# Hybrid brain-computer interface using error-related potential and reinforcement learning

Aline Xavier Fidêncio[1,2,3,4]*, Felix Grün[2,4], Christian Klaes[3] and
Ioannis Iossifidis[2]

[1]Faculty of Electrical Engineering and Information Technology, Ruhr University Bochum, Bochum,
Germany, [2]Robotics and BCI Laboratory, Institute of Computer Science, Ruhr West University of
Applied Sciences, Mülheim an der Ruhr, Germany, [3]KlaesLab, Department of Neurosurgery, University
Hospital Knappschaftskrankenhaus, Ruhr University Bochum, Bochum, Germany, [4]Faculty of
Computer Science, Ruhr University Bochum, Bochum, Germany

Brain-computer interfaces (BCIs) offer alternative communication methods for individuals with motor disabilities, aiming to improve their quality of life through external device control. However, non-invasive BCIs using electroencephalography (EEG) often suffer from performance limitations due to non-stationarities arising from changes in mental state or device characteristics. Addressing these challenges motivates the development of adaptive systems capable of real-time adjustment. This study investigates a novel approach for creating an adaptive, error-related potential (ErrP)-based BCI using reinforcement learning (RL) to dynamically adapt to EEG signal variations. The framework was validated through experiments on a publicly available motor imagery dataset and a novel fast-paced protocol designed to enhance user engagement. Results showed that RL agents effectively learned control policies from user interactions, maintaining robust performance across datasets. However, findings from the game-based protocol revealed that fast-paced motor imagery tasks were ineffective for most participants, highlighting critical challenges in real-time BCI task design. Overall, the results demonstrate the potential of RL for enhancing BCI adaptability while identifying practical constraints in task complexity and user responsiveness.

KEYWORDS

error-related potentials (ErrPs), adaptive brain-computer interface, BCI, reinforcement learning (RL), motor imagery (MI), EEG

## 1 Introduction

Rehabilitation and assistance systems can be used to improve life quality for patients living with motor impairments caused, for example, by an amputation, spinal cord injury, or stroke (Abiri et al., 2019; Soekadar et al., 2015; Kumar et al., 2019). Advances in the field of brain-computer interfaces (BCIs) provide patients with an alternative communication path to these systems. This is achieved through the direct decoding or classification of specific brain signals and their translation into appropriate control commands for the external systems. While different technologies can be used for neural signal acquisition, non-invasive electroencephalography (EEG) devices are widely applied due to their good temporal resolution, attractive price, and usability (Kumar et al., 2019).

BCIs can be developed based on different experimental paradigms and used to control different devices: from a cursor on the monitor to robotic arms, or wheelchairs, for example (Kumar et al., 2019). The experimental paradigms define, among others, which

kind of brain signal should be decoded, and common applications are based on event-related synchronization/desynchronization (ERS/ERD) modulations generated during motor imagination, steady-state visual evoked potentials (SSVEPs) or P300 potentials. For a comprehensive review of different paradigms, we refer the reader to Abiri et al. (2019).

However, the classification of brain signals is a challenging task, and a current limitation in developing high-performance BCI systems for long-term use is their decreasing performance over time due to the inherent non-stationarities in EEG data caused, for example, by changes in the subject's signals or the recording device itself, such as electrode placement and impedance. To address this problem in traditional BCIs, adaptive systems are proposed to dynamically adjust their behavior and parameters based on changes in the user's mental state, the environment, or the input data quality.

In the last years, several works have proposed using a specific brain signal elicited upon errors to improve BCIs. The so-called error-related potentials (ErrPs) can be elicited under different circumstances and measured with EEG (for a review, see Xavier Fidêncio et al., 2022; Chavarriaga et al., 2014; Kumar et al., 2019). In BCI research, seven different types of ErrPs are usually mentioned. Errors committed by the subject are called *response ErrPs* (Blankertz et al., 2002; van Schie et al., 2004). *Feedback ErrPs* are generated upon feedback about a choice made (Miltner et al., 1997; Chavarriaga et al., 2014) and *target ErrPs* can be generated by implementing unexpected changes in the task (Diedrichsen, 2005). However, in BCI paradigms it is more common to find the use of either *interaction ErrPs* (Ferrez and Millán, 2005, 2008), which are elicited when the interface wrongly interprets the user's input, or *observation ErrPs*, which are generated while the subject only observes a system over which they have no control perform a wrong action. *Execution* and *outcome ErrPs* are also reported as the neural response to unexpected movements (Diedrichsen, 2005; Spüler and Niethammer, 2015) or undesired outcomes (Krigolson et al., 2008; Spüler and Niethammer, 2015; Kreilinger et al., 2016), respectively. For an extensive review of each ErrP and experimental protocols used to elicit them see Xavier Fidêncio et al. (2022).

ErrPs have been combined with reinforcement learning (RL) in different studies to improve BCI performance (for a review, see Xavier Fidêncio et al., 2022). In the RL framework, an agent learns by trial and error while interacting with an environment (Sutton and Barto, 2018). The agent performs an action in the environment and receives a scalar numerical reward. Its goal is to maximize the cumulative reward, called return, and learn an optimal policy, that is mapping from inputs to actions. While supervised learning relies on ground-truth data being provided as a learning signal, RL agents need only a reward signal that indicates the quality of a policy to drive the learning process, making the use of ErrPs a natural fit, as they only represent the existence of an error, not what the expected outcome, action or observation would have been. In this work, because the agent's actions do not directly affect its next inputs, the setting is specifically a contextual bandit problem, where the agent's inputs do not represent states, as in the full RL problem, but rather a context that is independent between timesteps.

This study introduces a novel ErrP and RL-based BCI framework for the development of adaptive BCIs. The framework

applies reinforcement learning to learn the user intention directly from brain signals obtained with non-invasive recordings and uses the neural signature of errors measured in the form of interaction ErrPs to drive the learning. Our hypothesis is that the intrinsic interactive nature of the RL framework is particularly suitable for the development of such systems, inspired by the work of Kim et al. (2017). Moreover, as ErrPs are generated during human-system interaction upon BCI errors, it does not increase the mental load of the subject and directly constitutes a real-time feedback source for the RL agent.

To further validate the proposed framework, in this study, we also introduce a novel hybrid BCI paradigm using motor imagery and ErrPs. We propose a relatively fast-paced game to improve subjects' motivation and engagement. The hypothesis is that the gamified version of the commonly used cursor control task can increase subjects' motivation and interest in using the BCI, increasing overall performance (Škola et al., 2019; Atilla et al., 2024). Moreover, the increased game speed better aligns with real-time decision-making scenarios requiring faster reactions from participants.

The rest of this work is structured as follows: Section 2 reviews studies that have used ErrP for adaptive BCIs. We present the proposed ErrP-RL-based BCI framework and describe the datasets used in this study in Section 3. Section 4 presents our results. Finally, we conclude this work with a brief discussion and overview in Section 5.

## 2 Related work

Adaptive BCIs using ErrPs have been proposed previously. Llera et al. (2011) introduced an adaptive logistic regression based on interaction ErrPs. The weights of the classifier were modified based on the ErrP classification results. The approach was validated using both simulated and MEG data for a two-class MI paradigm, showing significant performance improvements compared to the baseline static classifier. Schiatti et al. (2019) later applied the same approach to MI data recorded with EEG. Mousavi et al. (2017) introduced a new strategy by directly combining the ErrP frequency-domain information and the MI-related modulations to improve the classification of MI trials. They used common spatial patterns (CSP) for feature extraction and linear discriminant analysis (LDA) for the classification of ErrPs and MI, combining the results with logistic regression, and observed significant improvements in performance with the proposed framework. This approach was further validated in an online follow-up study (Mousavi et al., 2020).

The ErrP information has often been used to validate the output of the BCI classifier. An online BCI-speller based on code-modulated visual evoked potentials (c-VEP) and ErrP was validated in Spüler et al. (2012). In this study, c-VEP trials were classified using a support vector machine (SVM) and a spatial filter (canonical correlation analysis). The ErrP information was used to label trials for the training dataset. Artusi et al. (2011) considered the classification of movement-related cortical potentials (MRCPs) into different motor tasks (e.g., slow vs. fast arm flexion). They also used the ErrP information to label trials before adding them to the training dataset for an SVM classifier. This approach was also used

recently by Tao et al. (2023) in a two-class MI task, using regularized common spatial patterns (R-CSP) for feature extraction and the combination of Fisher's discriminant analysis (FDA) and SVM for classification. Lastly, for the classification of MI data using k-NN, Haotian et al. (2023) also used the trials labeled based on the ErrP to create a dataset and, after applying cross-validation to evaluate MI classification improvement, they expanded the training dataset with the new trials. Chiang et al. (2021) used a similar approach to show the benefits of including the ErrP feedback information in the adaptation of a convolutional neural network (CNN) for the classification of steady-state visually evoked potentials (SSVEPs) for three classes, and Wang et al. (2024) for MI classification with four classes. In contrast to the previously mentioned studies, which included trials that elicited an ErrP after inverting the label, in these studies, the training dataset only included trials that did not elicit ErrPs. This is because with more than two classes, the presence of an ErrP is not enough to infer the correct label, it only indicates that the chosen label is not the correct one. All these studies reported improved performance when using the ErrP-based adaptation.

# 3  Materials and methods

We introduce a novel BCI framework using ErrP and reinforcement learning. Figure 1 shows the framework overview. Our hypothesis is that the intrinsic interactive nature of reinforcement learning agents is well-suited for the development of real-time adaptive BCIs. Moreover, ErrPs represent an intrinsic feedback source with no extra mental load given to the subject, as they are implicitly generated, even upon external error occurrences. Therefore, incorporating ErrP in the reinforcement learning framework as the reward is very straightforward. The setup is validated offline with a motor imagery paradigm for BCIs. We used a well-known open-source dataset for a two-class motor imagery protocol. Additionally, we propose a new BCI task designed to combine the MI and ErrP paradigms in a gamified setup and use part of the data we collected to also validate the proposed RL-based framework. Details on all these components are described in the following sections.

## 3.1  Open-source dataset: BCI competition IV dataset 2b

In this study, we used the open-source dataset 2b from the BCI Competition IV (Leeb et al., 2007). This dataset is widely used as a benchmark for the classification of motor imagery. The competition review describes the experimental protocols in detail (Tangermann et al., 2012). The dataset includes EEG and electrooculogram (EOG) data from nine subjects recorded in five sessions. In the competition, the first three sessions were intended for training, and the last two were for evaluation when validating proposed methods. However, different splits are commonly used in studies utilizing this dataset (for examples, see Ali et al., 2022). In summary, the experimental task was a cue-based screening paradigm. In the first two sessions, an arrow shown on the screen for 1.25 seconds indicated the MI task that the subjects should perform (either left or right hand, with the movement freely

chosen by each subject). Subjects had to imagine the corresponding hand movement for 4 seconds. Afterwards, a break of at least 1.5 seconds, followed by a randomized time of up to 1 second, was included (Tangermann et al., 2012). The last three sessions included smiley feedback. Subjects were instructed to move this smiley toward the left or right side according to the cue and keep the MI for as long as possible. A break and a random interval were also included at the end of these trials. EEG was recorded using three electrodes (C3, Cz, C4) with a sampling rate of 250 Hz. EOG was recorded using three monopolar electrodes. Supplementary Table S1 reports the number of trials recorded for each subject.

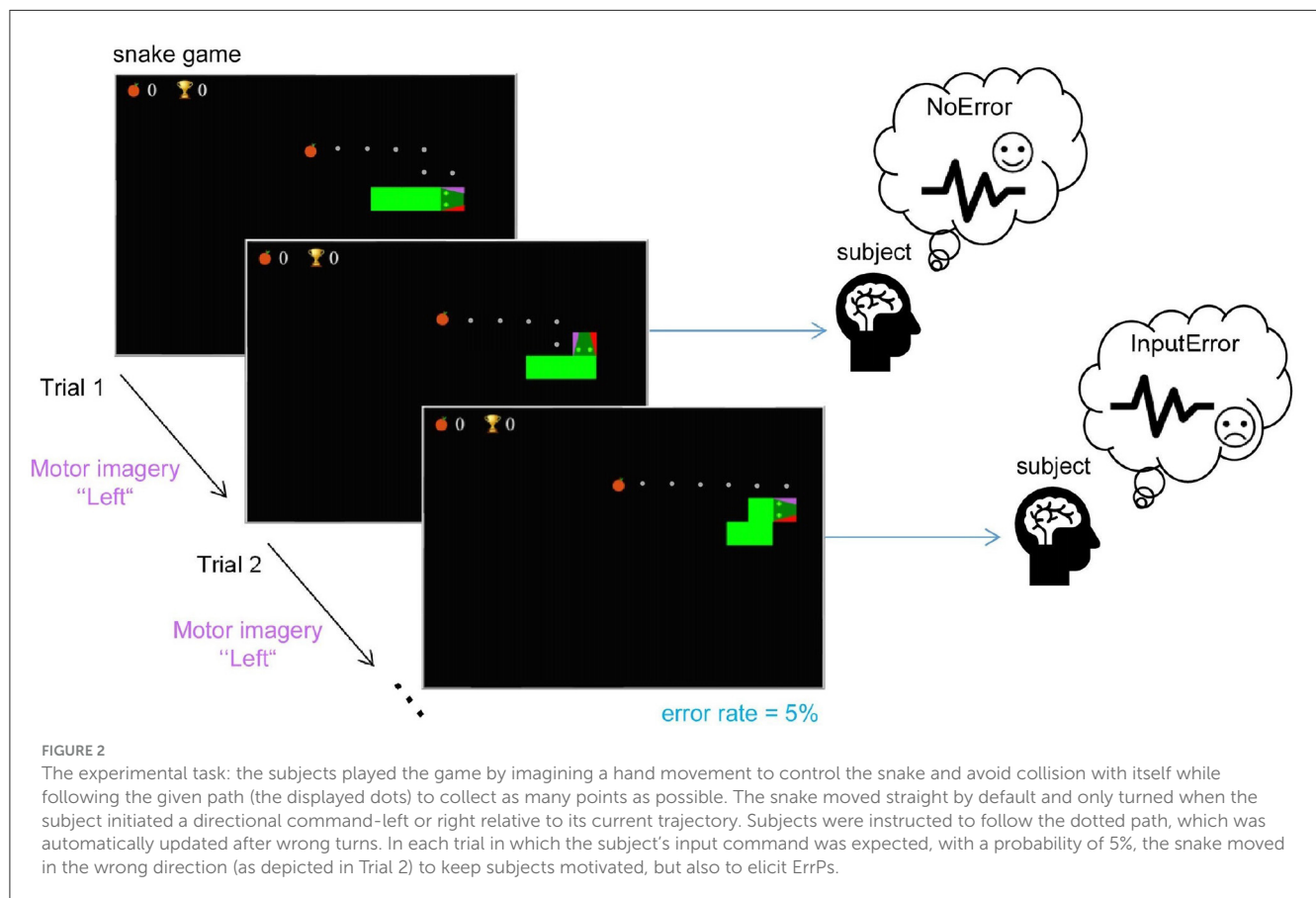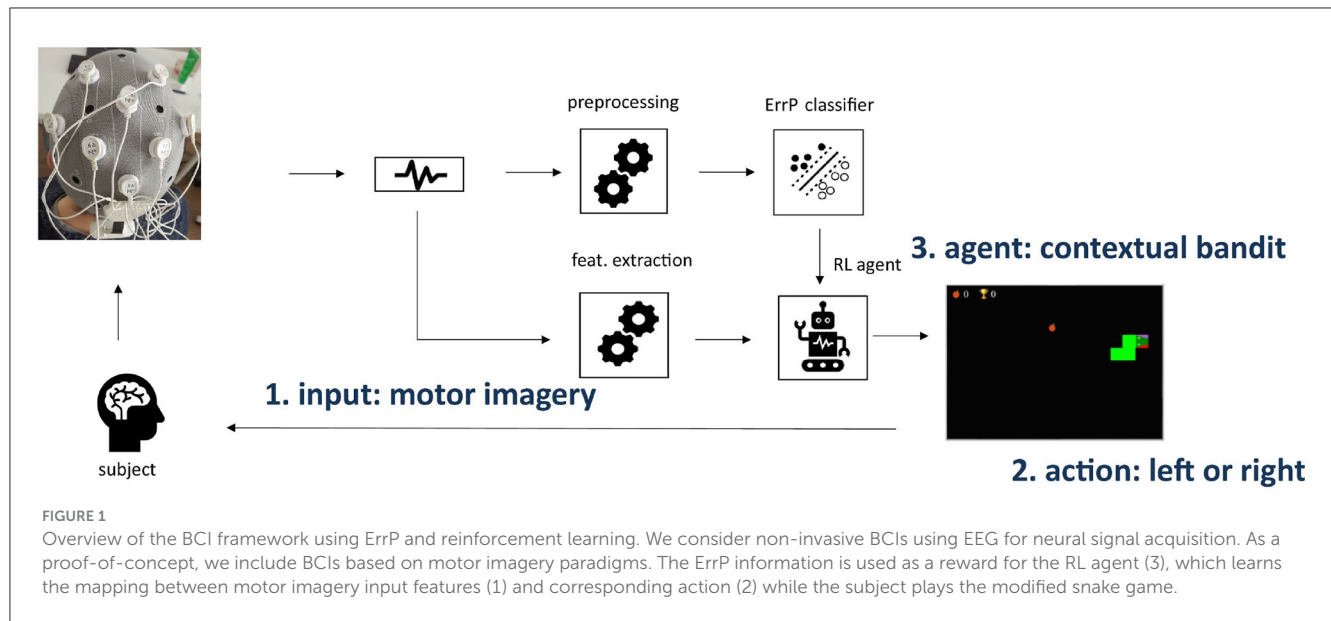## 3.2  Snake game BCI dataset

We implemented a new experimental protocol to validate the detection of motor imagery and ErrP-related neural signals. We developed a modified version of the snake game[1] (see Figure 2) to (1) propose an interactive game design controlled via motor imagery to keep subjects motivated and focused; (2) demonstrate the feasibility of detecting MI modulations with a fast-paced task; (3) demonstrate the feasibility of detecting interaction ErrPs in response to misinterpreted commands by the BCI at a low artificial error rate and simultaneously with the fast-paced MI control; (4) provide the basis for the development of an adaptive MI-based BCI using ErrP and reinforcement learning.

The study involving human participants was reviewed and approved by the Ethics Committee of the Medical Faculty of the Ruhr University Bochum. The participants provided their written informed consent prior to participation.

### 3.2.1  Experimental protocol

All subjects were instructed both verbally and with written instructions to imagine the movement of the left or right hand to interact with the game. Participants chose freely whether to imagine an open-hand gesture or squeezing a ball. The game included a given path from the snake to the fruit that subjects were instructed to follow. This ensured the ground-truth label for the MI trials (left or right hand). Note that, while subjects believed they were actively controlling the snake, we did not decode the MI data online in this study. We artificially introduced error trials to keep subjects motivated. With a 5% chance, the snake moved in the opposite direction as defined by the path. This also allowed us to in parallel demonstrate that interaction ErrPs can also be elicited with the proposed protocol (the commonly used error rate in ErrP studies is 20-30% for a review, see Xavier Fidêncio et al., 2022). We included a familiarization phase without artificial error activation to let subjects get used to the game and avoid self-made errors (wrong hand movement imagined). Moreover, subjects were familiar with the snake game itself as they first participated in a recording session with the keyboard version of the game. In this first session, the

---

1   game used as baseline: https://www.geeksforgeeks.org/snake-game-in-python-using-pygame-module/.

**FIGURE 1**
Overview of the BCI framework using ErrP and reinforcement learning. We consider non-invasive BCIs using EEG for neural signal acquisition. As a proof-of-concept, we include BCIs based on motor imagery paradigms. The ErrP information is used as a reward for the RL agent (3), which learns the mapping between motor imagery input features (1) and corresponding action (2) while the subject plays the modified snake game.



**FIGURE 2**
The experimental task: the subjects played the game by imagining a hand movement to control the snake and avoid collision with itself while following the given path (the displayed dots) to collect as many points as possible. The snake moved straight by default and only turned when the subject initiated a directional command-left or right relative to its current trajectory. Subjects were instructed to follow the dotted path, which was automatically updated after wrong turns. In each trial in which the subject's input command was expected, with a probability of 5%, the snake moved in the wrong direction (as depicted in Trial 2) to keep subjects motivated, but also to elicit ErrPs.

subjects used two keys on the keyboard to interact with the game. The data of this session is not included in this study.

The recordings were performed in three phases. In a pilot study to validate the game, we recorded one subject and used the data to define the preprocessing pipeline. As the results of this subject were very promising, we extended to six other subjects. As the

data analysis revealed significant modulations in the frequency domain as a response to the motor imagination as expected only for two subjects, we decided to include the Movement Imagery Questionnaire-3 (MIQ-3) (Williams et al., 2012) to assess the subject's ability to perform movement imagery before each recording. The study was then expanded to more subjects, and we

added a monetary compensation to attract participants. As none of the subjects had previous experience with motor imagery BCIs, we hypothesized that applying the questionnaire could help introduce the experimental protocol, further improving performance during the recordings. In this questionnaire, instructions are read to subjects to inform them about the movement they first had to physically perform and then imagine (using either internal visual, external visual, or kinesthetic imagery). After each mental task, they rated the ease/difficulty of performing the imagery on a 7-point Likert scale (1 - very hard to 7 - very easy to see/feel). Imagined movements included knee lift, jump, arm movement, and waist bending.

In total, thirty subjects participated in this study (thirteen males, age: 25.5 ± 3.78, one left-handed, all with normal or corrected-to-normal vision). We excluded the data from eight subjects due to artifacted EEG data and one subject because of excessive movements during recordings. The data of the remaining twenty-one subjects (ten male, age: 26.0 ± 4.21, one left-handed) were analyzed. Each subject performed up to 10 runs of 120 trials each, and between 600 and 900 trials were collected as recordings could be stopped at any time if subjects noticed they were losing focus. Subjects took a self-paced break between runs and were asked about continuing or ending the recording.

### 3.2.2 Data recording

EEG data was recorded with the Enobio wet EEG from Neuroelectrics at the following positions: FC1, FC2, C3, Cz, C4, CP1, CP2, and Pz. The sampling rate was set to 500 Hz, and CMS/DRL reference electrodes were fixed behind the right ear. All recordings were performed in a quiet room, and we turned off all electronic devices that were not required for the recording itself. The manufacturer's software uses a quality index (QI) instead of impedance control for the EEG channels. As recommended, recordings were only started when all channels showed a green or orange indicator. We implemented the use of a cotton swab soaked in skin-friendly disinfectant to remove hair between the electrode and scalp before gel application. Typically, all electrodes were green right after being filled with an appropriate amount of conductive gel.

### 3.2.3 Data analysis

EEG data was analyzed in MATLAB® using the open-source toolbox EEGLAB (Delorme and Makeig, 2004). We performed a manual artifacted data rejection on continuous data, as recommended in the EEGLAB documentation, to avoid spreading artifacts over good-quality data. The data was then filtered with a Hamming windowed sinc FIR filter in the range [0.5, 100] Hz. As recommended in EEGLAB, low and high-pass filters were also applied separately. Additionally, a notch filter was applied to reduce power line noise. Lastly, the data was epoched, and, if necessary, artifacted trials were removed. Trials were extracted in the time interval [-1.0, 2.0] seconds around the snake's movement onset. We excluded epochs containing automatic forward movements of the snake. Supplementary Table S2 reports the number of trials included for each subject.

## 3.3 Reinforcement learning agents

As stated in the introduction, the task faced by the agent in this setting, is not the full RL problem. This is because the agents input in timestep $t + 1$ does not depend on the agents action in timestep $t$. Specifically, the subject's EEG data is unaffected by the agent's classification of the data from the previous timestep.

In this study, we have applied the contextual bandit algorithm LinUCB (Li et al., 2010) and its deep-learning counterpart, NeuralUCB (Zhou et al., 2020), to validate our framework. The adoption of LinUCB was motivated by the results reported by Kim et al. (2017) and NeuralUCB was chosen for comparison. LinUCB models the reward of each action as a linear function of the given context. It builds a linear estimator for each action and chooses actions using an upper confidence bound strategy, favoring options that are some combination of promising and uncertain. NeuralUCB generalizes this approach by utilizing a neural network instead of a linear model, enabling it to learn complex, nonlinear mappings between context and expected rewards. In our setup, both agents used features extracted from EEG signals as their input context. While we opted for these contextual bandit agents, it is important to note that the framework could accommodate other types of contextual bandit agents as well. The agents were implemented in python and are publicly available.[2]

The agents receive input data derived from human EEG signals. Their task is to learn the mapping between motor imagery-related features and intended actions. This learning process receives direct feedback through interaction ErrPs, which can also be obtained from EEG data. A reward of 1 was assigned to an action when no ErrP was detected, and a reward of 0 was given otherwise. As in this study we only validate the framework offline, a perfect ErrPs classification was assumed using the true data labels which, in this case, are known beforehand. To simulate the online use of the proposed BCI framework, the motor imagery data was streamed trial-by-trial to the RL agents.

We used optuna (Akiba et al., 2019) to optimize the hyperparameters of the agents for each subject individually. For LinUCB, $\alpha$ was searched over (0.01, 0.1, 1, 2, 4, 10) and for NeuralUCB we optimized the network hidden size (16, 32, 64, 128, 256, 512), $\nu$ (0.1, 1, 10), $\lambda$ ($10^{-i}, i = 1, 2, 3, 4$), and the learning rate ($2 \times 10^{-i}, 5 \times 10^{-i}, i = 1, 2, 3, 4$), as in Zhou et al. (2020).

The motor imagery features were extracted from the EEG data using continuous wavelet transform (CWT) following the methods used in several studies (Ali et al., 2022; Lee and Choi, 2019). With CWT, we can obtain a time-frequency representation of the EEG trials. The feature extraction was implemented in python using the library MNE-python (Gramfort et al., 2013), and the Morlet wavelet was applied to a two-second epoch extracted from cue-onset for the snake dataset and 0.5 seconds from cue-onset for the open-source dataset. Then, we extracted both mu [6, 13] Hz and beta [17, 30] Hz band power from the wavelet coefficients. This procedure results in a two-dimensional feature matrix (number of samples in frequency and time axes, respectively) with different dimensions for the frequency of the mu and beta bands. To achieve

---

2   LinUCB: https://www.kaggle.com/code/phamvanvung/linucb/notebook | NeuralUCB: https://github.com/uclaml/NeuralUCB/.

equal representation and avoid bias toward one frequency band, we resized the feature matrices to $15 \times 32$ using cubic interpolation as done in Ali et al. (2022) and Lee and Choi (2019). These features were extracted for all channels available (C3, Cz, C4) and combined into a single, flattened vector for the RL agent.

# 4 Results

## 4.1 BCI competition IV dataset 2b

The performance of the contextual bandit agents for all subjects in the open-source dataset is illustrated in Figure 3. We ran the agents for multiple random seeds and report the average results. As the LinCUB agent is CPU-bound, it was computationally expensive, and we only executed five seeds. NeuralUCB benefits from GPU computation, and we were able to use ten seeds. The results show that, for most subjects, both agents can learn the mapping of motor imagery input features into actions with reasonable accuracy. A two-sided Wilcoxon signed-rank test for accuracy shows no statistically significant difference between the two agents (two-sided $p$-value: 0.91). For two subjects (B02 and B03), both agents perform close to the chance level. This can be explained by the lack of feature separability between the two motor imagery classes for these two subjects, which is directly reflected in the performance of supervised learning approaches in these datasets as well (see Ali et al., 2022).

We also evaluated the performance of the agents on the training data by splitting the trials in two (first- vs. second-half of the trials) and testing whether its performance improved over time. One-sided Wilcoxon signed-rank tests showed that, for both agents, the accumulated number of errors was significantly smaller in the second half of the training session compared to the first half ($p = 0.002$, for both agents). This further indicates that the agents were able to learn while interactively receiving new motor imagery trials as input.

## 4.2 Snake game BCI dataset

### 4.2.1 Neurophysiological analysis of the MI data

The mean event-related changes in spectral power compared to the pre-stimulus baseline for the first subject used to validate the snake game with motor imagery control are depicted in Figure 4.[3]

We used EEGLAB to generate the event-related spectral perturbation (ERSP) image. The following parameters were set: wavelet cycle parameter ($[2, 0.1]$), pre-stimulus baseline $[-750, -500]$ ms, and frequency range of $[3, 30]$ Hz. We look directly at channels C3 and C4, as MI-related modulations can be measured at electrodes located over the sensorimotor cortex (Abiri et al., 2019). The ERSP images for this subject show the expected motor imagery-related modulation in the contralateral hemisphere. For example, for the right-hand motor imagination (MIRight), an event-related desynchronization (ERD) is visible in the C3 electrode. This ERD is more pronounced in the frequency ranges

within the expected mu $[6, 13]$ Hz and beta $[17, 30]$ Hz bands (Abiri et al., 2019). For each channel, we additionally show the statistical comparison of the ERSPs for the two experimental conditions, which highlight that the differences observed in mu and beta ranges are statistically significant (permutation test with 800 permutations and using false-discovery rate to correct for multiple comparisons. In Figures 4–6, as well as in Appendix Figures 9, 10, $p$-values below the significance level of 0.05 are shown in red).

As described in Section 3, we used the data recorded with this pilot subject to define the preprocessing pipeline and validate the feasibility of the experimental task for detecting motor imagery-related neural activity. Given the results shown in Figure 4, we extended the study to more subjects. Figure 5 shows the mean event-related changes in spectral power for all subjects included in the first extended study ($n = 7$), including S07. While ERSP plots show some contralateral ERD for the experimental conditions, the statistical comparison does not show significant differences.

Finally, Figure 6 shows the mean event-related changes in spectral power for all subjects included in the final extended study ($n = 21$). Also in this case, while some ERD is visible, the differences across experimental conditions for all subjects are not significant as seen for the pilot subject (S07).

The questionnaires to assess the motor imagery abilities were analyzed after the recordings, and the scores are summarized in Supplementary Table S3. If we consider a threshold for the total motor imagery ability score at 75% of the maximum (score of 15.21), only four subjects scored below this level. However, we found that a high score in the questionnaire did not imply significant ERSP modulations. Subject S07 did not reach the highest score and, still, for no other subject in this dataset such significant ERSP modulations were observed. Subject S12 obtained a low score. However, statistically significant differences were found in the ERSP modulations, especially during left-hand motor imagination. On the other hand, subject S29 reached the highest score in the questionnaire but no significant ERSP modulations can be seen (see Appendix Figures 9, 10). Nonetheless, we believe that applying the questionnaire helped improving subjects' understanding of the concept of motor imagination and how it can be performed. However, one further aspect to consider is the increased overall experimental time when such questionnaire is applied.

Overall, the results obtained with the pilot subject S07 demonstrate the feasibility of using the proposed fast-paced motor imagery paradigm. On the other hand, it is intriguing to us that only one particular subject performed remarkably well. We extended the study first to seven and then to thirty subjects to validate the protocol with a broader audience with the expectation of having more subjects with such significant modulations. It is unclear to us how this particular subject differs from the others such that classification performance is so outstanding.

### 4.2.2 Neurophysiological analysis of the ErrPs

As ErrPs are fronto-centrally located, higher amplitudes are expected at channels FCz and Cz (Xavier Fidêncio et al., 2022). The pre-stimulus interval $[-0.2, 0]$ s was used for baseline correction and the ErrPs are calculated as the difference of error trials minus correct trials. Figure 7 shows the ErrP grand averages for the correct

---

3   Please note that this subject, assigned the ID S07, also participated in a previous ErrP-only study.
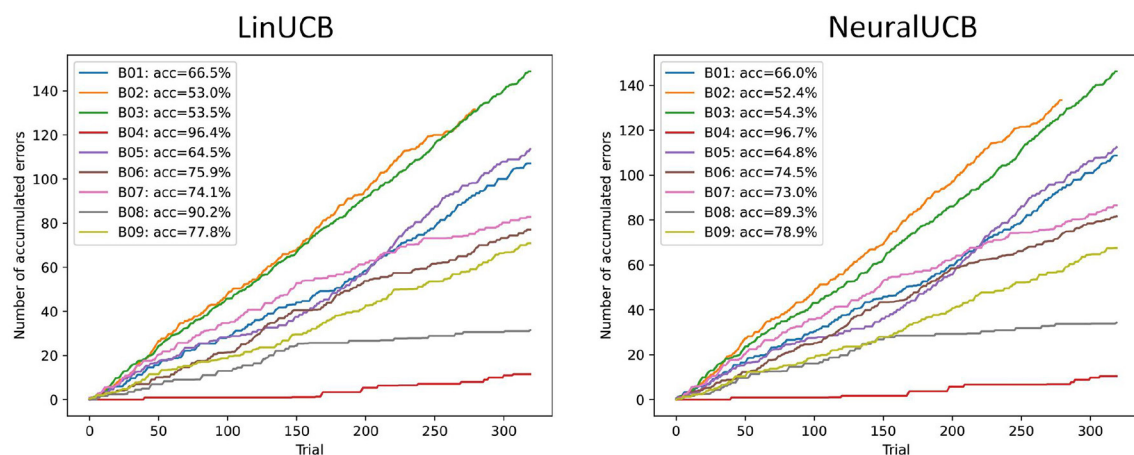
**FIGURE 3**
Performance of two contextual bandit agents (LinUCB and NeuralUCB) in the evaluation sessions for the open-source dataset ($n = 9$). Results are averaged over different seeds (five and ten, respectively). The plots show the accumulated number of errors across all trials. The accuracy is calculated based on the final accumulated regret. Results show that both agents perform reasonably well for all except two subjects (B02, B03). There is no statistically significant difference in the performance of both agents (two-sided Wilcoxon signed rank test, $p = 0.91$).
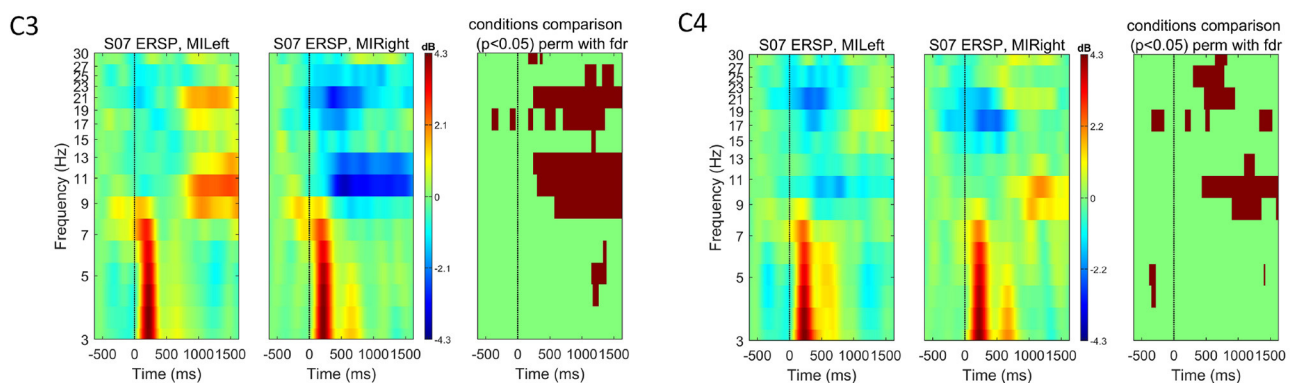


**FIGURE 4**
Event-related spectral perturbation (ERSP) for one subject (S07) at channels C3 and C4 for both left and right-hand motor imagery tasks. The color bars show the color and power spectral density in dB. For each channel, we used EEGLAB to compare the two experimental conditions (left versus right) and show in the right-most plot the permutation results (800 permutations, using false-discovery rate correction for multiple comparisons, significant $p$-values shown in red for $\alpha = 0.05$). These plots highlight how motor-imagery-related spectral modulations could be measured for this subject, with modulations mostly visible in the frequency ranges of [10–13] Hz and [16–30] Hz.

and error conditions measured at channel Cz. The measured ErrP displays a positive peak at 200 ms, a negative peak at 252 ms, and a positive peak at 348 ms. The statistical comparison shows significant differences between the error and correct conditions. The observed waveform is consistent with existing literature on interaction ErrPs (Ferrez and Millán, 2005, 2008; Ferrez and Del R Millan, 2008; Ferrez and Millán, 2009). However, the expected negative peak between 430-550 ms is not clearly visible with this experimental protocol. In our previous study, when subjects interacted with the game via keypress, this component was also visible in the ErrP (for details, see Xavier Fidêncio et al., 2024).

### 4.2.3 Reinforcement learning results

Considering the results from the data analysis on the MI-related ERSP modulations, this study only applied the agent to a

subset of eight subjects. We selected subjects based on the presence of at least some significantly different modulations between left and right-hand motor imagination, indicating the potential for sufficient class separability. Furthermore, considering the higher computation time for running the CPU-bound LinUCB agent, with the implementation used in this study, we only applied the NeuralUCB agent for this dataset, as results on the open-source dataset were very similar between linear and neural UCB agents.

For each subject, the data was randomly shuffled, and we used a simple train-test split to create the training and evaluation datasets (80/20). Figure 8 shows the performance of the NeuralUCB agent in the evaluation datasets. The pilot subject (S07) achieves the highest accuracy. This was expected and simply reflects the quality of the input features and the higher class separability. In general, the results obtained for subject S07 support our hypothesis that (1) the proposed experimental protocol can be used to elicit MI-related
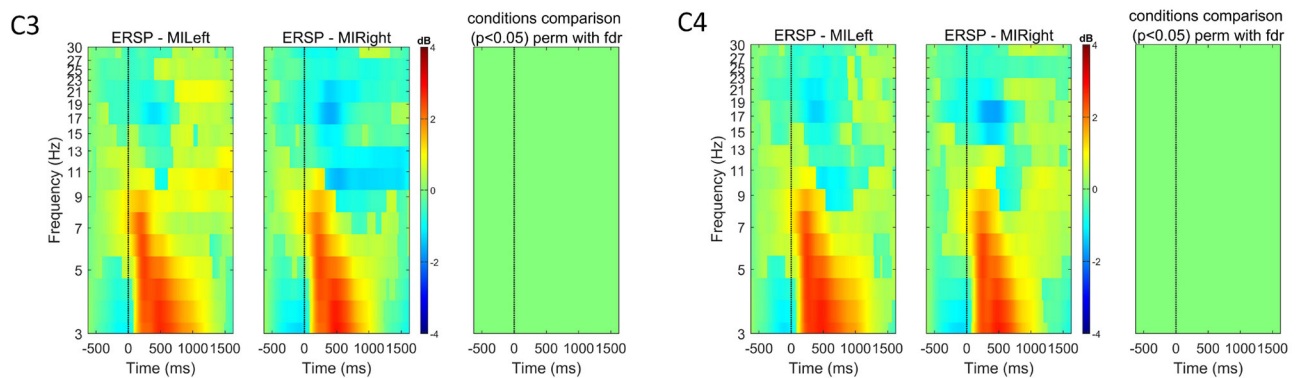
FIGURE 5

Event-related spectral perturbation (ERSP) for all subjects ($n = 7$) at channels C3 and C4 for both left and right-hand motor imagery tasks. The color bars show the color and power spectral density in dB. For each channel, we used EEGLAB to compare the two experimental conditions (**left** vs. **right**) and show in the right-most plot the permutation results (800 permutations, using false-discovery rate correction for multiple comparisons, for $\alpha = 0.05$). These plots show that, although some ERD is visible, the differences across experimental conditions are not significant when considering all subjects.
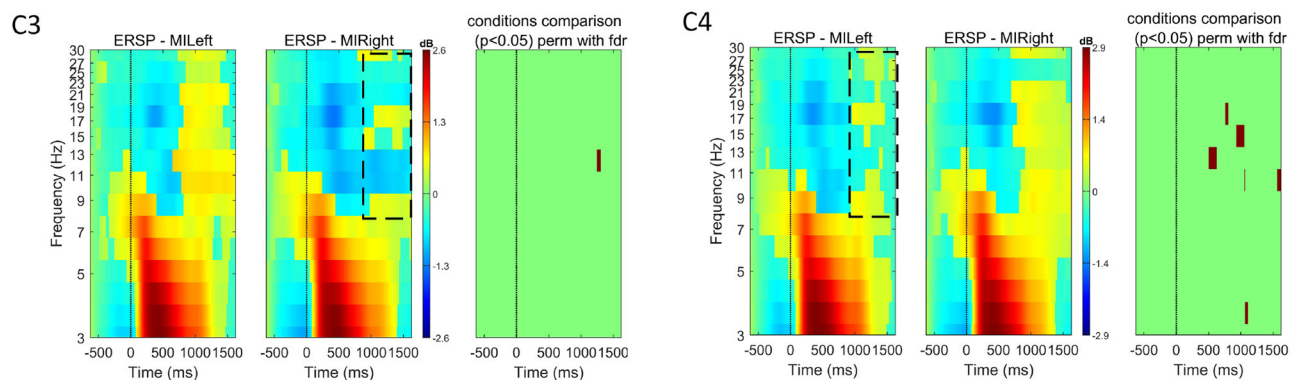


FIGURE 6

Event-related spectral perturbation (ERSP) for all subjects ($n = 21$) at channels C3 and C4 for both left and right-hand motor imagery tasks. The color bars show the color and power spectral density in dB. For each channel, we used EEGLAB to compare the two experimental conditions (left versus right) and show in the right-most plot the permutation results (800 permutations, using false-discovery rate correction for multiple comparisons, for $\alpha = 0.05$). These plots show that, some ERD is visible (see dotted areas for each condition in the contralateral hemisphere), the differences across experimental conditions are not significant when considering all subjects.
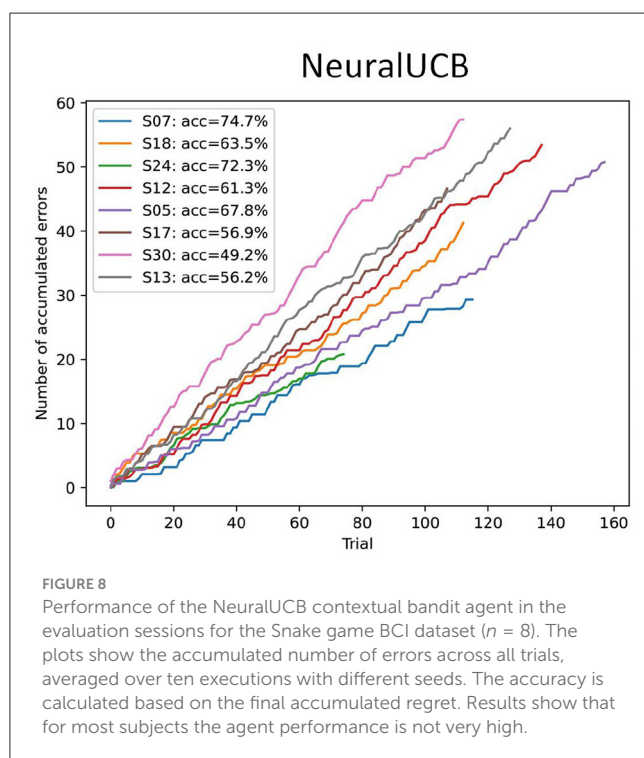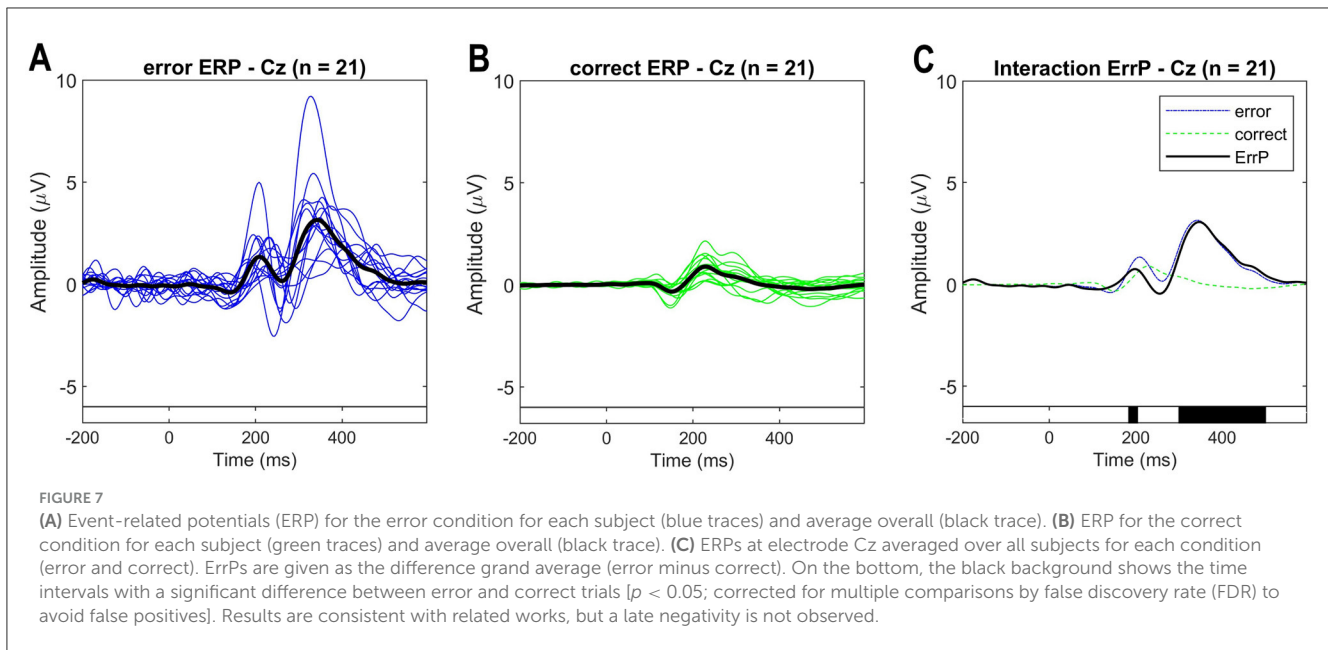
modulation and (2) a reinforcement learning agent can be used to learn the mapping between MI input features and intended action based on the feedback from the ErrP. On the other hand, the low-class separability obtained for most subjects with our experimental protocol directly impacts the agent's performance, highlighting the need for high-quality MI data to make the proposed ErrP-RL-based BCI framework feasible.

# 5  Discussion

This paper introduces a novel framework using error-related potentials (ErrPs) and reinforcement learning (RL) for the development of adaptive non-invasive brain-computer interfaces (BCIs). This study explores the use of a contextual bandit agent to learn the mapping between motor imagery-related features and intended actions, demonstrating the feasibility and effectiveness of this approach in interpreting and responding to neural signals

associated with motor imagery (MI). The learning framework was applied to both an open-source and an in-house dataset we recorded using a new experimental protocol. The results indicate that the agents are able to learn from the time-frequency domain features extracted from EEG recordings with reasonable accuracy using the simulated perfect ErrP classification as a reward. Moreover, the pilot study with the novel experimental task for using motor imagery and ErrPs in a fast-paced, interactive BCI suggests the feasibility of the introduced protocol.

Results with two contextual bandit algorithms (LinUCB and NeuralUCB) using the open-source BCI Competition IV dataset 2b for a two-class MI task show that both agents perform similarly and are able to learn the mappings between MI features extracted using the continuous wavelet transform, and the classes (left or right). The results obtained with selected subjects recorded with the novel MI-ErrP experimental task confirm the feasibility of the proposed RL-based MI framework, complementing the results obtained with the open-source dataset. However, as in other domains of machine

**FIGURE 7**
**(A)** Event-related potentials (ERP) for the error condition for each subject (blue traces) and average overall (black trace). **(B)** ERP for the correct condition for each subject (green traces) and average overall (black trace). **(C)** ERPs at electrode Cz averaged over all subjects for each condition (error and correct). ErrPs are given as the difference grand average (error minus correct). On the bottom, the black background shows the time intervals with a significant difference between error and correct trials [$p < 0.05$; corrected for multiple comparisons by false discovery rate (FDR) to avoid false positives]. Results are consistent with related works, but a late negativity is not observed.



**FIGURE 8**
Performance of the NeuralUCB contextual bandit agent in the evaluation sessions for the Snake game BCI dataset ($n = 8$). The plots show the accumulated number of errors across all trials, averaged over ten executions with different seeds. The accuracy is calculated based on the final accumulated regret. Results show that for most subjects the agent performance is not very high.

learning, reinforcement learning performance is highly dependent on the quality of features and, therefore, of the raw input data. The accuracies for individual subjects obtained in both datasets vary from close to chance level (e.g., for subjects B02 and B03 in the open-source dataset) to very high (e.g., for subject B04 in the open-source dataset). Nevertheless, as the aim of this study was to validate the overall feasibility of the proposed learning framework, rather than obtaining optimal feature extraction performance, and performing MI can be a challenging task,

we are confident that with higher quality MI features agents' performance can be further improved in the proposed framework. Different preprocessing and feature extraction pipelines can be investigated together with the proposed framework to improve overall discrimination accuracy.

In this work, we used BCIs with a binary output. However, as the framework is based on reinforcement learning (RL), we believe the setup can be easily extended to non-binary tasks and higher action spaces. This is, in fact, one of the main advantages of the proposed approach compared to the related works reviewed in Section 2. Most studies considered the binary case, where, upon ErrP detection, the true class label could be inferred as the opposite label (with some uncertainty due to the imperfect ErrP classification accuracy). In the few studies that considered more than two classes, only trials that did not elicit ErrPs were used, and trials for which the true label could not be inferred were discarded. In contrast, in an RL framework, such as the one proposed in this work, the agent can learn from every sample based on the reward received, regardless of the number of classes (actions) in the output. Hence, future work could extend this framework to more classes, for example, using the open-source dataset 2a from the BCI Competition IV for a four-class MI task (Tangermann et al., 2012). Moreover, the framework is not limited to the decoding of MI-related responses, and its application to other signals can also be evaluated, using benchmark datasets for offline validation or including online experiments to address real-time feasibility.

While RL shows potential for dealing with the non-stationarities in the EEG data due to its adaptability, its successful application in BCIs requires overcoming some challenges, such as noisy EEG data, limited training data, and the design of the reward function. In this work, the latter was addressed by including ErrPs as reward information, as it is intrinsically generated during interaction with BCIs. In the proposed framework, mistakes made by the BCI in the form of wrongly classified MI trials are expected to elicit interaction ErrPs, which can be used to provide

feedback for the agent in the form of reward (or penalty). On the other hand, as the learning results show, the agents' performance is significantly reduced if the classes are less clearly separable. This is the case for subjects B02 and B03 in the open-source dataset, for whom other works also report reduced classification accuracies and three out of eight subjects in the in-house collected dataset. Therefore, future work should further investigate different feature extraction methods to improve input data quality for the agents, and further validate the learning framework to establish performance boundaries and minimal requirements.

In this study, we also introduce a novel experimental paradigm for the development of a hybrid BCI using motor imagery and ErrPs. An interactive snake game was proposed to increase the subject's motivation and mitigate issues like boredom and reduced attention that commonly happens in repetitive BCI tasks. The increased game speed was defined based on previous studies and feedback from subjects, who were much less interested when playing the slower game. Moreover, a shorter interval better aligns with real-time decision-making scenarios that would require a faster reaction from subjects than commonly used slower-paced MI tasks. Data collected with twenty-one subjects show the feasibility of eliciting ErrPs under a low error rate of 5% while subjects perform MI in a fast-paced task. Data analysis of the MI data shows the expected mu and beta band modulations, but significant differences across experimental conditions (left- versus right-hand motor imagination) are only visible for a few subjects, with one subject performing particularly well.

While the two-second step deviates from the commonly used timings (usually 3-5 seconds), the proposed protocol provided insights into the feasibility of fast-paced MI tasks and the real-life deployment of such BCIs. On the other hand, such a short inter-trial interval might also not have been sufficient for most subjects to account for awareness of the game state, cognitive preparation for the MI task, and execution of the MI task with robust neural activation, leading to the observed low-class separability and consequently reduced decoding accuracy. Moreover, while subjects might be more motivated by playing, the continuous nature of the game might increase the cognitive load too much, leading to faster mental fatigue or inconsistent performance over time. Furthermore, all recorded subjects had no previous experience with motor imagery. This can be learned and improved over several training sessions (Tao et al., 2023). Therefore, results indicate that future research should also validate the protocol with experienced subjects or include training sessions. Extending the experiment and analysis protocol this way should rule out some possible reasons for insufficient data quality, enabling systematic validation of the feasibility of the proposed experimental protocol. Another aspect that could be included is the evaluation of whether including a higher reaction time between feedback presentation and the start of the motor imagination is required. Moreover, even though we already included breaks between experimental blocks, we would like to evaluate the quality of the motor imagination with shorter blocks (e.g., only 20 trials per block instead of 120).

Another aspect of the proposed task to consider is that incorporating the pre-programmed paths in the game helped ensure true labels for the MI trials. This can also be particularly helpful for recording labeled data for training classifiers using supervised learning approaches, which are widely used in BCI development. On the other hand, if subjects doubt their influence on the game control, this can reduce their motivation and the quality of the MI data. We also artificially introduced ErrPs in the task. In the envisioned online BCIs, these signals are generated because the BCI misinterpreted the subject's intention. In both cases, the interest in the BCI can also reduce if too many error trials are spotted.

Lastly, in this study we analyzed the proposed framework offline and assumed a perfect ErrP classification. In reality, ErrP classification accuracy will most likely be lower. Our ongoing work also considers the systematic validation of the proposed framework considering different rates of ErrPs misclassification to understand the performance boundaries for a general ErrP-based RL framework for adaptive BCIs. Future work should validate the entire framework during online use, as this is the intended application of BCI systems.

In summary, the development of non-invasive BCIs using EEG data usually requires the design of subject-specific classifiers to decode the neural modulations of interest for each specific paradigm. Not only must these classifiers be calibrated before use, but their performance might also degrade over time due to non-stationarities in the EEG signals. Some works have proposed different re-calibration strategies to update the classifiers and account for changes during long-term BCI use. In this work, we proposed and validated a new framework based on error-related potentials (ErrPs) and reinforcement learning for the development of adaptive BCIs. We hypothesize that RL methods have the potential to deal with the non-stationarities of EEG signals and, by using the intrinsic ErrP generation as a reward, they can constitute the fundamental block for the development of adaptive BCIs.

## Data availability statement

The datasets presented in this article are not readily available because the dataset used for this study can be obtained from the corresponding author on a reasonable request. Requests to access the datasets should be directed to aline.xavierfidencio@rub.de.

## Ethics statement

The studies involving humans were approved by the Ethics Committee of Medical Faculty of the Ruhr University Bochum. The studies were conducted in accordance with the local legislation and institutional requirements. The participants provided their written informed consent to participate in this study.

## Author contributions

AX: Conceptualization, Formal analysis, Investigation, Methodology, Software, Visualization, Writing – original draft, Writing – review & editing. FG: Methodology, Writing – original draft, Writing – review & editing. CK: Conceptualization, Funding acquisition, Resources, Supervision, Writing – original draft, Writing – review & editing. II: Conceptualization, Funding acquisition, Resources, Supervision, Writing – original draft, Writing – review & editing.

## Funding

## Acknowledgments

## Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## Generative AI statement

The author(s) declare that no Gen AI was used in the creation of this manuscript.

## Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

## Supplementary material

The Supplementary Material for this article can be found online at: https://www.frontiersin.org/articles/10.3389/fnhum. 2025.1569411/full#supplementary-material

## References

Abiri, R., Borhani, S., Sellers, E. W., Jiang, Y., and Zhao, X. (2019). A comprehensive review of EEG-based brain-computer interface paradigms. *J. Neural Eng.* 16:011001. doi: 10.1088/1741-2552/aaf12e

Akiba, T., Sano, S., Yanase, T., Ohta, T., and Koyama, M. (2019). "Optuna: a next-generation hyperparameter optimization framework," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery Data Mining* (Anchorage, AK, USA: ACM), 2623–2631. doi: 10.1145/3292500.3330701

Ali, O., Saif-ur Rehman, M., Dyck, S., Glasmachers, T., Iossifidis, I., and Klaes, C. (2022). Enhancing the decoding accuracy of eeg signals by the introduction of anchored-stft and adversarial data augmentation method. *Sci. Rep.* 12:4245. doi: 10.1038/s41598-022-07992-w

Artusi, X., Niazi, I. K., Lucas, M.-F., and Farina, D. (2011). Performance of a simulated adaptive bci based on experimental classification of movement-related and error potentials. *IEEE J. Emer. Selected Topics Circ. Syst.* 1, 480–488. doi: 10.1109/JETCAS.2011.2177920

Atilla, F., Postma, M., and Alimardani, M. (2024). Gamification of motor imagery brain-computer interface training protocols: a systematic review. *Comput. Hum. Behav. Rep.* 16:100508. doi: 10.1016/j.chbr.2024.100508

Blankertz, B., Schäfer, C., Dornhege, G., and Curio, G. (2002). "Single trial detection of EEG error potentials: a tool for increasing BCI transmission rates," in *Artificial Neural Networks — ICANN 2002* (Berlin Heidelberg: Springer), 1137–1143. doi: 10.1007/3-540-46084-5_184

Chavarriaga, R., Sobolewski, A., and Millna, J. d. R. (2014). Errare machinale est: the use of error-related potentials in brain-machine interfaces. *Front. Neurosci.* 8:208. doi: 10.3389/fnins.2014.00208

Chiang, K.-J., Emmanouilidou, D., Gamper, H., Johnston, D., Jalobeanu, M., Cutrell, E., et al. (2021). "A closed-loop adaptive brain-computer interface framework: improving the classifier with the use of error-related potentials," in *2021 10th International IEEE/EMBS Conference on Neural Engineering (NER)* (Italy: IEEE), 487–490. doi: 10.1109/NER49283.2021.9441133

Delorme, A., and Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *J. Neurosci. Methods* 134, 9–21. doi: 10.1016/j.jneumeth.2003.10.009

Diedrichsen, J. (2005). Neural correlates of reach errors. *J. Neurosci.* 25, 9919–9931. doi: 10.1523/JNEUROSCI.1874-05.2005

Ferrez, P., and Millán, J. (2009). "EEG-based brain-computer interaction: improved accuracy by automatic single-trial error detection," in *Advances in Neural Information Processing Systems 20 (NIPS 2007)*, 8.

Ferrez, P. W., and Del R Millan, J. (2008). "Simultaneous real-time detection of motor imagery and error-related potentials for improved BCI accuracy," in *Proceedings of the 4th International Brain-Computer Interface Workshop and Training Course*, 7.

Ferrez, P. W., and Millán, J. d. R. (2005). "You are wrong!–automatic detection of interaction errors from brain waves," in *Proceedings of the 19th International Joint Conference on Artificial Intelligence*, (CONF).

Ferrez, P. W., and Millán, J. d. R. (2008). Error-related EEG potentials generated during simulated brain-computer interaction. *IEEE Trans. Biomed. Eng.* 55, 923–929. doi: 10.1109/TBME.2007.908083

Gramfort, A., Luessi, M., Larson, E., Engemann, D. A., Strohmeier, D., Brodbeck, C., et al. (2013). MEG and EEG data analysis with MNE-Python. *Front. Neurosci.* 7:267. doi: 10.3389/fnins.2013.00267

Haotian, X., Anmin, G., Jiangong, L., Fan, W., Peng, D., and Yunfa, F. (2023). Online adaptive classification system for brain-computer interface based on error-related potentials and neurofeedback. *Biomed. Signal Process. Control* 82:104554. doi: 10.1016/j.bspc.2022.104554

Kim, S. K., Kirchner, E. A., Stefes, A., and Kirchner, F. (2017). Intrinsic interactive reinforcement learning - Using error-related potentials for real world human-robot interaction. *Sci. Rep.* 7:17562. doi: 10.1038/s41598-017-17682-7

Kreilinger, A., Hiebel, H., and Muller-Putz, G. R. (2016). Single versus multiple events error potential detection in a BCI-controlled car game with continuous and discrete feedback. *IEEE Trans. Biomed. Eng.* 63, 519–529. doi: 10.1109/TBME.2015.2465866

Krigolson, O. E., Holroyd, C. B., Van Gyn, G., and Heath, M. (2008). Electroencephalographic correlates of target and outcome errors. *Exper. Brain Res.* 190, 401–411. doi: 10.1007/s00221-008-1482-x

Kumar, A., Gao, L., Pirogova, E., and Fang, Q. (2019). A review of error-related potential-based brain-computer interfaces for motor impaired people. *IEEE Access* 7, 142451–142466. doi: 10.1109/ACCESS.2019.2944067

Lee, H. K., and Choi, Y.-S. (2019). Application of continuous wavelet transform and convolutional neural network in decoding motor imagery brain-computer interface. *Entropy* 21:1199. doi: 10.3390/e21121199

Leeb, R., Lee, F., Keinrath, C., Scherer, R., Bischof, H., and Pfurtscheller, G. (2007). Brain-computer communication: motivation, aim, and impact of exploring a virtual apartment. *IEEE Trans. Neural Syst. Rehabilit. Eng.* 15, 473–482. doi: 10.1109/TNSRE.2007.906956

Li, L., Chu, W., Langford, J., and Schapire, R. E. (2010). "A contextual-bandit approach to personalized news article recommendation," in *Proceedings of the 19th International Conference on World Wide Web*, 661–670. doi: 10.1145/1772690.1772758

Llera, A., Van Gerven, M., Gómez, V., Jensen, O., and Kappen, H. (2011). On the use of interaction error potentials for adaptive brain computer interfaces. *Neural Netw.* 24, 1120–1127. doi: 10.1016/j.neunet.2011.05.006

Miltner, W. H., Braun, C. H., and Coles, M. G. (1997). Event-related brain potentials following incorrect feedback in a time-estimation task: evidence for a "generic" neural system for error detection. *J. Cogn. Neurosci.* 9, 788–798. doi: 10.1162/jocn.1997.9.6.788

Mousavi, M., Koerner, A. S., Zhang, Q., Noh, E., and de Sa, V. R. (2017). Improving motor imagery BCI with user response to feedback. *Brain-Comput. Interf.* 4, 74–86. doi: 10.1080/2326263X.2017.1303253

Mousavi, M., Krol, L. R., and De Sa, V. R. (2020). Hybrid brain-computer interface with motor imagery and error-related brain activity. *J. Neural Eng.* 17:056041. doi: 10.1088/1741-2552/abaa9d

Schiatti, L., Barresi, G., Tessadori, J., King, L. C., and Mattos, L. (2019). "The effect of vibrotactile feedback on ERRP-based adaptive classification of motor imagery," in *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)* (Berlin, Germany: IEEE), 6750–6753. doi: 10.1109/EMBC.2019.8857192

Škola, F., Tinková, S., and Liarokapis, F. (2019). Progressive Training for Motor Imagery Brain-Computer Interfaces Using Gamification and Virtual Reality Embodiment. *Front. Hum. Neurosci.* 13:329. doi: 10.3389/fnhum.2019.00329

Soekadar, S. R., Birbaumer, N., Slutzky, M. W., and Cohen, L. G. (2015). Brain-machine interfaces in neurorehabilitation of stroke. *Neurobiol. Dis.* 83, 172–179. doi: 10.1016/j.nbd.2014.11.025

Spüler, M., and Niethammer, C. (2015). Error-related potentials during continuous feedback: using EEG to detect errors of different type and severity. *Front. Hum. Neurosci.* 9:155. doi: 10.3389/fnhum.2015.00155

Spüler, M., Rosenstiel, W., and Bogdan, M. (2012). Online adaptation of a c-VEP brain-computer interface (BCI) based on error-related potentials and unsupervised learning. *PLoS ONE* 7:e51077. doi: 10.1371/journal.pone.0051077

Sutton, R. S., and Barto, A. G. (2018). *Reinforcement Learning: An Introduction.* Cambridge: MIT press.

Tangermann, M., Müller, K.-R., Aertsen, A., Birbaumer, N., Braun, C., Brunner, C., et al. (2012). Review of the BCI competition IV. *Front. Neurosci.* 6:55. doi: 10.3389/fnins.2012.00055

Tao, T., Jia, Y., Xu, G., Liang, R., Zhang, Q., Chen, L., et al. (2023). Enhancement of motor imagery training efficiency by an online adaptive training paradigm integrated with error related potential. *J. Neural Eng.* 20:016029. doi: 10.1088/1741-2552/acb102

van Schie, H. T., Mars, R. B., Coles, M. G., and Bekkering, H. (2004). Modulation of activity in medial frontal and motor cortices during error observation. *Nat. Neurosci.* 7, 549–554. doi: 10.1038/nn1239

Wang, J., Wang, W., Su, J., Wang, Y., and Hou, Z.-G. (2024). Toward a fast and robust MI-BCI: online adaptation of stimulus paradigm and classification model. *IEEE Trans. Instrum. Meas.* 73, 1–12. doi: 10.1109/TIM.2024.3488147

Williams, S. E., Cumming, J., Ntoumanis, N., Nordin-Bates, S. M., Ramsey, R., and Hall, C. (2012). Further validation and development of the movement imagery questionnaire. *J. Sport Exerc. Psychol.* 34, 621–646. doi: 10.1123/jsep.34.5.621

Xavier Fidêncio, A., Klaes, C., and Iossifidis, I. (2022). Error-related potentials in reinforcement learning-based brain-machine interfaces. *Front. Hum. Neurosci.* 16:806517. doi: 10.3389/fnhum.2022.806517

Xavier Fidêncio, A., Klaes, C., and Iossifidis, I. (2024). A generic error-related potential classifier based on simulated subjects. *Front. Hum. Neurosci.* 18:1390714. doi: 10.3389/fnhum.2024.1390714

Zhou, D., Li, L., and Gu, Q. (2020). "Neural contextual bandits with UCB-based exploration," in III, H. D. and Singh, A., editors, *Proceedings of the 37th International Conference on Machine Learning, volume 119 of Proceedings of Machine Learning Research* (PMLR), 11492–11502.

# Appendix



**FIGURE 9**
Event-related spectral perturbation (ERSP) for subject $S12$ at channels C3 and C4 for both left- and right-hand motor imagery tasks. The color bars show the color and power spectral density in dB. For each channel, we used EEGLAB to compare the two experimental conditions (**left** vs. **right**) and show in the right-most plot the permutation results (800 permutations, using false-discovery rate correction for multiple comparisons, significant $p$-values shown in red for $\alpha = 0.05$). These plots highlight how some motor-imagery-related spectral modulations could be measured for this subject, even though they obtained the lowest score on the motor imagery ability compared to other subjects in the same study.



**FIGURE 10**
Event-related spectral perturbation (ERSP) for subject $S12$ at channels C3 and C4 for both left- and right-hand motor imagery tasks. The color bars show the color and power spectral density in dB. For each channel, we used EEGLAB to compare the two experimental conditions (**left** vs. **right**) and show in the right-most plot the permutation results (800 permutations, using false-discovery rate correction for multiple comparisons, significant $p$-values shown in red for $\alpha = 0.05$). These plots highlight how no significant motor-imagery-related spectral modulations could be measured for this subject, even though they obtained the highest score on the motor imagery ability compared to other subjects in the same study.