# Deep Learning Architectures for Code-Modulated Visual Evoked Potentials Detection *

**Kiran Nair**
Department of Computer Science
California State University, Fresno
Fresno, CA, USA
kiranpnair8@mail.fresnostate.edu

**Hubert Cecotti**[†]
Department of Computer Science
California State University, Fresno
Fresno, CA, USA
hcecotti@csufresno.edu

## Abstract

Non-invasive Brain–Computer Interfaces (BCIs) based on Code-Modulated Visual Evoked Potentials (C-VEPs) require highly robust decoding methods to address temporal variability and session-dependent noise in EEG signals. This study proposes and evaluates several deep learning architectures, including convolutional neural networks (CNNs) for 63-bit m-sequence reconstruction and classification, and Siamese networks for similarity-based decoding, alongside canonical correlation analysis (CCA) baselines. EEG data were recorded from 13 healthy adults under single-target flicker stimulation. The proposed deep models significantly outperformed traditional approaches, with distance-based decoding using Earth Mover's Distance (EMD) and constrained EMD showing greater robustness to latency variations than Euclidean and Mahalanobis metrics. Temporal data augmentation with small shifts further improved generalization across sessions. Among all models, the multi-class Siamese network achieved the best overall performance with an average accuracy of 96.89%, demonstrating the potential of data-driven deep architectures for reliable, single-trial C-VEP decoding in adaptive non-invasive BCI systems.

*Keywords* Brain–Computer Interfaces · Deep Learning · Siamese Networks

## 1 Introduction

The rapid advancement of intelligent systems and neurotechnology has transformed how humans interact with machines, driving innovation across domains from healthcare to human augmentation. Among these emerging technologies, Brain-Computer Interfaces (BCIs) have established themselves as a revolutionary pathway enabling direct communication between the brain and external devices without requiring muscular output [1]. BCIs are increasingly recognized as part of the broader evolution of intelligent systems, bridging neuroscience, signal processing, and machine learning. Initially designed to assist motor-impaired users following a spinal cord injury, or those who are unable to speak, such as those suffering from locked-in syndrome, BCIs have since expanded into diverse areas, including neurorehabilitation [2, 3], cognitive enhancement [4], mental state monitoring, and entertainment applications [5]. Their ability to translate neural signals into actionable commands makes BCIs an integral component of the intelligent system ecosystem with significant potential for both clinical and non-clinical applications.

Despite notable progress, BCIs have yet to achieve the robustness, reliability, and user-friendliness required for widespread adoption. The main technical barriers include inconsistent performance across users and sessions, low signal-to-noise ratio (SNR), and sensitivity to artifacts such as eye blinks, muscle activity, and environmental interference [6, 7]. Electroencephalography (EEG) remains the most widely used modality for non-invasive BCIs because of its affordability, safety, and portability. However, EEG signals are inherently noisy and spatially blurred, requiring effective preprocessing

---

and robust decoding algorithms to achieve high accuracy. These challenges underscore the need for intelligent, data-driven approaches that can adapt to user variability and non-stationary signal conditions.

Among non-invasive paradigms, the Code-Modulated Visual Evoked Potential (c-VEP)[8, 9, 10, 11, 12] has emerged as one of the most promising approaches for high-speed and reliable communication. In c-VEP systems, each target is encoded by a unique pseudorandom binary sequence, resulting in distinct time-locked EEG responses when the target is attended. Thanks to the excellent correlation properties of these codes, c-VEPs enable high classification accuracy and information transfer rates [13]. Recent advances have demonstrated impressive communication speeds exceeding 200 bits/min with minimal calibration [14], establishing c-VEPs as strong candidates for real-world BCI deployment. Nevertheless, challenges persist, including subject-specific variability, temporal misalignment in neural responses, and the need for improved preprocessing to enhance spatial selectivity [15]. These issues motivate the development of advanced decoding methods and spatial enhancement strategies to ensure generalizable, session-independent performance.

This study proposes a systematic evaluation of both traditional and deep learning-based decoding strategies for c-VEP BCIs. Specifically, we compare classical feature extraction approaches such as Canonical Correlation Analysis (CCA) and correlation-based methods with Bayesian Linear Discriminant Analysis (BLDA) classifiers against modern neural architectures, including Convolutional Neural Networks (CNNs) and Siamese networks. Surface Laplacian filtering is integrated to enhance spatial resolution and suppress noise, while CNN outputs are evaluated using distance-based similarity metrics, including Euclidean, Mahalanobis, and Earth Mover's Distance (EMD). Together, these components aim to provide a unified framework for robust, high-speed, and calibration-efficient c-VEP decoding.

The main contributions of this paper are summarized as follows:

1. A comparative analysis of traditional correlation-based methods and deep learning architectures (CNN and Siamese networks) for c-VEP decoding.
2. The development of a CNN framework that reconstructs 63-bit code patterns from EEG and classifies targets using distance-based metrics (Euclidean, Mahalanobis, and EMD).
3. The analysis and comparison of single multi-class Siamese networks and multiple-classifier binary Siamese networks for evaluating cross-session generalization and temporal shift invariance in c-VEP decoding.
4. An analysis of data augmentation and decision combination with local temporal shift in c-VEP decoding.

The remainder of this paper is organized as follows. Section 2 reviews prior work related to deep learning in BCIs, c-VEP-based paradigms, spatial filtering, and similarity learning. Section 3 describes the experimental setup for obtaining the signals and all the classifiers. Section 4 presents the performance of the different classifiers. Finally, the results are discussed in Section 5, while Section 6 summarizes the findings.

## 2   Related Works

Deep learning has profoundly influenced EEG-based BCI research by enabling data-driven feature extraction from raw neural signals. [16] introduced EEGNet, a compact CNN architecture that generalizes across paradigms such as P300, SSVEP, and motor imagery, setting a benchmark for lightweight yet robust models. [17] demonstrated that deep convolutional networks could learn interpretable spatiotemporal filters directly from EEG, outperforming traditional CSP-based approaches. More recent studies have refined these architectures with multi-branch 1D CNNs for frequency–time feature fusion [18], as well as transformer-based attention mechanisms for dynamic feature representation [19]. Some studies demonstrated that deep learning methods, such as convolutional neural networks (CNNs), could effectively model electroencephalography (EEG) signals for P300-based detection, marking a milestone in end-to-end neural decoding [20, 21, 22, 23]. Subsequent work emphasized the benefits of current-source density and surface Laplacian filtering for motor-imagery BCIs, demonstrating measurable gains in spatial resolution and classification performance [24, 25]. Such models not only improve accuracy but also reduce reliance on handcrafted features and extensive calibration, marking a paradigm shift toward generalizable, explainable neural decoding.

Convolutional Neural Networks (CNNs) are a class of artificial neural networks widely employed in deep learning for their ability to learn hierarchical feature representations directly from data [26]. Inspired by the organization of the human primary visual cortex [27], CNNs consist of stacked convolutional, pooling, and fully connected layers that progressively extract increasingly abstract features from structured inputs such as images or multi-dimensional signals. Their effectiveness has been demonstrated across diverse domains, including computer vision [28, 29], robotics [30], chemistry [31], and astronomy [32].

Within the family of visual evoked potentials, c-VEP-based BCIs have become a focal point for high-speed, reliable communication. Early concepts by Sutter [33] introduced pseudorandom code sequences to evoke distinct neural

responses, later re-established by [34] using non-invasive EEG. [35] enhanced c-VEP decoding through spatiotemporal beamforming, achieving over 170 bits/min, while [14] demonstrated cross-subject transfer learning that achieved 250 bits/min with under one minute of calibration. [15] provided a comprehensive review of current trends in code modulation, template correlation, and machine-learning-based decoding. Further studies have optimized code design [36], explored multi-target spellers [37], and proposed novel coding strategies such as chaotic or burst stimuli to enhance user comfort and signal distinctiveness [19]. [38] reported near-universal usability with 97.8% accuracy across 86 subjects, mitigating the long-standing issue of BCI illiteracy.

In addition to code-level advances, signal preprocessing plays a critical role in improving EEG reliability [39, 40]. Surface Laplacian (SL) and current source density (CSD) transformations have been shown to enhance spatial resolution by isolating local cortical activity while suppressing distant sources. Carvalhaes and de Barros [41] formalized the theoretical foundations of SL methods, and [42] demonstrated their practical benefits for VEP-based BCIs. [24] indicated that SL filtering improved classification accuracy by 3–5%, particularly under low SNR or sparse electrode conditions. These findings establish Laplacian filtering as an essential preprocessing step for modern EEG pipelines.

Recent developments [43, 44, 45, 46, 47] have extended deep learning in BCIs beyond conventional classification into similarity-based learning. Siamese and twin-network architectures enable metric learning for EEG, allowing models to capture invariant relationships across sessions or users. [48] introduced a Siamese CNN for motor imagery BCIs using one-vs-one scaling, while [49] employed Siamese networks for cross-subject EEG data augmentation. Multiscale Siamese CNNs [50] and recent attention-based twin networks [51] have further advanced the extraction of discriminative EEG embeddings. Such metric-learning approaches are particularly relevant for c-VEP decoding[52, 53, 54, 55, 56], where robustness to latency shifts and inter-session variability is critical.

The literature reveals a convergence of innovations in c-VEP paradigms, spatial pre-processing, and deep metric learning. These works provide the scientific foundation and motivation for developing a unified decoding framework that combines spatial enhancement, correlation-based analysis, and deep neural architectures to achieve reliable, session-independent c-VEP classification.

# 3 Methods

## 3.1 Experimental Protocol

### 3.1.1 Subjects

The experimental protocol was designed to evaluate the performance of the proposed Code-Modulated Visual Evoked Potential (C-VEP)-based Brain–Computer Interface (BCI). The study was approved by the Committee for the Protection of Human Subjects (CPHS) - Institutional Review Board (IRB) at California State University, Fresno, and conducted in accordance with the principles of the Helsinki Declaration of 1975, as revised in 2000.

Thirteen healthy participants (10 males and 3 females; mean age = $23.54 \pm 6.74$ years) voluntarily took part in the experiment. All subjects reported normal or corrected-to-normal vision and no prior history of neurological disorders. Each participant provided written informed consent before the experiment began. During data collection, participants were seated comfortably in a quiet, well-lit room, approximately 60 cm from the computer screen displaying the visual stimuli(see Fig. 1).

### 3.1.2 Signal Acquisition

EEG signals were recorded using a BioSemi ActiveTwo amplifier system with a sampling rate of 512 Hz. A total of $N_s = 8$ active electrodes were positioned over the occipital and parietal regions according to the international 10–20 system[57, 58] at locations: O1, O2, Oz, Pz, P3, P4, PO7, PO8, as shown in Fig. 4. These electrode sites were selected for their known sensitivity to visual evoked potentials. The ground and reference electrodes were connected to the BioSemi common mode sense (CMS) and driven right leg (DRL) electrodes, respectively. A total of five sessions were recorded per participant ($N_{\text{sessions}} = 5$).

Given the screen refresh rate of 60 Hz and the $K$-bit stimulation sequence ($K = 63$), each visual epoch corresponded to 538 time points ($N_t = 538$), equivalent to approximately 1.05 s per trial. For each subject, $N_s$ electrodes were used, and six stimulus classes ($N_{\text{classes}} = 6$) corresponding to circular code shifts of 0, 8, 16, 24, 32, and 40 bits (see Fig. 3) were considered.
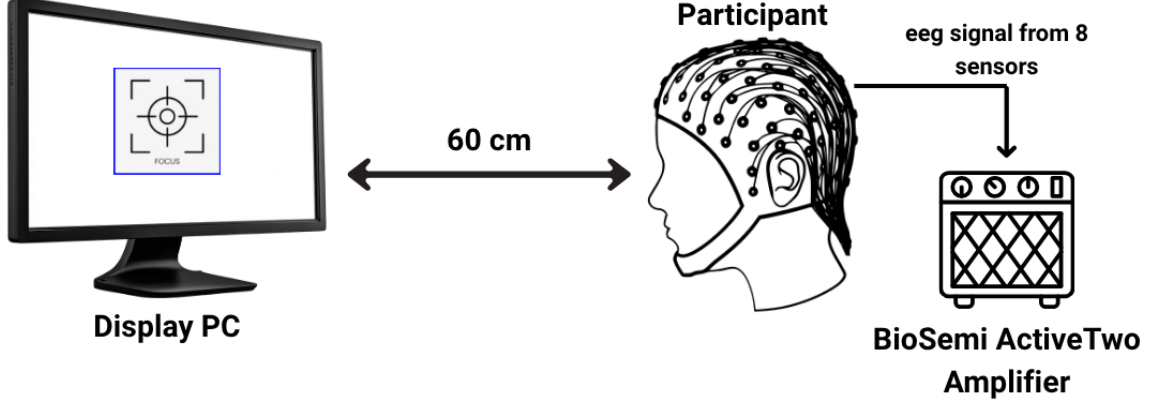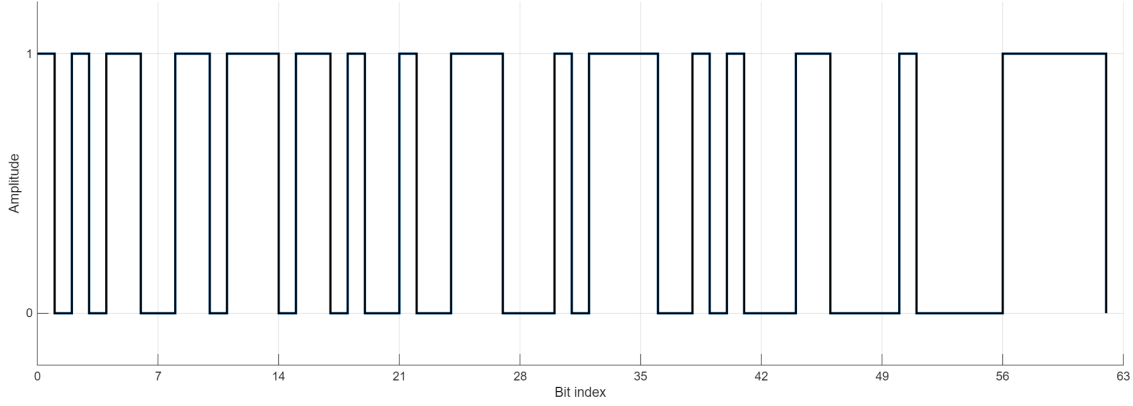
Figure 1: Schematic of the experimental setup.



Figure 2: Binary representation of the $K$-bit $m$-sequence stimulus pattern used for C-VEP paradigm. The sequence alternates between 0 and 1 according to the pseudorandom maximal-length code, forming the base temporal modulation for the visual stimulus.

### 3.1.3 Experimental design

The calibration interface (see Fig. 5) consisted of a single box flickering on a 27-inch screen (LG27GL83A, 1920 × 1080, 350 $cd/m^2$, 60 Hz refresh rate). The stimulus was 192 × 192, corresponding to about 6.22 $cd/m^2$. The flickering pattern was modulated using a K-bit m-sequence code [38], alternating between an on/off state (black/white image on the screen). The m-sequence code(see Fig. 2) used was: '10101100110111011010010011000101111100101000110 0001000001111110', where alternated between two states: a calibration image displayed for code 1 and a black screen for code 0. This distinct pattern ensured robust stimulation for generating C-VEP signals [59]. The calibration interface, which displayed a flickering single box, was designed and implemented in MATLAB and Psychtoolbox to ensure precise timing and synchronization of the L-bit m-sequence visual stimuli [60].

### 3.2 Signal Processing

The raw EEG signals were first scaled to microvolts ($\mu V$) and preprocessed to enhance signal quality. Each channel was detrended to remove slow-varying baseline drifts and low-frequency fluctuations, ensuring a zero-mean baseline across all epochs. This step was essential for eliminating slow DC shifts and stabilizing the data for subsequent feature extraction and classification. Following detrending, a bandpass filter with cutoff frequencies of $f_c^{low} = 0.5$ Hz and $f_c^{high} = 42.66$ Hz was applied to preserve the frequency range relevant to visual evoked potentials while attenuating muscle noise and low-frequency artifacts. The filtered signals were then segmented into stimulus-locked epochs of $N_t$ samples.
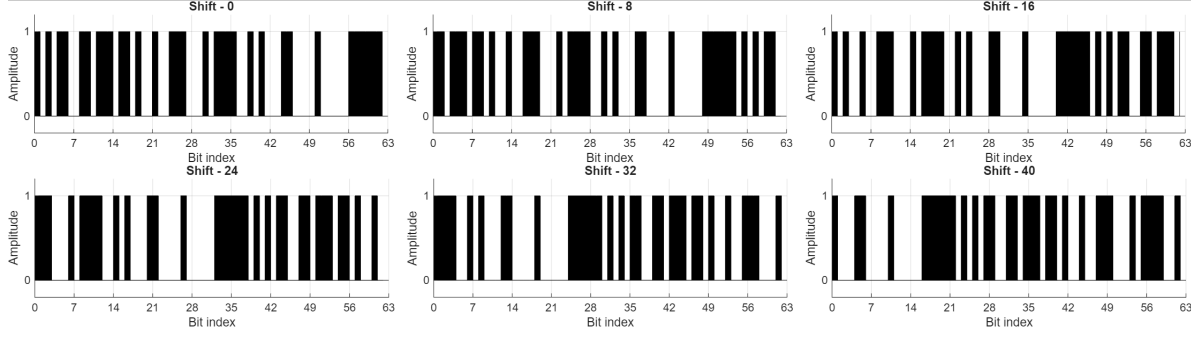
Figure 3: Binary representation of the $K$-bit $m$-sequence and its five circularly shifted variants (shifts of 8, 16, 24, 32, and 40 bits).
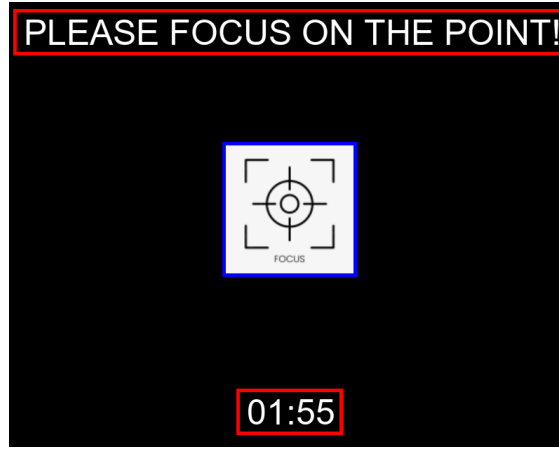


Figure 4: Visual stimuli presented on the screen for recording the calibration data.

### 3.2.1 Surface Laplacian Filtering with Great-Circle Distance

Sensors are arranged on a spherical head model to account for scalp curvature. In this configuration, the distances between electrodes are calculated using the great-circle distance (orthodromic distance) [61, 62, 63], ensuring that spatial relationships are preserved on the scalp surface. This step is fundamental for determining the influence of neighboring electrodes during signal processing. We denote by $d_{ij}^{\text{GCD}}$ the great-circle distance between sensor $i$ and sensor $j$.

The great-circle distance between two points $(\phi_i, \lambda_i)$ and $(\phi_i, \lambda_i)$ on a sphere of radius $r$ is given by:

$$d_{ij}^{\text{GCD}} = r \cdot \Delta\sigma, \tag{1}$$

where

$$\Delta\sigma = \text{atan2}\left(\sqrt{a_1 + a_2}, \, a_3\right), \tag{2}$$

with

$$\begin{aligned}
a_1 &= (\cos\phi_j \sin\Delta\lambda)^2, \\
a_2 &= (\cos\phi_i \sin\phi_j - \sin\phi_i \cos\phi_j \cos\Delta\lambda)^2, \\
a_3 &= \sin\phi_i \sin\phi_j + \cos\phi_i \cos\phi_j \cos\Delta\lambda.
\end{aligned} \tag{3}$$

With $\Delta\lambda = \lambda_j - \lambda_i$ being the difference in longitude between the two sensors. This is based on the Vincenty formula for great-circle distance.

The weights $w_{ij}$ determine the influence of neighboring electrodes and are computed using the following equation:
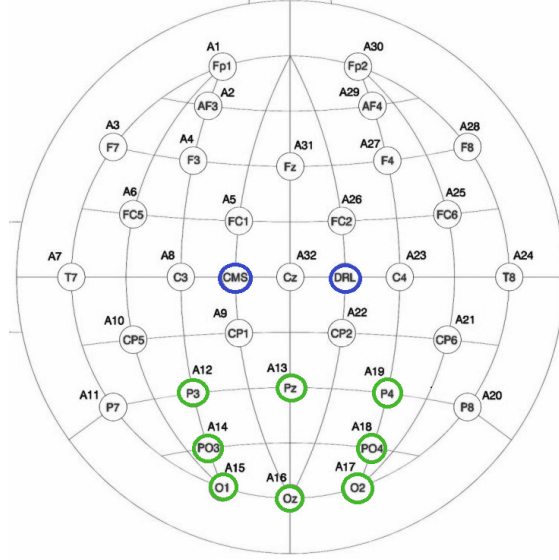
$$w_{ij} = \frac{1}{d_{ij}^{\text{GCD}}} \tag{4}$$

Figure 5: BioSemi 32+2 electrode layout. The highlighted electrodes(marked in blue(8 electrodes) and green(ground)) indicate those used for C-VEP signal acquisition.

The Laplacian-filtered signal $S'_i$ is then estimated as:

$$S'_i = S_i - \sum_{j \neq i} w_{ij} \cdot S_j \tag{5}$$

where $S_i$ is the signal at the target electrode $i$ and $S_j$ is the signal from neighboring electrodes $1 \leq j \leq N_s$, $i \neq j$, where $N_s$ is the number of sensors. This transformation enhances local neural activity by reducing the effect of distant sources. The filter applies weights to electrodes within a defined radius of 1 (in normalized units).

### 3.3 Feature extraction

#### 3.3.1 Canonical Correlation Analysis (CCA)

Canonical Correlation Analysis (CCA) was employed to extract correlation-based features that quantify the relationship between the recorded EEG signals and reference templates corresponding to each visual stimulus shift[64, 65, 66]. The resulting canonical correlation coefficients $\rho = [\rho_1, \ldots, \rho_{n_f}]$ capture the linear associations between the EEG projections and stimulus templates across electrodes.

Because one of the eight EEG channels exhibited high redundancy or excessive noise, the effective rank of the data matrix was reduced to seven. Consequently, only seven canonical correlations were computed per trial, reflecting the number of independent spatial components. To maintain a uniform feature length across all samples, an additional zero value was appended, yielding consistent eight-dimensional vectors per shift. Concatenating these across all six code models produced a $1 \times 48$ feature vector for each example. The resulting CCA features captured spatially distinct stimulus-specific correlations while ensuring dimensional consistency for subsequent Bayesian Linear Discriminant Analysis (BLDA) classification.

#### 3.3.2 Correlation Coefficients

Correlation-based feature extraction was employed to quantify the similarity between the recorded EEG responses and the expected m-sequence templates associated with each visual stimulus [67]. For each subject, $N_s$ electrodes were used, and $N_{classes}$ stimulus classes corresponding to circular code shifts of 0, 8, 16, 24, 32, and 40 bits were considered. Each class contained eight template signals, resulting in a total of $N_s \times N_s \times N_{classes} = 384$ correlation features per trial.

Each data block comprised 114 trials($N_{examples} = 114$), each representing an L-bit stimulation period. After applying six circular shifts to represent the $N_{classes}$ classes, a total of 684 examples were generated per block ($N_{examples}$ per class). The extracted correlation features were organized into training and testing matrices, $X_{\text{train}}$ and $X_{\text{test}}$, each of

dimension $684 \times 384$. These correlation coefficients served as discriminative representations of the neural response patterns corresponding to each visual stimulus, facilitating robust classification across subjects and sessions.

### 3.3.3 Deep Learning Features

The early convolutional layers of the CNN were used to learn discriminative spatial–temporal representations from raw EEG signals and it consisted of 11 layers with approximately 8.6k learnable parameters. The network began with an input layer of size $N_s \times N_t \times 1$, corresponding to the eight spatial channels and $N_t$ temporal samples. A spatial convolution layer with $N_s \times 1$ kernels was applied to capture inter-channel dependencies and enhance localized spatial information.

Two subsequent temporal convolutional blocks progressively extracted dynamic temporal patterns from the EEG. The first block used $1 \times 11$ kernels with 16 filters, followed by a ReLU activation, $1 \times 2$ max pooling, and a 20% dropout for regularization. The second block employed $1 \times 14$ kernels with 32 filters, also followed by ReLU activation, $1 \times 2$ max pooling, and a 30% dropout. These initial layers generated compact, informative spatial–temporal feature maps that were later used for classification and similarity learning.

## 3.4 Classifiers

### 3.4.1 Bayesian Linear Discriminant Analysis (BLDA)

BLDA extends traditional Linear Discriminant Analysis by introducing a Bayesian regularization framework that automatically balances model complexity and generalization performance[68, 69, 70]. A multiclass BLDA approach was implemented using a one-vs-rest strategy, where an independent binary classifier was trained for each class. Each BLDA model estimates class-conditional distributions and computes posterior probabilities based on learned discriminant functions, providing robust decision boundaries even in high-dimensional and small-sample settings. During testing, the class with the highest posterior score was selected as the predicted label, while normalized class probabilities served as confidence estimates.

### 3.4.2 Distances

The first CNN architecture was designed to reconstruct the 63-bit stimulus code directly from single-trial EEG recordings. Each input trial consisted of an $N_s \times N_t$ segment representing eight spatial EEG channels across $N_t$ temporal samples. The network output was a $K$-dimensional vector corresponding to the binary stimulation sequence used in the c-VEP paradigm. A sigmoid activation was applied to each output neuron to produce continuous values in the range [0, 1], enabling bitwise decoding and distance-based evaluation against the ground-truth $6K$-bit templates.

The network was trained using the Adam optimizer with an initial learning rate of $1 \times 10^{-3}$ and an $L_2$ regularization factor of $1 \times 10^{-4}$ to prevent overfitting. Training was performed for up to 40 epochs with a mini-batch size of 64, and the data were shuffled at each epoch to improve generalization. Validation was conducted after each epoch using a held-out set $(X_{vdl}, Y_{vdl})$, and the root mean square error (RMSE) was monitored as the primary performance metric.

The network consisted of 2 more layers added to 3.3.3 with 127.1k added learnable parameters. The first few layers as mentioned in 3.3.3 (as in Fig. 6)produced compact and informative spatial–temporal feature maps. These feature maps were flattened and passed through a fully connected layer with $K$ neurons, one corresponding to each bit in the stimulation code. Finally, a sigmoid activation layer generated the output vector of size $K \times 1$, representing the reconstructed stimulus sequence.

**Euclidean Distance Evaluation** To evaluate network performance, a distance-based decoding approach was employed, comparing the CNN-reconstructed bit sequence with the reference code templates for each visual stimulus. For each trial $i$, the network outputs a continuous-valued vector $\hat{\mathbf{y}}_i \in [0, 1]^K$, where $K = 63$ represents the number of bits in the c-VEP stimulation code. This vector can be interpreted as the predicted probability or confidence associated with each bit position.

The Euclidean distance between the reconstructed output $\bar{\mathbf{y}}_i$ (thresholded or mean-centered version of $\hat{\mathbf{y}}_i$) and each reference code template $\mathbf{c}^{(N_{classes})}$ was computed as

$$d_{i,N_{classes}}^{(2)} = \left\| \bar{\mathbf{y}}_i - \mathbf{c}^{(N_{classes})} \right\|_2 = \sqrt{\sum_{k=1}^{K} \left( \bar{y}_{i,k} - c_b^{(N_{classes})} \right)^2}. \tag{6}$$

Each distance $d_{i,N_{classes}}^{(2)}$ quantifies the dissimilarity between the reconstructed sequence and the expected codeword for $N_{classes}$. The predicted class for trial $i$ is determined as the one with the minimum Euclidean distance to the output

vector:

$$N^{\star}_{classes_i} = \underset{N_c=1,\ldots,N_{classes}}{\mathrm{argmin}} d^{(2)}_{i,N_c} = \underset{N_c}{\mathrm{argmin}} \sum_{k=1}^{K}\big(\bar{y}_{i,k} - c^{(N_c)}_k\big)^2, \tag{7}$$

where the square root is omitted without affecting the decision outcome.

This metric provides an intuitive, computationally efficient means of classifying reconstructed bit patterns, as smaller Euclidean distances indicate greater similarity between the predicted output and the true stimulus code.

**Mahalanobis Distance Evaluation.** The Mahalanobis distance was employed to account for feature correlations and varying variance among the reconstructed bits. This approach provides a more statistically informed similarity measure between the reconstructed output and the reference code templates. For each trial $i$, the CNN output vector $\hat{\mathbf{y}}_i \in [0,1]^K$ was compared against each class codeword $\mathbf{c}^{(N_c)}$, using a shared covariance matrix estimated from the reconstructed outputs across all samples. The covariance matrix $\mathbf{\Sigma}$ was regularized to ensure numerical stability via a shrinkage factor $\lambda = 0.1$, forming

$$\mathbf{\Sigma}_{\mathrm{reg}} = (1-\lambda)\mathbf{\Sigma} + \lambda\,\alpha\,\mathbf{I}, \tag{8}$$

where $\alpha$ denotes the mean of the diagonal elements of $\mathbf{\Sigma}$ and $\mathbf{I}$ is the identity matrix. The Mahalanobis distance between the reconstructed output and each code template was then computed as

$$d^{(M)}_{i,N_c} = \sqrt{(\bar{\mathbf{y}}_i - \mathbf{c}^{(N_c)})^{\top}\mathbf{\Sigma}^{-1}_{\mathrm{reg}}(\bar{\mathbf{y}}_i - \mathbf{c}^{(N_c)})}. \tag{9}$$

The predicted class label was assigned according to the minimum Mahalanobis distance:

$$N^{\star}_{classes_i} = \underset{N_c=1,\ldots,N_{classes}}{\mathrm{argmin}} d^{(M)}_{i,N_c}. \tag{10}$$

**Earth Mover's Distance (EMD) Evaluation.** The EMD was employed as a distribution-based metric. Unlike Euclidean or Mahalanobis distances, which measure pointwise differences, EMD quantifies the minimal cumulative "cost" required to transform one probability distribution into another, providing a more perceptually meaningful measure of dissimilarity for structured signals such as the $K$-bit c-VEP codes. For each trial $i$, the CNN output $\hat{\mathbf{y}}_i \in [0,1]^K$ and each reference template $\mathbf{c}^{(N_c)}$ were first normalized to form discrete probability mass functions (PMFs),

$$\mathbf{p}_i = \frac{\max(0,\,\hat{\mathbf{y}}_i)}{\sum_b \max(0,\,\hat{y}_{i,b})}, \qquad \mathbf{q}_{N_c} = \frac{\mathbf{c}^{(N_c)}}{\sum_b c^{(N_c)}_b}, \tag{11}$$

ensuring that $\sum_b p_{i,b} = \sum_b q_{N_c,b} = 1$. The one-dimensional EMD between $\mathbf{p}_i$ and each class template $\mathbf{q}_{N_c}$ was computed using the cumulative distribution formulation,

$$d^{(\mathrm{EMD})}_{i,N_c} = \frac{1}{K}\sum_{k=1}^{K}\big|\,\mathrm{CDF}_{p_i}(k) - \mathrm{CDF}_{q_{N_c}}(k)\big|, \tag{12}$$
$$= \mathrm{mean}\big(\big|\mathrm{cumsum}(\mathbf{p}_i - \mathbf{q}_{N_c})\big|\big).$$

The predicted class was determined by selecting the template with the minimum EMD value:

$$N_{classes_i}{}^{\star} = \underset{N_c=1,\ldots,N_{classes}}{\mathrm{argmin}} d^{(\mathrm{EMD})}_{i,N_c}. \tag{13}$$

By comparing cumulative distributions rather than individual bits, EMD captures global structural deviations between reconstructed and reference codes.

**Constrained-EMD Evaluation.** Contrained-EMD introduces a movement radius $R$, restricting the transport of probability mass between bit positions to a local neighborhood, thereby emphasizing short-range structural correspondence in the reconstructed codes. For each trial $i$, the predicted vector $\hat{\mathbf{y}}_i \in [0,1]^K$ was normalized into a probability mass function (PMF), and compared against the set of reference code templates $\mathbf{c}^{(N_c)}$, each also normalized to form $\mathbf{q}_{N_c}$. The resulting optimization seeks the minimum transport cost required to transform $\mathbf{p}_i$ into $\mathbf{q}_{N_c}$, under the constraint that transport is only allowed between positions whose bit indices differ by at most $R$:

$$\underset{\mathbf{X}\geq 0}{\min} \sum_{k,j}|k-j|\,X_{k,j} \quad \text{s.t.} \quad \sum_j X_{k,j} = p_{i,k},$$
$$\sum_k X_{k,j} = q_{k,j}, \quad X_{k,j} = 0 \text{ if } |k-j| > R. \tag{14}$$

Here, $\mathbf{X}$ denotes the transport flow matrix between the distributions $\mathbf{p}_i$ and $\mathbf{q}_{N_c}$. The problem was solved using a linear programming formulation based on the dual-simplex algorithm with equality constraints enforcing mass conservation. When $R$ is large ($R \geq K - 1$), the formulation naturally reduces to the unconstrained 1-D EMD case based on cumulative distribution differences.

The constrained-EMD distance for each class $N_c$ was computed as the minimal transport cost normalized by the total flow, and the predicted label was assigned to the class with the smallest constrained distance:

$$N^{\star}_{classes_i} = \underset{N_c=1,\ldots,N_{classes}}{\mathrm{argmin}} \; d^{(\mathrm{CEMD})}_{i,N_c}. \tag{15}$$

### 3.4.3 CNN

A CNN-based $N_{classes}$ classification architecture was designed to classify EEG responses directly into six distinct classes corresponding to the six circular shifts (0, 8, 16, 24, 32, and 40 bits) of the $K$-bit m-sequence used in the c-VEP stimulation.

The network was trained using the Adam optimizer with an initial learning rate of $1 \times 10^{-3}$ and an $L_2$ regularization coefficient of $1 \times 10^{-4}$ to mitigate overfitting. The model was trained for 40 epochs with a mini-batch size of 64, and data was shuffled at each epoch to ensure robust learning. The training employed categorical cross-entropy as the loss function and accuracy as the primary evaluation metric. The network comprised 3 more layers added to 3.3.3 with approximately 12.1k added learnable parameters. The feature maps from 3.3.3 were flattened into a one-dimensional representation and passed to a fully connected layer with six neurons. The Softmax activation at the output produced a probability distribution over $N_{classes}$, enabling direct classification of the attended stimulus (see Fig. 6).

The trained network generated predictions for the test set via forward propagation, producing a $N_{classes} \times N_{examples}$ matrix of class probabilities for each trial. The predicted class label $\hat{y}$ for each trial was obtained as the index of the highest confidence value.
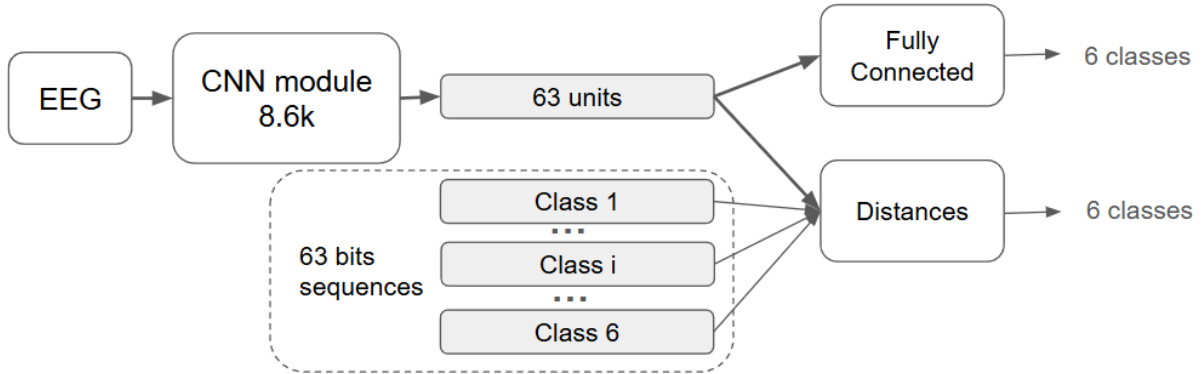


Figure 6: Schematic representation of the CNN models used for direct classification of the responses and to be used with distances.

### 3.4.4 Siamese Networks

**Single multi-class Siamese (Twin) Network**  The third architecture is a convolutional Siamese (twin) network configuration designed to learn a similarity metric between pairs of EEG trials from different classes. Unlike previous architectures that performed classification directly, this model learns a shared embedding space in which trials from the same stimulus class produce similar representations, while trials from different classes are mapped farther apart. Each input pair consisted of two single-trial EEG segments of size $N_s \times N_t$, processed through identical CNN Modules with shared weights. The Siamese model output was a single probability representing the likelihood that the two inputs originated from the same c-VEP code shift.

The subnetwork consisted of one additional layer to 3.3.3, with approximately 127.1k additional parameters followed by the same spatial–temporal design as the previous CNNs. The feature maps from 3.3.3 were flattened and projected into a $K$-dimensional embedding space, which served as the compact feature representation for each EEG trial.

The similarity estimation between the two inputs was computed as the absolute difference between their embeddings, followed by a fully connected layer and a sigmoid activation producing a scalar similarity probability (see Fig 7). The network was trained using the binary cross-entropy loss between predicted similarity scores and binary pair labels (1 for
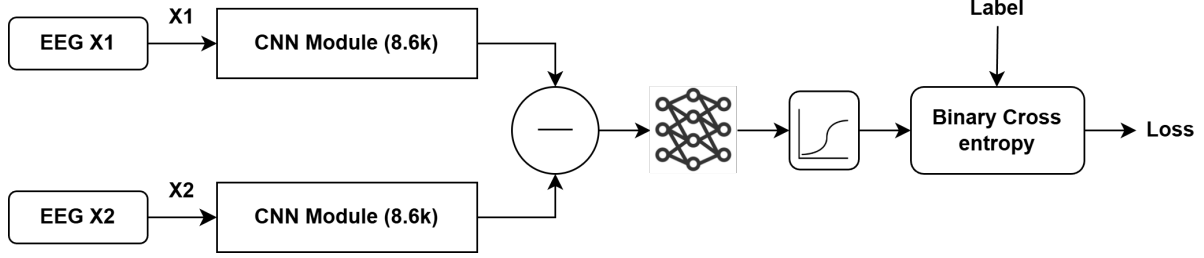
Figure 7: Schematic representation of the Siamese (Twin) CNN architecture. The model consists of two identical CNN Modules that process paired EEG inputs ($X_1$, $X_2$) with shared weights.

matching class, 0 otherwise). Training was conducted for 30 epochs with a batch size of 256 randomly balanced pairs per iteration. The Adam optimizer was used with a learning rate of $1 \times 10^{-3}$, and the loss was backpropagated through both subnetworks simultaneously to update the shared parameters.

**Multi-classifier binary Siamese Network**   The final configuration extended the convolutional Siamese framework to a multi-classifier setting, where separate Siamese models were trained for each of the $N_{classes}$. While the underlying architecture remained identical to the Single multi-class Siamese network described in Section 3.4.4, the training procedure was modified to produce class-specific similarity models that collectively enable multi-class discrimination. Each individual model was trained to distinguish between matching and non-matching EEG trial pairs associated with a single target code shift, effectively learning a dedicated embedding subspace for that particular class.

For each class $k \in \{1, \ldots, N_{classes}\}$, a Siamese model was trained using EEG trial pairs drawn from within and outside the class. The training data consisted of balanced positive (same-class) and negative (different-class) pairs, ensuring robust similarity learning. Training was conducted for 30 epochs with a batch size of 256 pairs, using the Adam optimizer with a learning rate of $1 \times 10^{-3}$.

During inference, the trained set of six Siamese models was evaluated independently. For each test pair, the similarity probability was computed for all $N_{classes}$, each yielding a binary classification output.

### 3.5   Performance Analysis

The performance of the proposed methods was evaluated using a five-fold cross-validation scheme based on $N_{sessions}$ per subject. In each fold, one recording consisting of $N_{examples}$ was used for testing, while the remaining $N_{sessions} - 1$ ($4 \times N_{examples}$) were used for training. The results reported in the subsequent tables represent the mean accuracy and standard deviation computed across the five folds for each subject.

## 4   Results

### 4.1   Grand Average Plots

To visualize the consistency and reliability of neural responses across participants and sessions, grand-average analyses were performed. Fig. 8 depicts time-domain plots of EEG signals from eight selected electrodes showing the mean response across trials . The shaded envelopes represent the standard error across participants, illustrating the trial-to-trial variability and consistency of neural responses. These plots highlight the stability of the signal features used for template-based classification in the C-VEP BCI system.

Fig 9 provides an overview of the signals and the reconstructed $K$-bit outputs obtained from the proposed CNN-based model. Each subplot corresponds to a single participant and displays the CNN's $K$-bit output, averaged across all trials and sessions. The dark gray line indicates the mean reconstructed signal, the shaded gray area represents the $\pm$SD range, and the dashed trace denotes the original $K$-bit $m$-sequence used for stimulation. The final panel shows the overall grand average and variability across subjects. The averaged waveforms highlight the network's ability to capture the temporal structure of the stimulation code, offering insight into the spatial–temporal patterns that underpin classification performance.
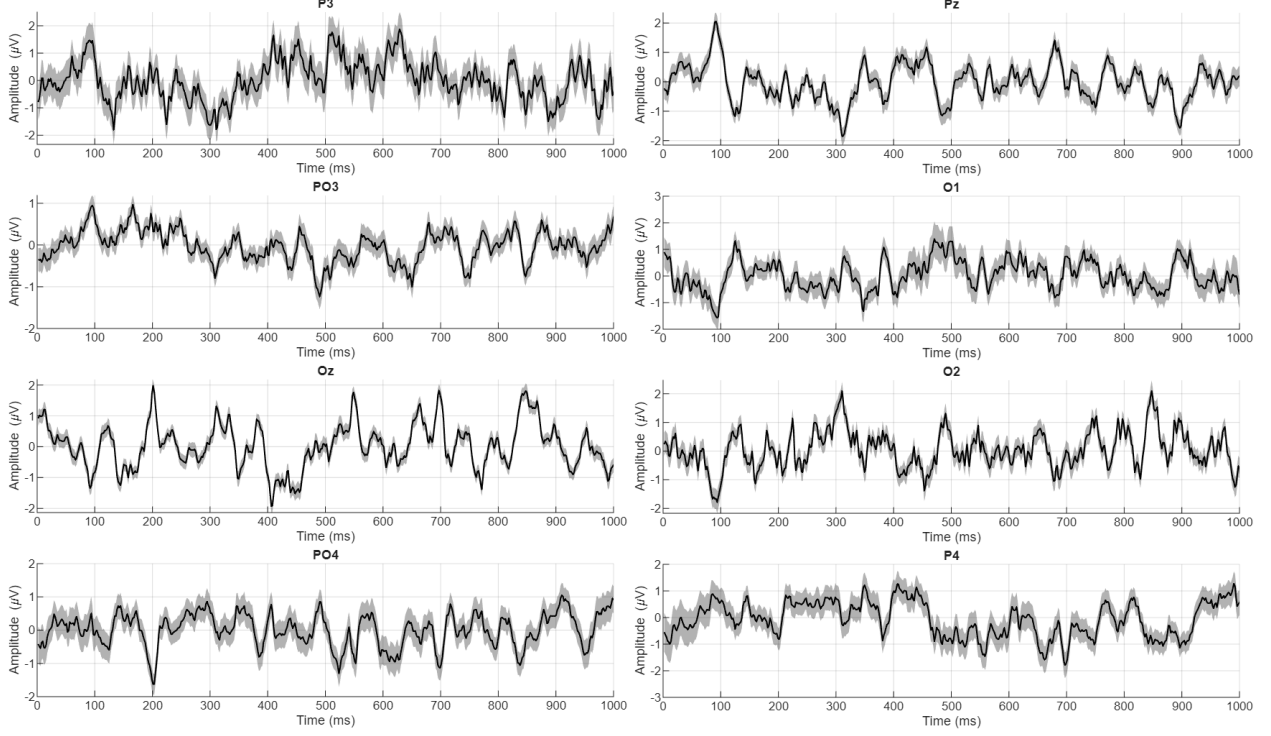
10

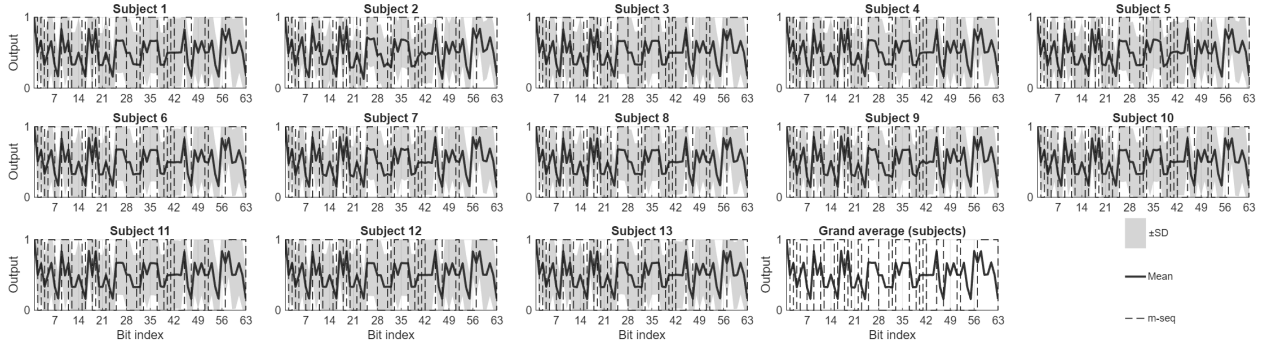Figure 8: Time-domain plots of EEG signals from eight selected electrodes.



Figure 9: Grand-average of the CNN's $K$-bit outputs across all subjects.

## 4.2 Classifier Performance

The classification accuracy across 13 participants using the four methodological categories, correlation-based, canonical correlation-based, CNN-based, and Siamese-based networks, is summarized in Table 1, Table 2, Table 3, and Table 4, respectively. Among the conventional baselines, Corr+BLDA achieved a mean accuracy of $81.64 \pm 14.02\%$, outperforming CCA+BLDA ($62.59 \pm 16.57\%$), confirming the advantage of correlation-driven feature representation with Bayesian regularization for C-VEP decoding.

In contrast, all deep learning models exhibited substantially higher performance, demonstrating the efficacy of end-to-end feature learning. The CNN-based $K$-bit models achieved accuracies of $93.86 \pm 6.16\%$ (Euclidean), $92.38 \pm 8.34\%$ (Mahalanobis), $93.88 \pm 6.38\%$ (EMD), and $91.30 \pm 7.99\%$ (Constrained-EMD). The CNN+$N_{classes}$ classifier reached $91.42 \pm 7.93\%$, indicating strong generalization despite reduced output dimensionality. The Siamese networks yielded $87.29 \pm 8.00\%$ (Single) and $96.89 \pm 2.74\%$ (multiplewith the latter achieving the highest overall mean accuracy overall. These results highlight that deep architectures, particularly distance-based CNNs and the Siamese multiple-classifier variant, significantly outperform traditional BLDA methods.

11

Table 1: Classification accuracy (%) comparison between feature extraction using Correlation coefficients and classification using BLDA and feature extraction using CCA and classification using BLDA across subjects.

| Subject | Corrcoef+BLDA | CCA+BLDA |
|---|---|---|
| 1 | $88.72 \pm 6.83$ | $59.39 \pm 6.96$ |
| 2 | $63.36 \pm 25.35$ | $59.79 \pm 6.73$ |
| 3 | $99.24 \pm 0.38$ | $88.68 \pm 4.23$ |
| 4 | $91.94 \pm 4.00$ | $79.14 \pm 4.48$ |
| 5 | $50.21 \pm 9.38$ | $66.42 \pm 9.54$ |
| 6 | $90.11 \pm 2.19$ | $67.92 \pm 6.93$ |
| 7 | $73.42 \pm 17.38$ | $45.00 \pm 11.20$ |
| 8 | $77.89 \pm 3.97$ | $32.14 \pm 4.29$ |
| 9 | $71.12 \pm 28.93$ | $50.35 \pm 17.49$ |
| 10 | $84.66 \pm 4.94$ | $49.66 \pm 7.11$ |
| 11 | $87.81 \pm 3.13$ | $53.89 \pm 3.83$ |
| 12 | $97.17 \pm 1.40$ | $80.90 \pm 5.64$ |
| 13 | $85.65 \pm 6.30$ | $83.33 \pm 3.26$ |
| **Mean** | $\mathbf{81.64 \pm 8.78}$ | $\mathbf{62.59 \pm 7.05}$ |
| **SD** | $\mathbf{14.02 \pm 9.25}$ | $\mathbf{16.57 \pm 3.89}$ |

Table 2: Performance comparison (%) of the $K$-bit reconstruction CNN using different distance evaluation metrics.

| Subject | Euclidean Evaluation | Mahalanobis Evaluation | EMD Evaluation | Constrained EMD Evaluation |
|---|---|---|---|---|
| 1 | $99.39 \pm 0.72$ | $98.86 \pm 1.14$ | $98.04 \pm 1.29$ | $97.60 \pm 2.14$ |
| 2 | $83.21 \pm 29.22$ | $71.39 \pm 32.64$ | $80.88 \pm 27.26$ | $70.08 \pm 29.87$ |
| 3 | $99.68 \pm 0.39$ | $99.62 \pm 0.43$ | $99.44 \pm 0.36$ | $99.27 \pm 0.36$ |
| 4 | $98.48 \pm 1.05$ | $96.90 \pm 2.05$ | $95.91 \pm 1.59$ | $95.35 \pm 2.52$ |
| 5 | $86.42 \pm 6.69$ | $87.03 \pm 7.24$ | $83.75 \pm 7.38$ | $85.25 \pm 6.66$ |
| 6 | $96.98 \pm 0.97$ | $95.99 \pm 1.89$ | $95.11 \pm 1.11$ | $94.73 \pm 1.78$ |
| 7 | $84.95 \pm 21.99$ | $84.83 \pm 23.63$ | $83.96 \pm 20.97$ | $84.46 \pm 20.95$ |
| 8 | $94.30 \pm 3.06$ | $93.54 \pm 1.89$ | $91.26 \pm 3.54$ | $91.11 \pm 2.66$ |
| 9 | $86.81 \pm 22.04$ | $85.09 \pm 24.34$ | $93.13 \pm 25.12$ | $83.57 \pm 23.92$ |
| 10 | $95.06 \pm 3.55$ | $94.71 \pm 3.17$ | $93.45 \pm 3.42$ | $93.25 \pm 3.30$ |
| 11 | $97.57 \pm 0.80$ | $97.22 \pm 0.48$ | $96.40 \pm 0.91$ | $96.59 \pm 0.62$ |
| 12 | $99.09 \pm 0.76$ | $99.04 \pm 0.76$ | $98.77 \pm 0.82$ | $98.65 \pm 0.91$ |
| 13 | $98.30 \pm 0.77$ | $98.22 \pm 0.64$ | $97.16 \pm 1.12$ | $97.37 \pm 1.10$ |
| **Mean** | $\mathbf{93.86 \pm 7.08}$ | $\mathbf{92.38 \pm 7.57}$ | $\mathbf{93.88 \pm 7.84}$ | $\mathbf{91.30 \pm 7.55}$ |
| **SD** | $\mathbf{6.16 \pm 10.18}$ | $\mathbf{8.34 \pm 10.24}$ | $\mathbf{6.38 \pm 8.87}$ | $\mathbf{7.99 \pm 8.98}$ |

Table 3: Classification accuracy (%) of the $N_{classes}$ CNN across all subjects.

| Subject | CNN + $N_{classes}$ |
|---|---|
| 1 | $97.43 \pm 1.75$ |
| 2 | $77.17 \pm 30.88$ |
| 3 | $99.30 \pm 0.30$ |
| 4 | $95.58 \pm 2.28$ |
| 5 | $79.70 \pm 10.14$ |
| 6 | $94.79 \pm 1.76$ |
| 7 | $80.37 \pm 23.20$ |
| 8 | $92.49 \pm 2.80$ |
| 9 | $84.82 \pm 24.26$ |
| 10 | $94.42 \pm 3.63$ |
| 11 | $97.19 \pm 0.73$ |
| 12 | $98.51 \pm 0.82$ |
| 13 | $96.73 \pm 1.22$ |
| **Mean** | $\mathbf{91.42 \pm 7.98}$ |
| **SD** | $\mathbf{7.93 \pm 10.76}$ |

Table 4: Performance comparison of Siamese network configurations: Single multi-class Siamese vs. Multiple-classifier binary Siamese.

| Subject | Multi-class Siamese | Multi-classifier Siamese |
|:---:|:---:|:---:|
| 1 | $92.67 \pm 3.60$ | $98.39 \pm 0.93$ |
| 2 | $75.64 \pm 15.17$ | $92.73 \pm 5.28$ |
| 3 | $97.90 \pm 1.58$ | $99.50 \pm 0.21$ |
| 4 | $88.60 \pm 4.02$ | $97.35 \pm 0.94$ |
| 5 | $69.05 \pm 5.55$ | $89.63 \pm 3.67$ |
| 6 | $85.24 \pm 3.31$ | $96.84 \pm 0.70$ |
| 7 | $84.96 \pm 7.38$ | $93.52 \pm 5.60$ |
| 8 | $84.99 \pm 3.04$ | $96.30 \pm 1.20$ |
| 9 | $87.33 \pm 10.24$ | $94.97 \pm 5.97$ |
| 10 | $93.27 \pm 2.92$ | $97.27 \pm 0.93$ |
| 11 | $95.53 \pm 1.73$ | $98.38 \pm 0.23$ |
| 12 | $93.34 \pm 6.60$ | $99.05 \pm 0.35$ |
| 13 | $86.27 \pm 1.65$ | $97.17 \pm 1.23$ |
| **Mean** | $\mathbf{87.29 \pm 5.14}$ | $\mathbf{96.89 \pm 1.99}$ |
| **SD** | $\mathbf{8.00 \pm 3.95}$ | $\mathbf{2.74 \pm 2.05}$ |

A non-parametric Friedman test was used to compare the nine methods across 13 subjects, using each subject's mean classification accuracy as the dependent variable. The analysis revealed a significant difference among the classifiers, $\chi^2(8) = 85.78$, $p < 0.001$, indicating that not all models performed equivalently. The corresponding Kendall's coefficient of concordance ($W = 0.825$) suggested a large effect size, confirming strong agreement in performance ranking across participants. The average ranks were as follows: CCA+BLDA (1.08), Corr+BLDA (2.23), Siamese (multi-class) (3.62), CNN–$K$ bit (Constrained-EMD) (4.46), CNN+$N_{classes}$ (4.77), CNN–$K$ bit (EMD) (5.62), CNN–$K$ bit (Mahalanobis) (6.85), CNN–$K$ bit (Euclidean) (8.15), and Siamese (multiple-classifier) (8.23).

Post-hoc pairwise comparisons were performed using the Wilcoxon signed-rank test with Bonferroni correction ($\alpha_{\text{adj}} = 0.0014$) to control for multiple testing. Significant differences ($p < 0.05$ after correction) were observed between the traditional BLDA-based baselines and most CNN and Siamese models, confirming that the deep learning approaches produced statistically distinct results. In particular, Corr+BLDA and CCA+BLDA yielded significantly lower performance ($p = 0.0088$) compared with nearly all CNN variants and the Siamese multiple-classifier model. Among the CNN–$K$ bit configurations, the Constrained-EMD and EMD distance metrics did not differ significantly, whereas both were superior to the Euclidean ($p = 0.0088$) and Mahalanobis ($p = 0.0088$) metrics. The CNN+$N_{classes}$ model performed comparably to the Constrained-EMD variant, showing no statistically significant difference between them.

Overall, the statistical analysis demonstrates that the proposed deep architectures outperform the classical correlation- and CCA-based BLDA methods. Within the deep learning group, the CNN —$K$ bit (Constrained-EMD) and Siamese (multiple-classifier) networks exhibited the most consistent performance across participants. In contrast, the Siamese (multi-class) network achieved significantly lower accuracies ($p = 0.0088$) than almost all other methods. These findings confirm that incorporating spatial–temporal feature learning with distance-based decoding yields a more robust and generalizable representation of C-VEP responses.

### 4.3 Data Augmentation and classifier combination

To evaluate the robustness of the proposed deep learning models against temporal variability in EEG signals, a series of controlled data augmentation experiments were conducted. Temporal augmentation was performed by shifting the EEG signal in the time domain by a fixed number of samples to simulate small latency variations in the visual response. Specifically, each trial was shifted by $\pm 1$, $\pm 2$, $\pm 4$, and $\pm 8$ time points, and new training and testing sets were generated for each shift condition. Three augmentation strategies were examined: (i) Train Augmentation (TA), where only the training data were augmented; (ii) Test Combination (TC), where the score from different shifted inputs are averaged; and (iii) Train Augmentation and Test Combination (TA&TC), where the training set was augmented and the scores were combined. The baseline condition (NA) corresponds to models trained and evaluated without temporal shifting.

This augmentation framework was applied consistently across all deep learning architectures. As illustrated in Figs. 10, 11, 12a, and 12b, temporal augmentation up to $\alpha = 4$ samples improved overall generalization and reduced performance variability across sessions. In contrast, excessive temporal shifting ($\alpha = 8$) led to a gradual decline in accuracy for most methods, particularly when augmentation was applied only to the testing data (TC).
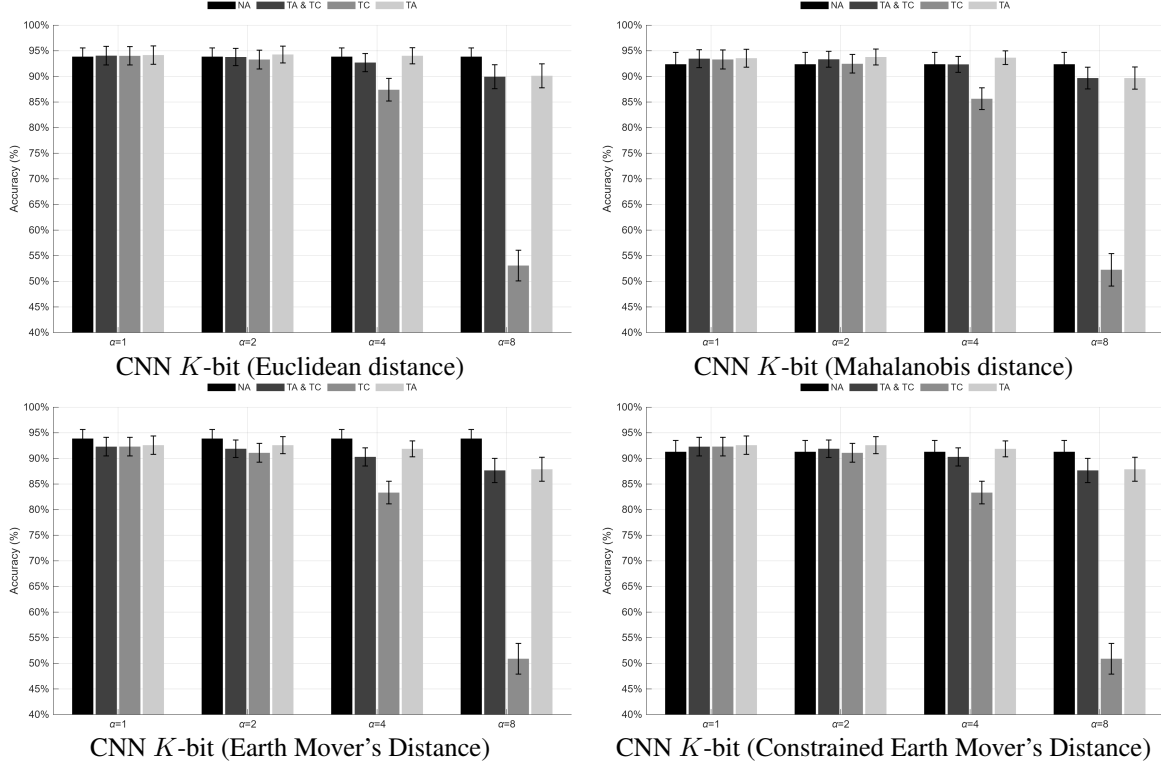
Figure 10: Performance comparison of CNN $K$-bit models under different distance metrics and temporal augmentation conditions.

Across all models, the first column of each accuracy matrix represents the baseline (NA) condition, followed by TA&TC, TC, and TA configurations. Under moderate temporal shifts ($\alpha \leq 4$), most models demonstrated stable or improved performance compared to the non-augmented baseline. For the CNN $K$-bit models, accuracies remained above 92% for all shift levels up to $\alpha = 4$, with peak values observed for the Euclidean ($94.17 \pm 6.44\%$) and Mahalanobis ($93.79 \pm 5.59\%$) variants under TA. The EMD and Constrained-EMD models achieved comparable stability, maintaining approximately $92.6 \pm 5$–6% accuracy at $\alpha = 2$–4, before declining to around 87–88% at $\alpha = 8$ when TC was applied. The CNN+$N_{classes}$ exhibited a similar pattern, improving from $91.42 \pm 7.93\%$ (NA) to $92.12 \pm 6.84\%$ under TA at $\alpha = 2$, but dropping sharply to $54.29 \pm 12.15\%$ at $\alpha = 8$ for TC-only conditions. The Siamese networks were notably resilient: the multi-class model increased from $87.29 \pm 8.00\%$ (NA) to $90.22 \pm 5.34\%$ under TA at $\alpha = 2$, while the multiple-classifier variant maintained consistently high performance near $96.2 \pm 2.7\%$ up to $\alpha = 4$, followed by a marked decline to $82.62 \pm 13.38\%$ at $\alpha = 8$ for TC and TA&TC. Overall, augmentation of the training data (TA and TA&TC) consistently improved temporal robustness, whereas applying shifts solely to the test set (TC) reduced accuracy and increased variance across subjects.

## 5 Discussion

We investigated the performance of multiple deep learning and classical models for single-trial decoding of C-VEPs. The approaches varied in their principles and levels, with discriminant approaches that classify input signals directly or reproduce signals for use in a density-based solution as they are compared with templates at a later stage. In deep learning approaches, layers dedicated to high-level feature extraction were common, and we observed substantial differences in how the later stages of the network were used.

The experiments demonstrated that convolutional and Siamese network architectures significantly outperform traditional BLDA and CCA-based classifiers. Among all approaches, the CNN $K$-bit models and Siamese networks achieved superior accuracy, exceeding 90% across subjects, while the classical Corr+BLDA and CCA+BLDA baselines achieved $81.64 \pm 14.02\%$ and $62.59 \pm 16.57\%$, respectively. The improvements achieved by deep learning methods highlight their ability to learn complex temporal–spatial dependencies from raw EEG signals, capturing subtle nonlinearities that conventional linear models cannot represent.
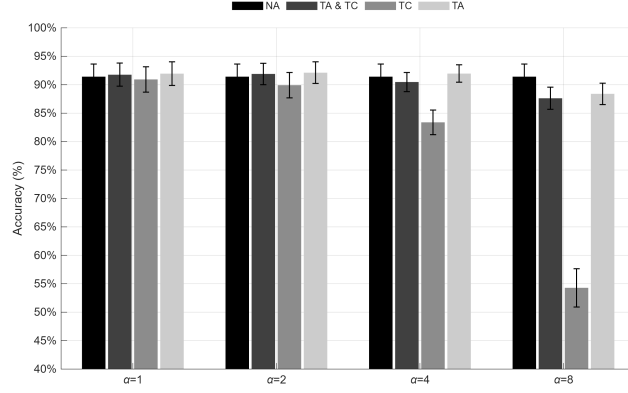
Figure 11: Performance of the CNN $+N_{classes}$ model under various temporal augmentation conditions and shift magnitudes.



(a) Single multi-class Siamese.



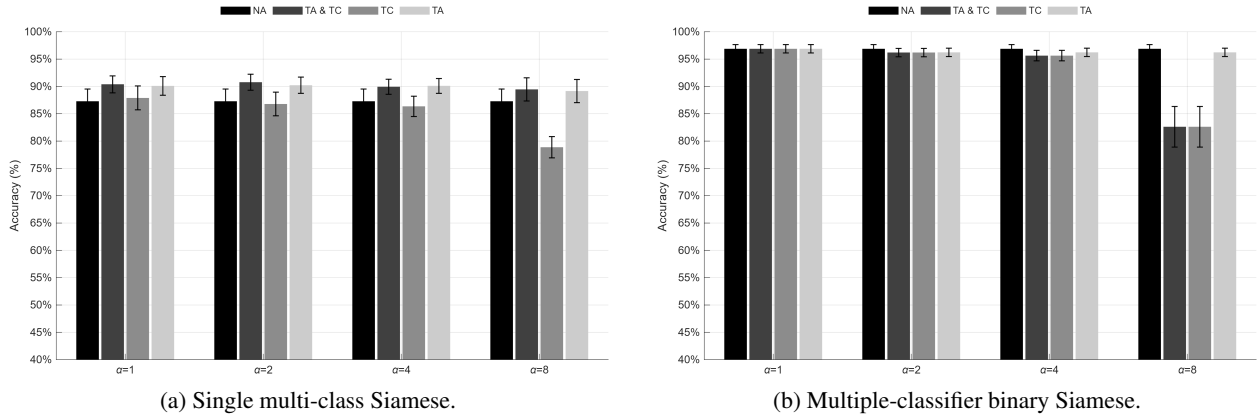(b) Multiple-classifier binary Siamese.

Figure 12: Performance of Siamese network models under different temporal augmentation conditions and shift magnitudes. Comparison between the Single multi-class Siamese and Multiple-classifier binary Siamese architectures.

The comparative analysis among CNN distance metrics revealed that reconstruction-based models decoded the $K$-bit m-sequence most effectively when coupled with perceptually consistent measures such as the EMD and its constrained variant. Both CNN $K$-EMD ($93.88 \pm 6.38\%$) and CNN $K$-C-EMD ($91.30 \pm 7.99\%$) exhibited greater robustness to inter-session variability compared to Euclidean ($93.86 \pm 6.16\%$) and Mahalanobis ($92.38 \pm 8.34\%$) metrics. This finding suggests that EMD-based decoding better preserves temporal structure and accounts for small phase distortions inherent in EEG recordings. Similarly, the CNN $+N_{classes}$ achieved $91.42 \pm 7.93\%$ accuracy, indicating that sequence-level learning can be effectively replaced by categorical learning without substantial performance degradation. The Siamese architectures further demonstrated the capacity of feature-space similarity learning for EEG-based classification, with the multi-class variant ($87.29 \pm 8.00\%$) providing a compact yet effective representation and the multiple-classifier model achieving the highest overall accuracy of $96.89 \pm 2.74\%$.

Temporal data augmentation analysis revealed that moderate temporal shifting enhanced model generalization by simulating variations in physiological response latency. Augmentation up to $\alpha = 4$ samples improved stability for both CNN and Siamese models, particularly under TA and TA&TC configurations, while larger shifts ($\alpha = 8$) caused performance degradation, especially in the TC-only setup. For instance, the CNN $K$-bit (Euclidean) model maintained $94.17 \pm 6.44\%$ under TA at $\alpha = 1$, whereas accuracy dropped to $53.08 \pm 10.81\%$ at $\alpha = 8$ for TC. Similarly, the Siamese multi-class model improved from $87.29 \pm 8.00\%$ (NA) to $90.22 \pm 5.34\%$ (TA, $\alpha = 2$), and the Siamese multiple-classifier variant remained nearly invariant ($\approx 96.2 \pm 2.7\%$) for $\alpha \leq 4$, followed by a sharp decline ($82.62 \pm 13.38\%$) at $\alpha = 8$. These results confirm that augmenting the training data with small temporal perturbations enhances temporal robustness. In contrast, excessive misalignment between training and testing signals introduces inconsistent phase information, which is detrimental to decoding accuracy.

Compared with prior EEG-based C-VEP and SSVEP decoding studies, which typically report 70–85% accuracy using correlation or canonical methods, the proposed deep learning architectures demonstrate state-of-the-art single-trial decoding performance. The use of m-sequence reconstruction and distance-based similarity learning provides a novel

framework that integrates both temporal fidelity and representational flexibility. Moreover, the inclusion of augmentation and inter-session validation establishes the practical relevance of these models for robust, session-independent BCI design.

Although the results are promising, several limitations should be acknowledged. The number of electrodes was limited to eight sites in the occipital and parietal regions, potentially limiting spatial resolution. The models were trained and evaluated within-subject using leave-one-session-out cross-validation, and cross-subject generalization remains to be examined. Future work will aim to extend these findings toward subject-independent training and closed-loop real-time C-VEP BCIs. The integration of multimodal signals, such as electrooculography (EOG), may further enhance decoding reliability and usability in practical neurotechnology systems.

# 6 Conclusion

We proposed and compared multiple deep learning approaches for EEG-based code-modulated visual evoked potential (C-VEP) decoding, including convolutional neural networks (CNNs) with both class-level and sequence-level outputs combined with different distance metrics, as well as Siamese network architectures. The results demonstrated that the proposed deep learning models significantly outperform traditional correlation-based methods, with the CNN $K$-bit models and Siamese networks achieving superior accuracy. Among all models, the multiple-classifier Siamese network achieved the highest overall performance with an accuracy of 96.89%, followed by distance-based CNN variants that attained accuracies in the range of 91–94%. In particular, distance-based decoding using Earth Mover's Distance (EMD) and constrained EMD improved temporal alignment and model generalization, while temporal data augmentation, even with small shifts, enhanced stability against latency variations. These findings highlight the advantages of end-to-end representation learning in capturing the complex temporal–spatial structure of EEG signals, representing a significant step toward practical, reliable, and adaptive non-invasive C-VEP brain–computer interfaces.

## 6.1 Acknowledgment

# References

Xiaorong Gao, Yijun Wang, Xiaogang Chen, and Shangkai Gao. Interface, interaction, and intelligence in generalized brain–computer interfaces. *Trends in cognitive sciences*, 25(8):671–684, 2021.

Ravikiran Mane, Tushar Chouhan, and Cuntai Guan. Bci for stroke rehabilitation: motor and beyond. *Journal of neural engineering*, 17(4):041001, 2020.

J.L. Sirvent Blasco, E. Iáñez, A. Úbeda, and J.M. Azorín. Visual evoked potential-based brain–machine interface applications to assist disabled people. *Expert Systems with Applications*, 39(9):7908–7918, 2012. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2012.01.110. URL `https://www.sciencedirect.com/science/article/pii/S0957417412001285`.

Alain de Cheveigné and Dorothée Arzounian. Robust detrending, rereferencing, outlier detection, and inpainting for multichannel data. *NeuroImage*, 172:903–912, 2018.

Hooman Nezamfar, Seyed Sadegh Mohseni Salehi, and Deniz Erdoğmuş. Stimuli with opponent colors and higher bit rate enable higher accuracy for c-vep bci. *IEEE Signal Processing in Medicine and Biology Symposium*, 2015. doi: 10.1109/SPMB.2015.7405476.

Christoph Guger. State-of-the-art in bci research: Bci award 2010. *InTech eBooks*, 2011. doi: 10.5772/15017.

P Suveetha Dhanaselvam and C Nadia Chellam. A review on preprocessing of EEG signal. In *2023 International Conference on Bio Signals, Images, and Instrumentation (ICBSII)*, pages 1–7. IEEE, 2023.

Martin Spüler and Simone Kurek. Alpha-band lateralization during auditory selective attention for brain–computer interface control. *Brain-Computer Interfaces*, 5(1):23–29, 2018. doi: 10.1080/2326263X.2017.1415496. URL `https://doi.org/10.1080/2326263X.2017.1415496`.

Xinjie He, Brendan Z. Allison, Ke Qin, Liang Wei, Xingyu Wang, Andrzej Cichocki, and Jing Jin. Leveraging transfer superposition theory for stablestate visual evoked potential cross-subject frequency recognition. *IEEE transactions on bio-medical engineering*, 2024. doi: 10.1109/TBME.2024.3406603.

Martin Spüler. A brain-computer interface (bci) system to use arbitrary windows applications by directly controlling mouse and keyboard. *EMBC*, 2015. doi: 10.1109/EMBC.2015.7318554.

Sebastian Nagel, W. Rosenstiel, and Martin Spüler. Random visual evoked potentials (rvep) for brain-computer interface (bci) control. *GBCIC*, 2016. doi: 10.3217/978-3-85125-533-1-64.

Víctor Martínez-Cagigal, Eduardo Santamaría-Vázquez, Sergio Pérez-Velasco, Diego Marcos-Martínez, Selene Moreno-Calderón, and Roberto Hornero. Non-binary m-sequences for more comfortable brain–computer interfaces based on c-veps. *Expert Systems with Applications*, 232:120815, 2023. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2023.120815. URL `https://www.sciencedirect.com/science/article/pii/S0957417423013179`.

Ayas Kiser, Atilla Cantürk, and Ivan Volosyak. Towards secure transaction authentication using a cvep-based bci. *Lecture notes in computer science*, 2025. doi: 10.1007/978-3-032-02728-3_43.

Yuan Miao, Lei Zhang, Jing Jin, Xingyu Wang, and Andrzej Cichocki. High-performance c-vep bci with minimal calibration using cross-subject transfer learning. *Expert Systems with Applications*, 249:123679, 2024.

Víctor Martínez-Cagigal, Jordy Thielen, Eduardo Santamaria-Vazquez, Sergio Pérez-Velasco, Peter Desain, and Roberto Hornero. Brain–computer interfaces based on code-modulated visual evoked potentials (c-vep): a literature review. *Journal of Neural Engineering*, 18(6):061002, 2021.

Vernon J. Lawhern, Amelia J. Solon, Nicholas R. Waytowich, Stephen M. Gordon, Chou Po, and Brent J. Lance. EEGNet: a compact convolutional neural network for EEG-based brain–computer interfaces. *Journal of Neural Engineering*, 15(5):056013, 2018. doi: 10.1088/1741-2552/aace8c.

Robin T. Schirrmeister, Jonathan T. Springenberg, Lukas D. Fiederer, Martin Glasstetter, Katharina Eggensperger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. Deep learning with convolutional neural networks for EEG decoding and visualization. *Human Brain Mapping*, 38(11):5391–5420, 2017. doi: 10.1002/hbm.23730.

J Zhang and K Li. A multi-view CNN encoding for motor imagerycsignals. *Biomedical Signal Processing and Control*, 85, June 2023. ISSN 1746-8094. doi: https://doi.org/10.1016/j.bspc.2023.105063. URL `https://eprints.whiterose.ac.uk/id/eprint/199684/`. © 2023 Elsevier Ltd. This is an author produced version of an article published in Biomedical Signal Processing and Control. Uploaded in accordance with the publisher's self-archiving policy.

Wenlong Wang, Baojiang Li, Haiyan Wang, Xichao Wang, Yuxin Qin, Xingbin Shi, and Shuxin Liu. EEG-FMCNN: A fusion multi-branch 1d convolutional neural network for eeg-based motor imagery classification. *Medical & Biological Engineering & Computing*, 62(1):107–120, 2024. doi: 10.1007/s11517-023-02931-x.

Hubert Cecotti and Akim Graeser. Convolutional neural networks for P300 detection with application to brain–computer interfaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(3):433–445, 2011. doi: 10.1109/TPAMI.2010.125.

S Abinayaa and S.S. Sridhar. An explainable hybrid bio-inspired feature selection framework using rsvp-evoked p300 eeg signals for identity authentication. *Expert Systems with Applications*, 298:129674, 2026. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2025.129674. URL `https://www.sciencedirect.com/science/article/pii/S0957417425032890`.

Xu Yan, Tingting Zhang, Yi Le, Yuhang Shi, Yitian Zhang, Xin Chen, and Yi Mao. Region-focused cnn with dynamic adaptive graph attention network for stereogram evoked eeg recognition. *Expert Systems with Applications*, 287:127859, 2025. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2025.127859. URL `https://www.sciencedirect.com/science/article/pii/S0957417425014812`.

Yu Zhang, Yu Wang, Guoxu Zhou, Jing Jin, Bei Wang, Xingyu Wang, and Andrzej Cichocki. Multi-kernel extreme learning machine for eeg classification in brain-computer interfaces. *Expert Systems with Applications*, 96:302–310, 2018. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2017.12.015. URL `https://www.sciencedirect.com/science/article/pii/S0957417417308291`.

Dheeraj Rathee, Haider Raza, Girijesh Prasad, and Hubert Cecotti. Current source density estimation enhances the performance of motor-imagery-related brain–computer interface. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 25(12):2461–2471, 2017.

Yike Sun, Yuhan Li, Yuzhen Chen, Chen Yang, Jingnan Sun, Liyan Liang, Xiaogang Chen, and Xiaorong Gao. Efficient dual-frequency ssvep brain-computer interface system exploiting interocular visual resource disparities. *Expert Systems with Applications*, 252:124144, 2024. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2024.124144. URL `https://www.sciencedirect.com/science/article/pii/S0957417424010108`.

Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, Cambridge, MA, 2017. ISBN 978-0-262-03561-3.

Kalanit Grill-Spector and Rafael Malach. The human visual cortex. *Annual Review of Neuroscience*, 27(1):649–677, 2004. doi: 10.1146/annurev.neuro.27.070203.144220.

Hubert Cecotti, Miguel A. Leray, Miguel A. Maldonado-Bonilla, Kikumini Singh, and Abuzer Yakaryılmaz. Grape detection with convolutional neural networks. *Expert Systems with Applications*, 159:113588, 2020. doi: 10.1016/j.eswa.2020.113588.

Dan C. Ciresan, Ueli Meier, and Jürgen Schmidhuber. Multi-column deep neural networks for image classification. In *2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3642–3649. IEEE, 2012. doi: 10.1109/CVPR.2012.6248110.

Joseph Redmon and Anelia Angelova. Real-time grasp detection using convolutional neural networks. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1316–1322. IEEE, 2015. doi: 10.1109/ICRA.2015.7139361.

Junshui Ma, Robert P. Sheridan, Andy Liaw, George E. Dahl, and Vladimir Svetnik. Deep neural nets as a method for quantitative structure–activity relationships. *Journal of Chemical Information and Modeling*, 55(2):263–274, 2015. doi: 10.1021/ci500747n.

Sander Dieleman, Kyle W. Willett, and Joni Dambre. Rotation-invariant convolutional neural networks for galaxy morphology prediction. *Monthly Notices of the Royal Astronomical Society*, 450(2):1441–1459, 2015. doi: 10.1093/mnras/stv632.

Erich E Sutter. The brain response interface: communication through visually-induced electrical brain responses. *Journal of Microcomputer Applications*, 15(1):31–45, 1992.

Guangyu Bin, Xiaorong Gao, Yijun Wang, Bo Hong, and Shangkai Gao. Vep-based brain-computer interfaces: time, frequency, and code modulations [research frontier]. *IEEE Computational Intelligence Magazine*, 4(4):22–26, 2009.

Benjamin Wittevrongel, Elia Van Wolputte, and Marc M Van Hulle. Code-modulated visual evoked potentials using fast stimulus presentation and spatiotemporal beamformer decoding. *Scientific reports*, 7(1):15037, 2017.

Qingguo Wei, Siwei Feng, and Zongwu Lu. Stimulus specificity of brain–computer interfaces based on code modulation visual evoked potentials. *PLOS ONE*, 11(5):e0156416, 2016. doi: 10.1371/journal.pone.0156416.

Shuang Liu, Kun Wang, Wei Zhou, Chao Chen, Qiang Liu, and Xiaorong Gao. A multi-target brain–computer interface based on code modulated visual evoked potentials. *PLOS ONE*, 13(8):e0202478, 2018. doi: 10.1371/journal.pone.0202478.

Ivan Volosyak, Aya Rezeika, Mihaly Benda, Felix Gembler, and Piotr Stawicki. Towards solving of the illiteracy phenomenon for vep-based brain-computer interfaces. *Biomedical Physics & Engineering Express*, 6(3):035034, 2020.

Pramod Gaur, Ram Bilas Pachori, Hui Wang, and Girijesh Prasad. A multi-class eeg-based bci classification using multivariate empirical mode decomposition based filtering and riemannian geometry. *Expert Systems with Applications*, 95:201–211, 2018. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2017.11.007. URL https://www.sciencedirect.com/science/article/pii/S0957417417307492.

Ji-Wung Han, Soyeon Bak, Jun-Mo Kim, WooHyeok Choi, Dong-Hee Shin, Young-Han Son, and Tae-Eui Kam. Meta-eeg: Meta-learning-based class-relevant eeg representation learning for zero-calibration brain–computer interfaces. *Expert Systems with Applications*, 238:121986, 2024. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2023.121986. URL https://www.sciencedirect.com/science/article/pii/S0957417423024880.

Claudio Carvalhaes and J Acacio De Barros. The surface laplacian technique in eeg: Theory and methods. *International Journal of Psychophysiology*, 97(3):174–188, 2015.

Dennis J McFarland. The advantages of the surface laplacian in brain–computer interface research. *International Journal of Psychophysiology*, 97(3):271–276, 2015.

Felix Gembler and Ivan Volosyak. A novel dictionary-driven mental spelling application based on code-modulated visual evoked potentials. *Comput.*, 2019. doi: 10.3390/COMPUTERS8020033.

Asghar Zarei, Babak Mohammadzadeh Asl, Asghar Zarei, and Babak Mohammadzadeh Asl. Automatic detection of code-modulated visual evoked potentials using novel covariance estimators and short-time eeg signals. *Computers in Biology and Medicine*, 2022. doi: 10.1016/J.COMPBIOMED.2022.105771.

Milán András Fodor, Hannah Herschel, Atilla Cantürk, Gernot Heisenberg, and Ivan Volosyak. Evaluation of different visual feedback methods for brain—computer interfaces (bci) based on code-modulated visual evoked potentials (cvep). *Brain Science*, 2024. doi: 10.3390/BRAINSCI14080846.

Mohammadreza Behboodi, Amin Mahnam, Hamid Reza Marateb, and Hossein Rabbani. Optimization of visual stimulus sequence in a brain-computer interface based on code modulated visual evoked potentials. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 2020. doi: 10.1109/TNSRE.2020.3044947.

Jing Zhao, Jiaxin Li, Xinrui Wang, Qian Zhang, Zheng Li, and Zhenhu Liang. Discriminating brainwave patterns of different control and non-control states for enhancing asynchronous brain-computer interfaces. *Expert Systems with Applications*, 252:124145, 2024. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2024.124145. URL https://www.sciencedirect.com/science/article/pii/S095741742401011X.

Soroosh Shahtalebi, Amir Asif, and Arash Mohammadi. Siamese neural networks for EEG-based brain–computer interfaces. *arXiv*, abs/2002.00904, 2020.

Rongrong Fu, Yaodong Wang, and Chengcheng Jia. Data augmentation for cross-subject EEG features using siamese neural network. *Biomedical Signal Processing and Control*, 73:103614, 2022. doi: 10.1016/j.bspc.2022.103614.

Huijuan Jiang, Yimin Hou, Yihua Bai, Zehan Zheng, and Xiaojun Wu. A multiscale siamese convolutional neural network with cross-channel fusion for motor imagery decoding. *Journal of Neuroscience Methods*, 363:109426, 2021. doi: 10.1016/j.jneumeth.2021.109426.

Xiaoshan Zhou, Carol C. Menassa, and Vineet R. Kamat. Siamese network with dual attention for eeg-driven social learning: Bridging the human-robot gap in long-tail autonomous driving, 2025. URL https://arxiv.org/abs/2504.10296.

Xiaochen Ye, Chen Yang, Yonghao Chen, Yijun Wang, Xiaorong Gao, and Hongxin Zhang. Multisymbol time division coding for high-frequency steady-state visual evoked potential-based brain-computer interface. *IEEE transactions on neural systems and rehabilitation engineering*, 2022. doi: 10.1109/TNSRE.2022.3183087.

Daiki Aminaka, Shoji Makino, and Tomasz M. Rutkowski. Svm classification study of code-modulated visual evoked potentials. *2015 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA)*, 2015. doi: 10.1109/APSIPA.2015.7415435.

Rafael Grigoryan and Alexander Kaplan. High-speed brain-computer communication interface based on code-modulated visual evoked potentials. *TARGETED ONCOTHERAPY*, 2019. doi: 10.24075/BRSMU.2019.019.

Felix Gembler, Piotr Stawicki, Aya Rezeika, Mihaly Benda, and Ivan Volosyak. Exploring session-to-session transfer for brain-computer interfaces based on code-modulated visual evoked potentials. *2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, 2020. doi: 10.1109/SMC42975.2020.9282826.

Hongze Zhao, Yuanfang Chen, Weihua Pei, Hongda Chen, and Yijun Wang. Towards online applications of eeg biometrics using visual evoked potentials. *Expert Systems with Applications*, 177:114961, 2021. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2021.114961. URL https://www.sciencedirect.com/science/article/pii/S0957417421004024.

Richard W Homan, John Herman, and Phillip Purdy. Cerebral location of international 10–20 system electrode placement. *Electroencephalography and Clinical Neurophysiology*, 66(4):376–382, 1987. ISSN 0013-4694. doi: https://doi.org/10.1016/0013-4694(87)90206-9. URL https://www.sciencedirect.com/science/article/pii/0013469487902069.

Andrew Morley, Lizzie Hill, and A Kaditis. 10-20 system eeg placement. *European Respiratory Society, European Respiratory Society*, 2016.

Richard Inger, Jonathan Bennie, Thomas W Davies, and Kevin J Gaston. Potential biological and ecological effects of flickering artificial light. *PloS one*, 9(5):e98631, 2014.

Mario Kleiner, David Brainard, and Denis Pelli. What's new in psychtoolbox-3? 2007.

R Bullock. Great circle distances and bearings between two locations. *MDT, June*, 5:1–3, 2007.

Emilio Porcu, Moreno Bevilacqua, and Marc G Genton. Spatio-temporal covariance and cross-covariance functions of the great circle distance on a sphere. *Journal of the American Statistical Association*, 111(514):888–898, 2016.

Carl Carter. Great circle distances. *SiRF White Paper*, 2002.

Anna Bykhovskaya and Vadim Gorin. Canonical correlation analysis: review, 2025. URL https://arxiv.org/abs/2411.15625.

Su-Na Zhao, Yingxue Cui, Guangxin Zhao, Liying Jiang, and Xin-Zhong Chang. A new cca-based method for improving ssvep-based bci system classification. In *2023 WRC Symposium on Advanced Robotics and Automation (WRC SARA)*, pages 432–437, 2023. doi: 10.1109/WRCSARA60131.2023.10261825.

Guangyu Bin, Xiaorong Gao, Zheng Yan, Bo Hong, and Shangkai Gao. An online multi-channel ssvep-based brain-computer interface using a canonical correlation analysis method. *Journal of neural engineering*, 6:046002, 07 2009. doi: 10.1088/1741-2560/6/4/046002.

Patrick Schober, Christa Boer, and Lothar A Schwarte. Correlation coefficients: appropriate use and interpretation. *Anesthesia & analgesia*, 126(5):1763–1768, 2018.

Ibrahim Gokcen and Jing Peng. Comparing linear discriminant analysis and support vector machines. *Lecture Notes in Computer Science*, 2002.

Zhidong Shen, Yuhao Zhang, and Weiying Chen. A bayesian classification intrusion detection method based on the fusion of pca and lda. *Secur. Commun. Networks*, 2019.

D. Huang, C. Xiang, and S. S. Ge. Feature extraction for face recognition using recursive bayesian linear discriminant. *2007 5th International Symposium on Image and Signal Processing and Analysis*, 2007.