

Improving Intention Detection in Single-Trial Classification through Fusion of EEG and Eye-tracker Data

Xianliang Ge*, Yunxian Pan, Sujie Wang, Linze Qian, Jingjia Yuan, Jie Xu, Nitish Thakor, *Fellow, IEEE*, and Yu Sun*, *Senior Member, IEEE*

Abstract—Intention decoding is an indispensable procedure in hands-free human-computer interaction (HCI). Conventional eye-tracking system using single-model fixation duration possibly issues commands ignoring users' real expectation. In the current study, an eye-brain hybrid brain-computer interface (BCI) interaction system was introduced for intention detection through fusion of multi-modal eye-track and ERP (a measurement derived from EEG) features. Eye-track and EEG data were recorded from 64 healthy participants as they performed a 40-min customized free search task of a fixed target icon among 25 icons. The corresponding fixation duration of eye-tracking and ERP were extracted. Five previously-validated LDA-based classifiers (including RLDA, SWLDA, BLDA, SKLDA, and STDA) and the widely-used CNN method were adopted to verify the efficacy of feature fusion from both offline and pseudo-online analysis, and optimal approach was evaluated through modulating the training set and system response duration. Our study demonstrated that the input of multi-modal eye-track and ERP features achieved superior performance of intention detection in the single trial classification of active search task. And compared with single-model ERP feature, this new strategy also induced congruent accuracy across different classifiers. Moreover, in comparison with other classification methods, we found that the SKLDA exhibited the superior performance when fusing feature in offline test (ACC=0.8783, AUC=0.9004) and online simulation with different sample amount and duration length. In sum, the current study revealed a novel and effective approach for intention classification using eye-brain hybrid BCI, and further supported the real-life application of hands-free HCI in a more precise and stable manner.

Index Terms—Single trial classification, eye-tracking, eye-brain-computer interface, event related potential, EEG

I. INTRODUCTION

This work was supported by the National Natural Science Foundation of China (81801785, 31800931), the Hundred Talents Program of Zhejiang University, by the Zhejiang University Global Partnership Fund (100000-11320), by the Zhejiang Lab (2019KE0AD01), and by the Space Medical Experiment Project of China Manned Space Program (HYZHXM03001). (*, *Corresponding author: X. Ge and Y. Sun*)

X. Ge, Y. Pan and J. Xu are with the Center for Psychological Sciences, Zhejiang University, 310027, Hangzhou, Zhejiang, China, (X. Ge, email: 0918082@zju.edu.cn).

S. Wang, L. Qian, J. Yuan are with the Key Laboratory for Biomedical Engineering of Ministry of Education of China, Department of Biomedical Engineering, Zhejiang University, 310027, Hangzhou, Zhejiang, China.

N. Thakor is with the Department of Biomedical Engineering, Johns Hopkins University School of Medicine, United States, and also with the Department of Biomedical Engineering, National University of Singapore, 117456, Singapore.

Y. Sun is with the Key Laboratory for Biomedical Engineering of Ministry of Education of China, Department of Biomedical Engineering, Zhejiang University, and also with the Zhejiang Lab, 310027, Hangzhou, China, (email: yusun@zju.edu.cn).

HANDS-free human-computer interaction (HCI) draws growing attention considering its convenience in various environments. As an efficient input modality in HCI, eye-gaze system was proposed in 1980s [1], [2]. By the use of eye-tracker, selection is typically achieved through gazing on the target item for a specific period. Traditional eye-tracking system is straightforward to implement, ensuring short calibration time [3] and faster target acquisition [4]. However, unexpected command was possibly issued when using the standalone eye-gaze system, leading to the Midas-touch problem [5], [6]. System cannot distinguish the spontaneous fixation from intended selection when dwelling time indeed exceeds the threshold. Besides, it is also inapplicable to adapt the threshold of dwell time to various ecological scenarios, so that the processing time of the user to the stimulus may overrun occasionally. Thus, in order to overcome this problem, other forms of input signals were introduced to assist the eye-gaze HCI system in intention decoding, such as special eye saccade and other physical movements [7]. But such additional motor activities will induce extra mental workload and distract the execution of main task, as well as being inconvenient to the disabled. Therefore, in order to enhance the fluency and robustness of the system, a natural and intuitive input and decoding approach is necessary.

Brain-computer interface (BCI) establishes a direct communication channel from brain signals [8], which can be used to detect the ongoing cognition, such as mental fatigue [9], emotion state [10] as well as intention [11]–[13]. This technology has been proposed to decode the cognitive information implicitly from users' mind, without additional interruption to the primary task [14]. Among the signal acquisition techniques of BCI, electroencephalography (EEG) has been widely adopted due to its numerous advantages in temporal resolution, usability and cost. As one of the most popular characteristics in EEG analysis, event-related potential (ERP) is generated by the neuron sensitive to the specified stimulus or events and broadly used to capture cognitive or sensory process [15]. Generally, ERP includes several time-locked components (i.e. N170, a negative waveform at around 170 ms post-stimulus; P300, a positive waveform at around 300 ms post-stimulus, etc.) that correspond to particular cognition states [16], which can serve as natural biomarkers of the user's intention of interaction. This special characteristic of ERP has been evaluated for device control [17], target and error detection [18], [19], and so on. As a potential substitute for

the “click” operation in HCI system, accumulating evidences confirmed that, when an intention of item selection emerged, negative potentials could be detected during conscious dwell time in the central [17] and parietal [20], [21] electrodes. One of the most popular paradigm in the EEG-based HCI research is P300 speller, which utilizes the uncommon event (flash of the target character) to induce P300 wave and decode user’s intention [19]. As for the free visual search task, Kaunitz et al found that when subjects detected the target among distractors, a robust sensory component of fixation event-related potentials emerged, and a single-trial analysis could differentiate the type of stimulus based on EEG signals [22]. Devillez et al also observed a P300 component for fixation of the target natural scene compared with the free viewing without any target [23].

In sum, the constituents of ERP contain abundant information towards personal intention and accompany gaze-based control intuitively in free search task, which can serve as the feature for distinguishing the intended selection from involuntary fixation. As a result, a combination of ERP and the eye-gaze input system could be complementary and provides more robust interaction experience. Accumulating evidence indicate that this hybrid BCI can satisfy the need for speed and accuracy simultaneously, overcoming the Midas-Touch problem of eye-tracker and inter-person variability of BCI protocol. For example, Kalika et al fused the eye-gaze data into a P300 speller pipeline and reported an improved classification accuracy and declined flash number for character spelling [24]. Choi et al utilized the gaze position to contract a 12×12 character matrix into a 3×3 one, and highlighted this smaller navigation area to enhance the decoding performance of P300 speller [25]. However, to the best of our knowledge, most hybrid BCI researches about the free visual search task only focus on the decoding of EEG signal, or take the EEG and eye-tracking for separate control purposes. Few of them attempted to analyze these two modalities in parallel and fuse them as input for intention classification.

This research gap inspired our study to find out whether the alliance of input data from eye-tracker and BCI could facilitate the performance of intention detection in single-trial classification for active search task. Specifically, a self-designed HCI paradigm was proposed in which participants were required to search and identify a target icon among a total of 25 icons for each trial. We analyzed the fixation duration of each stimuli and corresponding ERP components, and these two inputs served as features for Target/Non-target classification. In the previous studies, LDA was extensively used for ERP detection owing to the satisfactory performance and simplicity [26], [27], while it also accompanied with the disadvantages of high noise sensitivity, poor inter-person generalization and the need of large training sample [28]. Therefore, multiple improved LDA classifiers (including RLDA, SWLDA, BLDA, SKLDA, and STDA) with divergent edges were adopted for evaluating their performance in this task. And CNN, the most prevalent deep learning framework in the study of BCI [29], was also taken for comparison. Our study demonstrated that the fusion of concurrent ERP and fixation duration induced a superior performance over single feature in target intention decoding among all classifiers. Pseudo-online validation was

further conducted to explore the proper amount of training set, response time of the system and optimal classification approach, so as to provide additional support for the practical application in various real-life HCI scenarios.

II. METHODS AND MATERIALS

A. Subjects

The study sample consisted of 70 university students (male / female = 35 / 35) from the Zhejiang University, China. All participants were aged between 17 and 29 years (mean age = 22.4 ± 2.3 years) and reported normal or corrected-to-normal vision. Participants with chronic illness, sleep disorder, childhood history of ADHD, and long-term medication history were excluded during pre-screened telephone interview. Prior to the experiment date, the included subjects were required to obtain a full night of sleep (> 7 hours) for continuous 2 nights to minimise the effect of prior sleep restriction on neurobehavioral functions. Subjects consuming caffeine or alcohol, or undertaking strenuous exercise for 24 hour preceding the study were rescheduled. The study was approved by the Institutional Review Board of Zhejiang University and all participants signed informed consent prior to participation.

B. Experimental Protocol

The experimental protocol was a typical target identification task that was customized using C programming language (Fig. 1). Specifically, participants were requested to search for a target icon among multiple non-target icons as quick and accurate as possible. A total of 25 icons were presented on a screen (1920×1080 pixels) with the background color was set at [R, G, B] = [192, 192, 192], which were arranged in 5×5 . The experimental interface is shown in (Fig. 2). The size of each icon was set at 24×24 pixels, corresponding to a field of view (FOV) of $0.67^\circ \times 0.67^\circ$. The horizontal/vertical distance between each pair of adjacent icons was set at 100 pixels. If an icon was highlighted, its size would be enlarged to 1.5 times of the original size (i.e., 36×36 pixels, $\text{FOV} = 1^\circ \times 1^\circ$). During the target searching, a certain icon which the participants gazed at and the surrounding eight icons would be highlighted (Fig. 2).

In the experiment, a predefined target icon (60×60 pixels, $\text{FOV} = 1.67^\circ \times 1.67^\circ$) was initially presented for 3 s. A black fixation was presented for 1 s, indicating the start of each experimental trial. Then, after a 0.8 s of a blank screen, the search interface was presented with a randomly assigned 5×5 icon pattern and a timer was started. The cursor was hidden at the moment. Participants were required to search the target icon within 5 s in the search interface. The mouse cursor appeared on the screen after finishing the 5 s searching period. Meanwhile, all the icons on the display were masked with a dotted line. The participant was requested to move the cursor to the location of the masked target icon and click to confirm the selection. The program would proceed to the next trial starting with a black fixation. A short period of break (i.e., 5 s) was introduced after completing 15 trials while a long period of break (i.e., 10 s) was introduced after completing 30 trials. For each participant, a total of 240 trials were administrated and the entire experiment lasted approximately 40 min.

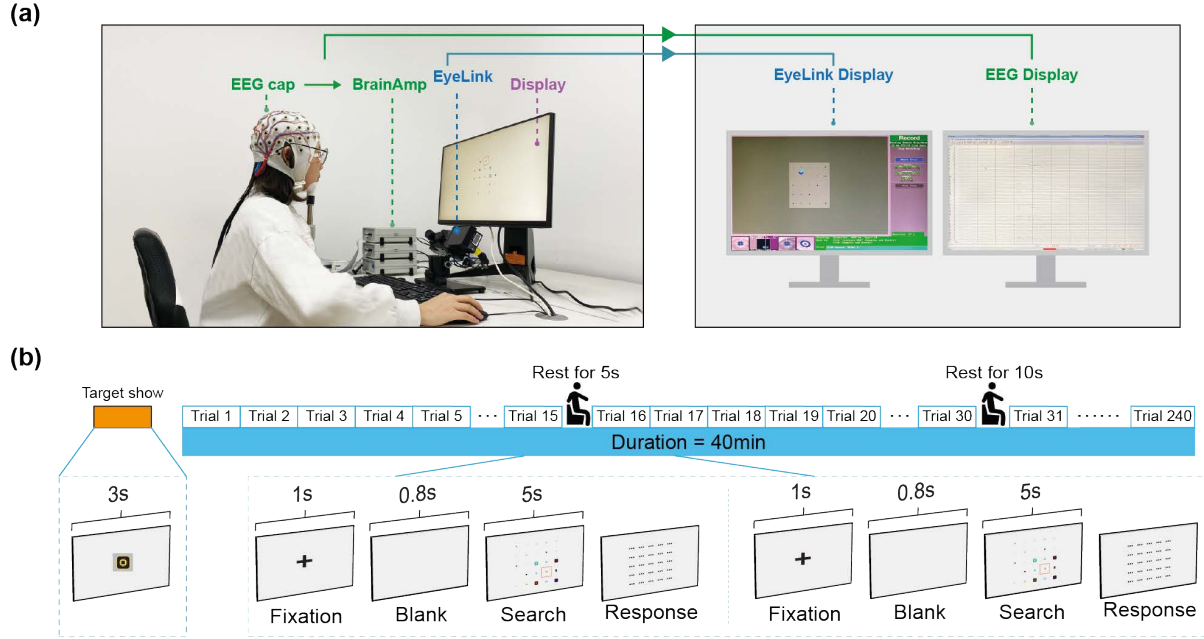


Fig. 1. A schematic diagram of the experimental protocol. (a) The setup of the experiment. EEG data was obtained from 64-channel BP system and eye-tracking data was collected by using EyeLink 1000 Plus system. (b) Each participant performed a 40-min target identification experiment, where the participant was asked to search the target icon as quick and accurate as possible and memorize and identify its location in the response period.

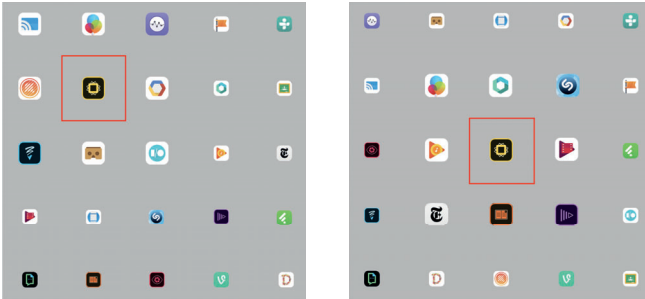


Fig. 2. Sample screenshots of visual target identification interface for two random trials. The icons surrounding the target icon as highlighted with a red box were enlarged to facilitate the identification.

C. Data Acquisition and Preprocessing

The EyeLink 1000 Plus eye-tracking system (model: SR Research, Ottawa, Canada) was used to record the eye-tracking data. The sample rate was set at 1000 Hz. The participants were seated 60 cm from the monitor with a FOV of $13.86^\circ \times 13.86^\circ$ using a chin support. Prior to the experiment, the eye-tracking system was calibrated for each participant. EEG data was recorded from a 64-channel EEG system (model: BrainAmp DC, Brain Products GmbH, Gilching, Germany) according to the international 10-20 system. In addition, horizontal and vertical electrooculograms (EOG) were recorded on the lateral to the outer canthi (HEOG) as well as above and below (VEOG) the right eye. Electrode impedance was kept below $10 \text{ k}\Omega$ throughout the whole experiment. Anti-aliasing was achieved with a band-pass filter (0.5 – 100 Hz) and additionally a 50 Hz notch filter was applied to avoid main interferences. Raw EEG and EOG signals were digitized at a

sampling rate of 500 Hz using FCz as the reference. Two subjects were excluded due to data recording issues.

In the analysis of eye-tracking data, we used the 100×100 px area centered on each icon as an area of interest (AOI) of the icon. When a fixation fell in an AOI including non-target icon, this fixation was classified as the fixation of the non-target icon; while a fixation fell in an AOI including target icon, this fixation was classified as the fixation of the target icon. The fixation duration was calculated by summing up the duration of all fixations in the AOI.

A standard EEG preprocessing pipeline was adopted here, which included FIR band-pass filtering (1 – 40 Hz), re-referencing to the average of all electrodes and ocular artifacts removal by removing the most correlated components to the EOG signals through independent component analysis (ICA) [30]. All preprocessing steps were performed using customized codes and the EEGLAB toolbox [31] in Matlab 2017b (The MathWorks Inc, US). Greater details of the preprocessing steps could be found in our previous studies [10].

The channel selection were based upon results in [26], which has been shown the optimal balance between classification performance and least number of electrodes. Here, Fz, Cz, Pz, Oz, P3, P4, PO7 and PO8 electrodes were included for the following classification. To extract features for target and non-target responses, the continuous EEG signals were segmented to 500 ms epochs with baseline correction by 100 ms interval before the icon presentation. Afterward each epoch was down-sampled to 32 Hz, that is 16 points for each channel and 128 points in total for all 8 channels.

D. Classification Algorithms

Several widely-used algorithms that were popular in the studies of ERP-BCI were adopted in the current work to assess the classification performance [32], [33], including regularized linear discriminant analysis (RLDA) [34], Stepwise linear discriminant analysis (SWLDA) [35], Bayesian linear discriminant analysis (BLDA) [36], Shrinkage linear discriminant analysis (SKLDA) [37], Spatial-Temporal discriminant analysis (STDA) [38], and Convolutional neural network (CNN) [39]. These algorithms were selected to cover the common categories of method for ERP-BCI, that is, concatenation of temporal points and spatial channels (RLDA, SWLDA, BLDA, SKLDA), adoption of spatial-temporal samples (STDA), and deep learning approach (CNN).

1) *Regularized LDA*: RLDA, a regularized version of LDA, is a popular technique for dimensionality reduction and feature extraction. It was originally introduced to solve the small sample size problem. The performance of RLDA technique depends upon the choice of the regularization parameter. In the current work, the regularization parameter was estimated using a deterministic approach according to [40]. This approach avoids the use of the heuristic cross-validation procedure for parameter estimation and improves the computational efficiency. Here the amount of regularization was set as $\lambda = 0.01$.

2) *Stepwise LDA*: SWLDA, another regularized version of LDA, has been shown to be superior in the case of small sample size due to its implementation of combined forward and backward stepwise analysis to select suitable features in the discriminant model. Briefly, model estimation for SWLDA is conducted in a greedy manner by iteratively inserting and removing features from the model based upon statistical tests until the maximal number of active variable is reached or no additional features satisfy the entry/removal criteria. Here, the criteria was set as $p_{ins} = 0.1$ and $p_{rem} = 0.15$ as recommended in [41].

3) *Bayesian LDA*: BLDA is a probabilistic method that based upon Bayesian regression and has been shown to outperform the original LDA method when only a small number of training sets was obtained or strong noise contamination in the data [42]. According to [36], the neurophysiological and experimental priors are employed explicitly by modeling the trial-level covariance and the weight vector covariance of LDA explicitly as linearly separable components with the relative contribution of each component is controlled by the hyperparameters that could be estimated via Restricted Maximum Likelihood.

4) *Shrinkage LDA*: Through adjusting the extreme eigenvalues of the covariance matrix towards the average eigenvalue, SKLDA improves the traditional LAD when using insufficient training samples. For high-dimensional data with only a few data points given (i.e., EEG data), the estimation for a covariance matrix may become imprecise, which may lead to a systematic error: large eigenvalues of the original covariance matrix are estimated too large, and small eigenvalues are estimated too small. Of note, shrinkage is a common remedy for compensating the systematic bias of the estimated covariance matrices and shrinkage parameter for high-dimensional feature spaces. In the current work, the shrinkage parameter was

set at 0.1 according to [43]. For details of SKLDA and its interpretation could be found in [33]

5) *Spatial-Temporal Discriminant Analysis*: STDA is a multiway extension of the LDA that tries to maximize the discriminant information between target and nontarget classes through finding two projection matrices from spatial and temporal dimensions collaboratively. Unlike the abovementioned different versions of LDA method where data were concatenated as input, through incorporating the spatial and temporal information, STDA reduces the feature dimensionality in the discriminant analysis and decreases the number of required training samples [38].

6) *Convolutional Neural Network*: CNN was initially used in computer vision and has gained substantial interest in BCI most recently for its superior performance. In this research, a five-layer CNN was developed for EEG pattern detection. The input of the network was a 2D space-time EEG signal with a size of 8×16 . It was followed by two paired layers, with each pair comprised a convolutional layer with batch normalization and a max-pooling layers. In the first convolutional layer, we utilized 32 kernels with a size of 1×5 for time domain convolution. While the second was used for spatial domain convolution, containing 32 kernels with a size of 8×1 , which equaled the number of EEG electrodes. After each convolution process, a ReLU function was employed for non-linearization. For max-pooling layers, they both utilized a pooling filter size of 1×2 to reduce computational complexity. After dropout process with a dropout rate of 0.3, the output of max-pooling layer was applied to two fully connected layer comprising 64 and 2 neurons respectively. In the decision step, the classification probability is determined by softmax function.

E. Offline Classification

In order to demonstrate that fusion of EEG and eye-tracker data would lead to superior performance in comparison with single EEG or eye-tracker modality, classification was initially performed using only EEG or eye-tracker data. Specifically, fixation duration corresponding to extracted epoch of one gaze was initially estimated and selected as input for LDA classifier as a benchmark. For EEG data, a 0 – 500 ms epoch after a gaze was cut out and selected as input for the offline classification. Of note, one trial might have multiple target and non-target samples (corresponding to the search for target) with the number of non-target samples larger than that of target samples. A cross-validation approach was initially employed to assess the performance of classifiers under different number of training samples using a reformatted balance data. Specifically, the training set (Target:Non-target = 1:1) was designed using sample number from 30 to 420 with a step of 30, while the testing set was randomly selected from the remaining samples and maintained a Target:Non-target = 1:1 fashion with a maximum amount. Of note, the same training and testing samples were applied on all classification algorithms to allow for fairness comparison. This procedure was repeated for 10 times, and the average area under curve (AUC) of the receiver operating characteristic (ROC) curves was computed

for the quantitative comparisons. Then, a separate 10-fold cross-validation approach was applied on the real data (on average, Target:Non-target \approx 1:2.3) to demonstrated that the fusion of multi-modal features induced a superior performance over single feature.

F. Pseudo-online Validation

1) *Online Classification*: A pseudo-online analysis was performed to validate the feasibility and practicability of implementing the decoding algorithm based upon our analysis framework. As at least 240 trials were performed for each subject in the experiment, the samples within the first 80 trials were utilized as training set whereas the remaining data were considered as testing set for assessing the performance of online classification. To avoid the imbalance of the sample amount between two classes in the training set and for the convenience of result analysis in the testing set, the number of the target and non-target classes was set to equal respectively. Of note, fixation duration longer than 500 ms would be redefined as 500 ms to ensure same duration of EEG data.

2) *Epoch threshold*: As it has been mentioned previously, gaze fixation duration was defined as the duration post a gaze for either target or non-target and EEG data between 0 – 500 ms was used as threshold for data extraction and the following classification. In order to assess the influence of different threshold on the classification, we have also used 300 ms to 800 ms with a step of 50 ms as threshold for the epoch extraction. For instance, for a predefined threshold (e.g., 400 ms), gaze fixation duration above the threshold would be redefined as the threshold value and the EEG data between 0 – 400 ms would be used as input. Besides, the number of training trials was also put into consideration as a factor contributing to online performance. In detail, the samples within the first 20 to 100 trials with a step of 10 trials were regarded as training set, while the remaining were used for testing. The ratio of Target and Non-target was rearranged to 1:1 in both training and testing set as well.

III. RESULTS

A. Behavioral Performance

Data from four participants were excluded for signs of poor motivation on the task, likely due to boredom experienced during the target identification experiment. Threshold for signs of poor motivation on the task was set if the error rate of the participant was 1 S.D. lower than the group average. Our final dataset thus consisted 64 participants (male / female = 29 / 35) and the following classification was conducted on these participants. Overall, the remaining participants performed the experiment well, as indicated by the relatively high detection rate (mean \pm S.D. = 98.12% \pm 1.33%). We had performed additional statistical analysis to assess the gender effect and found insignificant difference between males and females ($t_{62} = 0.084$, $p = 0.933$).

B. Characteristics of ERPs

The characteristics of ERPs were first analyzed and compared between target and non-target stimulus. Fig. 3 shows the

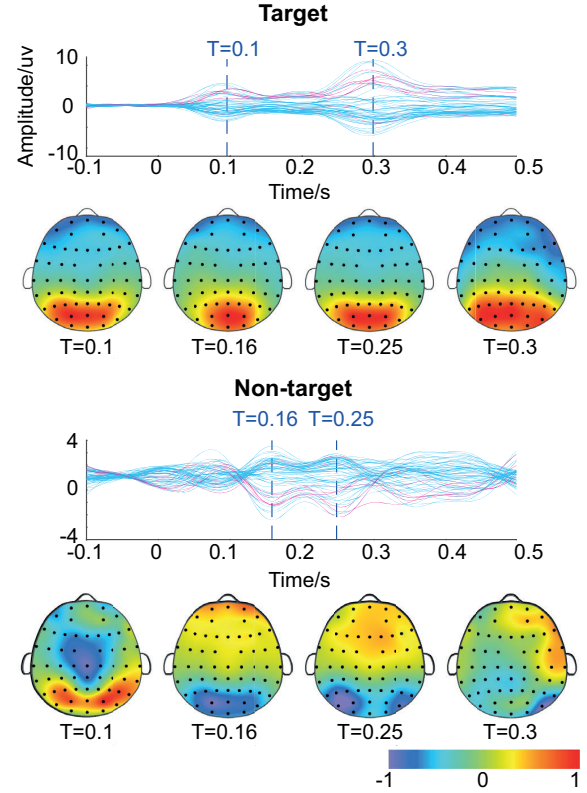


Fig. 3. Distributions of different ERP characteristics of Target and Non-target from a randomly-selected subject.

temporal and spatial differences between two kinds of ERPs for a randomly-selected subject. Specifically, the discriminant ERP features between target and non-target were restricted to the occipital areas post-stimulus. Hence, these evident differences between target and non-target serve as salient underlying features for the following classification algorithms. Moreover, the observed posterior differences were in line with the findings in [26] and justify the selection of the EEG channels (i.e., Fz, Cz, Pz, Oz, P3, P4, PO7 and PO8 in this work) for the classification algorithms.

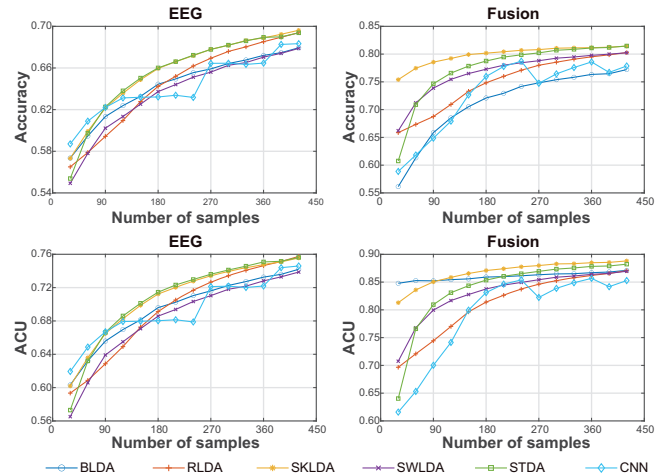


Fig. 4. Comparison of Accuracy and AUC across the employed methods under different number of training samples.

TABLE I
OFFLINE CLASSIFICATION PERFORMANCE ACROSS DIFFERENT ALGORITHMS

Algorithms	EEG		Fusion	
	Accuracy	AUC	Accuracy	AUC
RLDA	0.7793 ± 0.0403	0.8007 ± 0.0599	0.8802 ± 0.0551	0.9104 ± 0.0547
SWLDA	0.7754 ± 0.0382	0.7908 ± 0.0617	0.8761 ± 0.0534	0.9072 ± 0.0565
BLDA	0.6716 ± 0.0504	0.7933 ± 0.0632	0.7954 ± 0.0910	0.9066 ± 0.0573
SKLDA	0.7225 ± 0.0503	0.7838 ± 0.0615	0.8783 ± 0.0603	0.9004 ± 0.0608
STDA	0.7652 ± 0.0367	0.7789 ± 0.0583	0.8740 ± 0.0544	0.9049 ± 0.0568
CNN	0.7854 ± 0.0403	0.8053 ± 0.0622	0.8772 ± 0.0535	0.9028 ± 0.0512

Note: Values are presented as mean \pm S.D., Fusion indicates features from EEG and eye-tracker were fused to obtain the Accuracy and AUC.

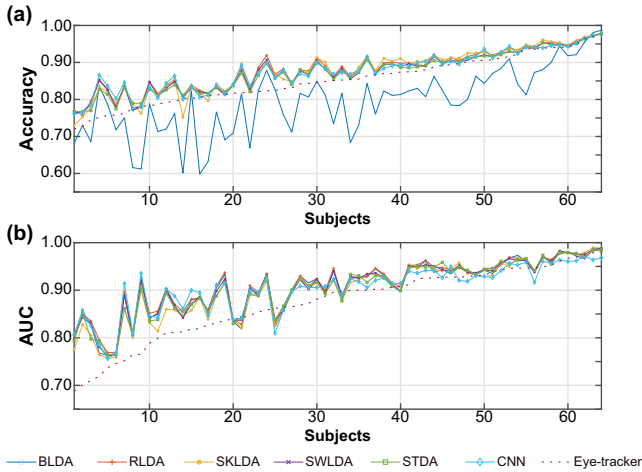


Fig. 5. (a) Accuracy and (b) AUC for multiple classifiers of each subject. The subject was sorted by the ascending of classification performance based upon the eye-tracker model.

C. Offline Classification

In the offline classification, we first assessed the performance of classifiers under different number of training samples. In line with previous study [32], we found that the classification performance was monotonically increased with the number of training samples (Fig. 4). We then assessed the classification performance when using features from eye-track and EEG data respectively. When using eye-track feature, we obtained the classification accuracy of 0.8527 ± 0.0638 and the AUC of 0.8734 ± 0.0746 that was served as benchmark. However, the classification performance across different algorithms using only EEG data is significantly lower than the benchmark probably due to the large inter-individual differences in single-trial EEG characteristics (Table I). Moreover, we found that through employing the features from both eye-track and EEG data, the classification performance was significantly improved and exhibited a superior performance compared to the benchmark for most of the subjects (Fig. 5). Further interrogation of the classification performance across six methods, we found BLDA exhibited relatively low accuracy in the fusion manner. Hence, the remaining five methods (i.e., RLDA, SWLDA, SKLDA, STDA, and CNN) were selected for the following pseudo-online validation.

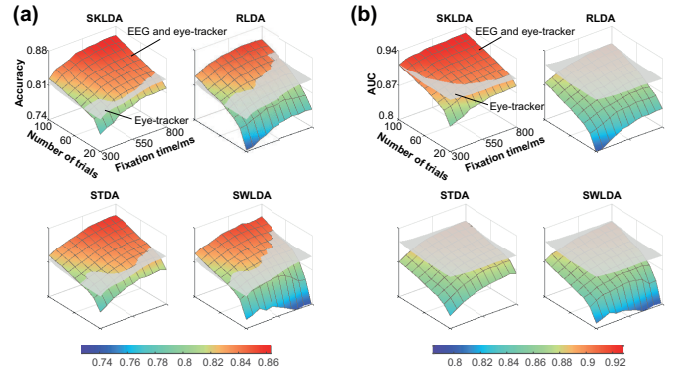


Fig. 6. Distribution of (a) Accuracy and (b) AUC across employed methods that exhibited different performance under number of training trials and durations of fixation. The gray surface indicated the performance when only using eye-tracking data that consider as the benchmark.

D. Pseudo-Online Classification

The performance of the pseudo-online classification was shown in (Table II). Again, we first obtained the performance benchmark from eye-track data (accuracy = 0.8277 ± 0.0862 and AUC = 0.9035 ± 0.0760). Similar to the offline results, the classification performance with single EEG feature was significantly lower than the benchmark for all methods. Although the performance was greatly improved by adding eye-tracker feature into classification, only SKLDA obtained both higher accuracy (0.8454 ± 0.0868) and higher AUC (0.9171 ± 0.0614) in a single-trial classification time of 2.3343 ± 0.0306 ms. Of note, the performance of CNN was significantly lower than the other four methods with a longer classification time. Following validation was therefore only performed on the remaining four methods.

To investigate how the number of training trials and the duration of fixation influence the pseudo-online classification performance, the distribution of accuracy and AUC with different settings for these four selected classifiers was shown in (Fig. 6). The rendered surface represented the classification performance obtained with fusion features, while the gray surface indicated the performance when only using eye-tracking data that consider as the benchmark. Similar to the offline results, the performance of all methods improved monotonically with increasing number of training trials. In

TABLE II
CLASSIFICATION PERFORMANCE ACROSS DIFFERENT ALGORITHMS DURING ONLINE VALIDATION

Algorithms	EEG			Fusion		
	Accuracy	AUC	Time (ms)	Accuracy	AUC	Time (ms)
RLDA	0.6550 ± 0.0560	0.7125 ± 0.0692	3.3722 ± 0.0544	0.8419 ± 0.0879	0.8890 ± 0.0778	3.2570 ± 0.0464
SWLDA	0.6466 ± 0.0535	0.7061 ± 0.0669	2.3832 ± 0.0243	0.8388 ± 0.0910	0.8858 ± 0.0820	2.3198 ± 0.0326
SKLDA	0.6595 ± 0.0558	0.7624 ± 0.0513	2.3773 ± 0.0348	0.8454 ± 0.0868	0.9171 ± 0.0614	2.3343 ± 0.0306
STDA	0.6543 ± 0.0568	0.7123 ± 0.0674	3.2896 ± 0.0687	0.8410 ± 0.0855	0.8914 ± 0.0791	3.2217 ± 0.0287
CNN	0.6592 ± 0.0557	0.7215 ± 0.0661	5.2227 ± 0.1007	0.7293 ± 0.0941	0.8426 ± 0.0790	6.2242 ± 0.0983

Note: Time indicates the duration for single-trial classification.

contrast, the classification performance exhibited a complex dependent level for the setting of fixation duration, i.e., the best performance was not always obtained using long fixation duration. Among the four methods, only SKLDA exhibited superior classification performance in most of the settings compared to the benchmark.

IV. DISCUSSION

In this study, we revealed an explicitly improved performance with the fusion of ERP and eye-tracking data in the single trial classification of free search task, both in offline and online analysis. Our previous research proved the effectiveness of block highlight display (BHD) eye-controlled technique [44], which was further validated by the high target-detection rate in the performance of active search paradigm in this study. The ERP components were also demonstrated in different spatial (occipital area) and temporal (100 and 300 ms) characteristics between target and non-target detection. Six widely-used classifiers were employed to verify the effectiveness of the hybrid BCI system. In the offline analysis, the classification approach with multi-modal inputs of fixation duration and ERP significantly outperformed the method with single-modal brain/eye features. Besides, the multiple versions of LDA approach, except for BLDA, were proved effective in the single-trial classification, showing a superiority in solving such problem with relatively simple inputs than sophisticated neural network structure. As for the online validation, we found the fused feature still provided a robust performance, while different classification approaches exhibited divergent adaptability towards limited number of samples and response time. And SKLDA ranked the top among all LDA classifiers from both the accuracy and application efficiency.

In the offline analysis, we observed that the fused feature provided a much higher accuracy and AUC in the 10-fold classification. When ERP was taken as the only feature, the classification performance of different classifiers was scattered. Although ERP signal could reflect cognitive processes with a high temporal resolution, its low amplitude and sensitivity to various artifacts made it hard to be extracted stably and also varied across different subjects. This was consistent with previous findings that the prediction performance using ERP was mixed in different studies. Specifically, SKLDA was proposed for single-trial classification of ERP-based BCI by Blankertz et al. [33], suggesting superior performance over ordinary LDA and SWLDA. Zhang developed STDA and

demonstrated its superiority among several forms of LDA methods (LDA, SWLDA and SKLDA) in ERP classification [38]. In our study, we found that among the LDA-based classifiers, RLDA outperformed other methods with an accuracy of 0.7793 and BLDA ranked the last with an accuracy of 0.6716. As the most popular deep learning method in ERP-related studies [29], CNN also has an exceptional performance with the top accuracy of 0.7854. Briefly, the performance of classifier varied across different studies and datasets for ERP-based classification.

However, when ERP and eye-tracker data were collectively adopted for classification input in the present study, the performance of divergent classifiers was greatly improved than only taking single-model ERP or eye-track feature. Besides, it was interesting to find that the accuracy and AUC of different classification algorithms became similar. In our experiment protocol, fixation duration tended to be higher when participants gazed on the target icon, consonant with the practical scenarios. The eye-tracker ensures a high temporal and spatial accuracy towards the gazing time and position, indicating a robust measure of underlying cognitive processes based on eye movement-related variables, such as fixation duration and saccade [45]. Therefore, compared with the single ERP signal, the fuse of fixation duration from eye-tracker provided a relatively stable criterion towards cognitive states without large inter-subject variability, so that the input formulation was highly adaptable and the prediction performance tended to be similar across different algorithms. On the other hand, when compared with the single fixation duration, the accuracy was also improved by feature fusing, which demonstrated that ERP was associated with ongoing intention of target selection to enhance the recognition performance of traditional eye-tracking system. Among different classification approaches, the performance of RLDA, SKLDA, SWLDA and STDA were close with the accuracy of 0.8802, 0.8783, 0.8761 and 0.8740 respectively. In addition, in the 10-fold offline analysis, CNN ranked the second when fusing eye and brain features, which was possibly due to the abundant training sample and unlimited processing time. When training set was massively cut down, most LDA-based classifiers outperformed CNN with fused feature. Furthermore, in the profile of classification performance with increasing training sample, we found that the SKLDA, STDA, SWLDA and RLDA provided an improving and robust accuracy for the fusion approach, and SKLDA was extraordinarily outstanding over other classifiers under various

scale of training set.

According to the review of Lotte et al, LDA is one of the most prevalent forms of classifiers for EEG-based BCI, especially for online and real-time processing [28]. The online analysis in the present study further validated the improvement effect of feature fusion in the selected classification algorithms. Notably, we tested the practicability of the hybrid system by modulating the amount of training set and the length of every input sample (i.e. system response time). Firstly, consistent with offline analysis, the accuracy and AUC maintained at a high level when 80 trials were taken for training and decreased with the contraction of training set. The amount of training sample determined the initial calibration time of the system and were directly related to the interaction convenience, but an extremely small training set, like 20 trials, also deteriorated the classification performance immensely. In addition, in order to investigate the effect of system response time to the decoding of interaction intention, we truncated the length of fixation duration and corresponding ERP epoch of each sample, for the sake of simulating different response time for intention recognition of the hybrid system and evaluating the classification performance. The accuracy and AUC declined when response duration decreased, with a slight slope from 800 ms to 400 ms and a substantial drop below 400 ms. From the perspective of eye-tracking data, using a shorter recognition time was harder to distinguish the intended selection from other conditions, such as long processing time towards stimulus for participants. Previous study showed that the dwell time of novices was typically between 450 and 1000 ms in gaze typing tasks, and decreased to 282 ms after repeated training [46]. In other eye selection related ERP studies, the fixation duration were usually set in a long threshold, such as 1000ms [18] or 2000ms [14]. Considering no pre-training was performed on our subjects, a proper recognition threshold above 400 ms could ensure a more accurate performance. From the aspect of ERP, it was obvious to find out that shorter EEG epoch after stimulus comprising less ERP components. Given that the predominate ERP was concentrated on the 100ms and 300 ms as shown in (Fig. 3), the epoch less than 400ms might induce the loss of crucial information especially when latency shift happened in P300 wave. However, the superiority of the fused feature was still observed compared with single fixation duration across the whole range of system response time (300 to 800 ms) as long as considerable training set was implemented. Single eye-tracker was natural to use and could reach a decent accuracy without much training, but only by a few calibrations, the fusion of ERP and fixation duration could outperform the former. Besides, consistent with offline analysis, SKLDA still maintained the best performance for the fused feature among all classification methods, showing the largest area over single eye-tracker classification performance in (Fig. 5), and a relatively less requirement of training set and response duration for the outperformance. Previous study proved that SKLDA remained effective when training set was insufficient [33], so that it was practical to be utilized in real-time scenarios. As for the computational time, we found that the single-trial process time was associated with the length of input signal, especially in RLDA. And the classifiers of

SKLDA and SWLDA required less computation cost in the practical applications.

In the present study, some factors should be considered when interpreting our results. Firstly, only young and healthy university students were included as the sample to test our protocol. But there is evidence suggesting that the performance of LDA declined largely for the ERP classification of the severely disabled in real-life applications [27]. Further study could extend the diversity of subjects to verify the robustness of fused features. Secondly, in the design of present interaction interface, the icons were arranged isolated and regularly with rigid distance as an array, to which the users would be accustomed with the progress of experiment. A real ecological application was supposed to display arbitrary targets, which provided a more oddball stimuli to the user and elicited stable ERP responses. At the same time, the stimulus was expressed in images in this study, while other type of interaction item, like words, was associated with different predominant ERP components [45]. Future study could introduce different forms of selection target to test the generalization of the hybrid system. In addition, this study observed a relatively limited increase of performance by fusing eye and ERP compared with single eye-tracker signal as the feature. Considering the computational complexity, only the time domain wave was taken from ERP in our pipeline. Other features could also be explored to overcome the shortcomings in ERP characteristics and optimize the performance of the hybrid interaction system for single trial classification.

V. CONCLUSION

In the current study, we introduced an eye-brain hybrid BCI interaction system and assessed the performance in a customized free visual search paradigm. In comparison with the single-model EEG or eye-track features, the proposed hybrid BCI system achieved better performance in both offline and online conditions. Furthermore, practical validation across six widely used classification methods showed that the SKLDA method could maintain superior performance under condition with few training set and fast response time. In sum, our study shed new insights on the approach of hands-free HCI and provided novel and practical solution to the intention detection in the real-world scenarios.

REFERENCES

- [1] C. Ware and H. H. Mikaelian, "An evaluation of an eye tracker as a device for computer input2," in *Proceedings of the SIGCHI/GI conference on Human factors in computing systems and graphics interface*, 1986, pp. 183–188.
- [2] T. E. Hutchinson, K. P. White, W. N. Martin, K. C. Reichert, and L. A. Frey, "Human-computer interaction using eye-gaze input," *IEEE Transactions on Systems, Man, and CBlankertzkybernetics*, vol. 19, no. 6, pp. 1527–1534, 1989.
- [3] R. J. Jacob, "The use of eye movements in human-computer interaction techniques: what you look at is what you get," *ACM Transactions on Information Systems (TOIS)*, vol. 9, no. 2, pp. 152–169, 1991.
- [4] D. M. Stampe and E. M. Reingold, "Selection by looking: A novel computer interface and its application to psychological research," in *Studies in visual information processing*. Elsevier, 1995, vol. 6, pp. 467–478.
- [5] L. E. Sibert and R. J. Jacob, "Evaluation of eye gaze interaction," in *Proceedings of the SIGCHI conference on Human Factors in Computing Systems*, 2000, pp. 281–288.

- [6] R. J. Jacob, "What you look at is what you get: eye movement-based interaction techniques," in *Proceedings of the SIGCHI conference on Human factors in computing systems*, 1990, pp. 11–18.
- [7] O. Špakov and P. Majaranta, "Enhanced gaze interaction using simple head gestures," in *Proceedings of the 2012 ACM Conference on Ubiquitous Computing*, 2012, pp. 705–710.
- [8] X. Gao, Y. Wang, X. Chen, and S. Gao, "Interface, interaction, and intelligence in generalized brain–computer interfaces," *Trends in Cognitive Sciences*, 2021.
- [9] B. T. Jap, S. Lal, P. Fischer, and E. Bekiaris, "Using eeg spectral components to assess algorithms for detecting fatigue," *Expert Systems with Applications*, vol. 36, no. 2, pp. 2352–2359, 2009.
- [10] G. N. Dimitrakopoulos, I. Kakkos, Z. Dai, H. Wang, K. Sgarbas, N. Thakor, A. Bezerianos, and Y. Sun, "Functional connectivity analysis of mental fatigue reveals different network topological alterations between driving and vigilance tasks," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 26, no. 4, pp. 740–749, 2018.
- [11] J. Atkinson and D. Campos, "Improving bci-based emotion recognition by combining eeg feature selection and kernel classifiers," *Expert Systems with Applications*, vol. 47, pp. 35–41, 2016.
- [12] H. Qi, Y. Xue, L. Xu, Y. Cao, and X. Jiao, "A speedy calibration method using riemannian geometry measurement and other-subject samples on a p300 speller," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 26, no. 3, pp. 602–608, 2018.
- [13] Y. Si, F. Li, K. Duan, Q. Tao, C. Li, Z. Cao, Y. Zhang, B. Biswal, P. Li, D. Yao *et al.*, "Predicting individual decision-making responses based on single-trial eeg," *NeuroImage*, vol. 206, p. 116333, 2020.
- [14] T. O. Zander and C. Kothe, "Towards passive brain–computer interfaces: applying brain–computer interface technology to human–machine systems in general," *Journal of Neural Engineering*, vol. 8, no. 2, p. 025005, 2011.
- [15] S. A. Hillyard and M. Kutas, "Electrophysiology of cognitive processing," *Annual Review of Psychology*, vol. 34, no. 1, pp. 33–61, 1983.
- [16] M. D. Rugg and M. G. Coles, "The erp and cognitive psychology: Conceptual issues," in *Electrophysiology of mind: Event-related brain potentials and cognition* (pp. 27–39). Oxford University Press, 1995.
- [17] M. Li, W. Li, L. Niu, H. Zhou, G. Chen, and F. Duan, "An event-related potential-based adaptive model for telepresence control of humanoid robot motion in an environment cluttered with obstacles," *IEEE Transactions on Industrial Electronics*, vol. 64, no. 2, pp. 1696–1705, 2016.
- [18] J. Protzak, K. Ihme, and T. O. Zander, "A passive brain–computer interface for supporting gaze-based human–machine interaction," in *International Conference on Universal Access in Human–Computer Interaction*. Springer, 2013, pp. 662–671.
- [19] E. Donchin, K. M. Spencer, and R. Wijesinghe, "The mental prosthesis: assessing the speed of a p300-based brain–computer interface," *IEEE Transactions on Rehabilitation Engineering*, vol. 8, no. 2, pp. 174–179, 2000.
- [20] G. F. Potts, M. Liotti, D. M. Tucker, and M. I. Posner, "Frontal and inferior temporal cortical activity in visual target detection: Evidence from high spatially sampled event-related potentials," *Brain Topography*, vol. 9, no. 1, pp. 3–14, 1996.
- [21] J. R. Wiersema, J. J. van der Meere, and H. Roeyers, "Developmental changes in error monitoring: an event-related potential study," *Neuropsychologia*, vol. 45, no. 8, pp. 1649–1657, 2007.
- [22] L. N. Kaunitz, J. E. Kamienskowski, A. Varatharajah, M. Sigman, R. Q. Quiroga, and M. J. Ison, "Looking for a face in the crowd: fixation-related potentials in an eye-movement visual search task," *NeuroImage*, vol. 89, pp. 297–305, 2014.
- [23] H. Devillez, N. Guyader, and A. Guérin-Dugué, "An eye fixation–related potentials analysis of the p300 potential for fixations onto a target object when exploring natural scenes," *Journal of Vision*, vol. 15, no. 13, pp. 20–20, 2015.
- [24] D. Kalika, L. Collins, K. Caves, and C. Throckmorton, "Fusion of p300 and eye-tracker data for spelling using bci2000," *Journal of Neural Engineering*, vol. 14, no. 5, p. 056010, 2017.
- [25] J.-S. Choi, J. W. Bang, K. R. Park, and M. Whang, "Enhanced perception of user intention by combining eeg and gaze-tracking for brain–computer interfaces (bcis)," *Sensors*, vol. 13, no. 3, pp. 3454–3472, 2013.
- [26] D. J. Krusienski, E. W. Sellers, D. J. McFarland, T. M. Vaughan, and J. R. Wolpaw, "Toward enhanced p300 speller performance," *Journal of Neuroscience Methods*, vol. 167, no. 1, pp. 15–21, 2008.
- [27] V. Martínez-Cagigal, J. Gomez-Pilar, D. Alvarez, and R. Hornero, "An asynchronous p300-based brain–computer interface web browser for severely disabled people," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 8, pp. 1332–1342, 2016.
- [28] F. Lotte, L. Bougrain, A. Cichocki, M. Clerc, M. Congedo, A. Rakotomamonjy, and F. Yger, "A review of classification algorithms for eeg-based brain–computer interfaces: a 10 year update," *Journal of Neural Engineering*, vol. 15, no. 3, p. 031005, 2018.
- [29] A. Craik, Y. He, and J. L. Contreras-Vidal, "Deep learning for electroencephalogram (eeg) classification tasks: a review," *Journal of neural engineering*, vol. 16, no. 3, p. 031001, 2019.
- [30] T.-P. Jung, S. Makeig, C. Humphries, T.-W. Lee, M. J. Mckeown, V. Iragui, and T. J. Sejnowski, "Removing electroencephalographic artifacts by blind source separation," *Psychophysiology*, vol. 37, no. 2, pp. 163–178, 2000.
- [31] A. Delorme and S. Makeig, "Eeglab: an open source toolbox for analysis of single-trial eeg dynamics including independent component analysis," *Journal of Neuroscience Methods*, vol. 134, no. 1, pp. 9–21, 2004.
- [32] X. Xiao, M. Xu, J. Jin, Y. Wang, T.-P. Jung, and D. Ming, "Discriminative canonical pattern matching for single-trial classification of erp components," *IEEE Transactions on Biomedical Engineering*, vol. 67, no. 8, pp. 2266–2275, 2019.
- [33] B. Blankertz, S. Lemm, M. Treder, S. Haufe, and K.-R. Müller, "Single-trial analysis and classification of erp components—a tutorial," *NeuroImage*, vol. 56, no. 2, pp. 814–825, 2011.
- [34] J. H. Friedman, "Regularized discriminant analysis," *Journal of the American statistical association*, vol. 84, no. 405, pp. 165–175, 1989.
- [35] N. R. Draper and H. Smith, *Applied regression analysis*. John Wiley & Sons, 1998, vol. 326.
- [36] X. Lei, P. Yang, and D. Yao, "An empirical bayesian framework for brain–computer interfaces," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 17, no. 6, pp. 521–529, 2009.
- [37] C. Vidaurre, N. Krämer, B. Blankertz, and A. Schlögl, "Time domain parameters as a feature for eeg-based brain–computer interfaces," *Neural Networks*, vol. 22, no. 9, pp. 1313–1319, 2009.
- [38] Y. Zhang, G. Zhou, Q. Zhao, J. Jin, X. Wang, and A. Cichocki, "Spatial-temporal discriminant analysis for erp-based brain–computer interface," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 21, no. 2, pp. 233–243, 2013.
- [39] H. Cecotti and A. Graser, "Convolutional neural networks for p300 detection with application to brain–computer interfaces," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 3, pp. 433–445, 2010.
- [40] A. Sharma and K. K. Paliwal, "A deterministic approach to regularized linear discriminant analysis," *Neurocomputing*, vol. 151, pp. 207–214, 2015.
- [41] D. J. Krusienski, E. W. Sellers, F. Cabestaing, S. Bayouthe, D. J. McFarland, T. M. Vaughan, and J. R. Wolpaw, "A comparison of classification techniques for the p300 speller," *Journal of Neural Engineering*, vol. 3, no. 4, p. 299, 2006.
- [42] U. Hoffmann, J.-M. Vesin, T. Ebrahimi, and K. Diserens, "An efficient p300-based brain–computer interface for disabled subjects," *Journal of Neuroscience Methods*, vol. 167, no. 1, pp. 115–125, 2008.
- [43] J. Schäfer and K. Strimmer, "A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics," *Statistical Applications in Genetics and Molecular Biology*, vol. 4, no. 1, 2005.
- [44] Y. Pan, X. Ge, L. Ge, and J. Xu, "Using eye-controlled highlighting techniques to support both serial and parallel processing in visual search," *Applied Ergonomics*, vol. in press, 2021.
- [45] T. Baccino, "Eye movements and concurrent event-related potentials: Eye fixation-related potential investigations in reading," in *The Oxford handbook of eye movements*, 2011.
- [46] P. Majaranta, U.-K. Ahola, and O. Špakov, "Fast gaze typing with an adjustable dwell time," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 2009, pp. 357–360.