ 7.9	145.8 $\pm$ 7.3	143.0 $\pm$ 6.8	138.1 $\pm$ 6.2
SSVEPformer	101.1 $\pm$ 7.4	139.7 $\pm$ 9.2	154.9 $\pm$ 8.9	160.5 $\pm$ 8.4	160.6 $\pm$ 7.9	161.8 $\pm$ 7.3	159.6 $\pm$ 6.6	157.3 $\pm$ 6.1	153.0 $\pm$ 5.6
TRC	137.6 $\pm$ 9.0	176.7 $\pm$ 9.0	189.9 $\pm$ 8.5	185.4 $\pm$ 7.6	183.2 $\pm$ 7.0	179.9 $\pm$ 6.4	171.9 $\pm$ 5.5	163.6 $\pm$ 5.0	157.7 $\pm$ 4.6
SESCNN	141.2 $\pm$ 9.5	178.5 $\pm$ 10.3	189.3 $\pm$ 9.2	186.4 $\pm$ 8.4	182.8 $\pm$ 7.2	179.1 $\pm$ 6.5	171.1 $\pm$ 5.9	165.8 $\pm$ 4.3	159.1 $\pm$ 4.8
ConsenNet	-	-	188.7 $\pm$ 9.0	-	191.4 $\pm$ 7.7	-	181.8 $\pm$ 6.2	-	-
Ours	173.4 $\pm$ 10.4	199.1 $\pm$ 10.1	202.7 $\pm$ 8.9	201.9 $\pm$ 7.3	197.4 $\pm$ 7.2	191.2 $\pm$ 6.8	182.9 $\pm$ 5.9	174.7 $\pm$ 4.5	164.0 $\pm$ 4.1

Vision Transformer (CvT) (Wu et al., 2021). Both models were tasked with generating extended temporal sequences from short EEG segments of 0.2 s, while keeping all modules in multistages identical. This controlled setup isolates the effect of the generator architecture on the overall decoding performance. Quantitative results, summarized in Table 16, demonstrate that ViT consistently achieves higher mean accuracy in two datasets. These findings substantiate the selection of ViT as the temporal generator in VIBE, confirming its advantage for extending short EEG segments while preserving fine-grained, frequency-specific information critical for high-precision SSVEP decoding.

Table 16: Comparison of ViT and CvT as temporal sequence generators. All other components in the VIBE architecture remain fixed. The results show that ViT consistently outperforms CvT in mean accuracy, supporting its selection as the temporal generator.

Model	Benchmark 0.2s	Benchmark 0.4s	BETA 0.2s	BETA 0.4s
ViT	65.5	83.8	54.1	70.8
CvT	62.1	82.5	51.0	69.1

## L ELECTRODE CONFIGURATION AND SENSITIVITY

The choice of EEG electrodes can significantly impact both the performance and practicality of BCI systems. Different electrodes capture overlapping but distinct spatial patterns of brain activity, and their number and placement can influence classification accuracy. To evaluate this, we tested our model using several commonly adopted electrode sets (Liu et al., 2021b):

- **Central occipital montage (Nch = 3):** Oz, O1, O2
- **Classical occipital montage (Nch = 9):** Pz, POz, PO3/4, PO5/6, Oz, O1/2
- **Occipital montage (Nch = 21):** Pz, P1/2, P3/4, P5/6, P7/8, POz, PO3/4, PO5/6, PO7/8, Oz, O1/2, CB1/2
- **Parietal-occipital montage (Nch = 30):** CPz, CP1/2, CP3/4, CP5/6, TP7/8, Pz, P1/2, P3/4, P5/6, P7/8, POz, PO3/4, PO5/6, PO7/8, Oz, O1/2, CB1/2
- **Full montage (Nch = 64):** all channels

Our results Tab 17 show that classification accuracy improves as the number of EEG channels increases from 3 to 21, reaching a peak of 73.0% with 21 channels. However, using all 64 channels does not further improve performance and may reduce system practicality. Considering the trade-off between accuracy and usability, we select a 9-channel configuration (classical occipital montage)

for the experiments in this study. In Table 18, we report the maximum achievable ITR for the five electrode channels selection configurations. Compared with TDCA, our method almost consistently attains higher ITR across all channel combinations.

Table 17: Datalength of 0.2s over two datasets, classification accuracy (%) and ITR for different numbers of EEG channels.

Nch	Benchmark Acc (%)	Benchmark ITR (bpm)	BETA Acc (%)	BETA ITR (%)
3	40.7	117.2	31.8	80.3
9	65.5	232.8	54.1	173.4
21	73.0	272.3	60.0	202.8
30	71.5	263.9	58.0	192.0
64	67.4	240.3	49.9	153.3

Table 18: Maximum ITR across all time lengths for different numbers of EEG channels on the Benchmark and BETA datasets. For each channel configuration, we report the maximum achievable ITR over all evaluated time windows. Results compare our VIBE model with the TDCA baseline.

Nch	Benchmark		BETA	
	VIBE	TDCA	VIBE	TDCA
3	142.5	153.6	119.3	101.3
9	263.8	244.5	202.7	173.6
21	290.7	281.0	216.2	198.0
30	289.0	281.7	213.0	195.6
64	267.3	275.3	189.5	181.9