

PURSUING THE IMPLEMENTATION OF A NEUROTUTOR: AN EEG-BASED CLASSIFICATION OF READING TYPES

H. Romero-Morales¹, J. N. Muñoz-Montes de Oca¹, A.A. Torres-García¹, L. Villaseñor-Pineda¹

¹ Biosignals Processing and Medical Computing Laboratory, Instituto Nacional de Astrofísica, Óptica y Electrónica, Puebla, Mexico

E-mail: jenny.munoz@inaoep.mx

ABSTRACT: Electroencephalogram (EEG)-based brain-computer interfaces (BCI) emerged as systems to aid impaired people in daily life. Nowadays, the number of applications and target users of BCI has increased, including those for education purposes. An example of these applications, called neurotutor, was posited in 2015 for improving students' learning process. As a first step towards developing a neurotutor, we analyzed the EEG responses related to two types of reading. Specifically, this work assessed whether a machine learning algorithm can distinguish accurately between both classes from features obtained from the signals using one of three wavelet-based techniques. Also, the impact of epoch length on classifier performance was assessed. The method performance was analyzed under two scenarios (intra-subject and inter-subject), outperforming previous work. The best average accuracies were $94.40 \pm 5.10\%$ and $54.40 \pm 6.7\%$ for intra-subject and inter-subject classification, respectively. Although the progress obtained for the intra-subject scenario is promising, several steps must be done to effectively implement a neurotutor, especially in inter-subject scenarios.

INTRODUCTION

BCIs are systems that leverage the neurons' electrical activity, to generate an alternative channel that does not depend on muscular or verbal outputs. Some BCIs' applications are rehabilitation systems, videogames, neuromarketing, and recently in education [1].

On the other hand, Intelligent Tutoring Systems (ITS) are computer assistive systems designed to provide adaptive, personalized content for students [2]. In 2015, a novel ITS and BCI application named *Neurotutor* was elucidated aiming to enhance students' learning experiences and tailoring the content to individual needs. [3].

A fundamental skill in education is reading, which serves as a critical gateway to learning and academic development. Moreover, reading fosters cognitive skills, such as information processing and inference,

self-learning, and analytic thinking.

Previous works have attempted to analyze EEG signals related to the reading process [4, 5]. Also, in [6, 7] a database of EEG signals associated with two reading states was collected and processed with a baseline method to differentiate the reading states. [8] explored this database to characterize the normal reading paradigm looking for patterns of event-related potentials.

In this work, a first step toward the development of an EEG-based neurotutor was made, by analyzing and processing the dataset collected in [6] to distinguish between two normally employed reading strategies (types), comprehension reading (NR), and Task-Specific Reading (TSR aka scanning). Particularly, this study evaluated the capability of machine learning algorithms to accurately differentiate between both reading tasks using three wavelet-based methods. Additionally, the impact of epoch-length on classification performance was assessed.

A *neurotutor* would benefit from assessing reading comprehension to adapt the contents based on the readability of a text. During NR, the student focuses on deriving information about the central themes of the text and drawing inferences. Whereas TSR is a reading strategy in which the reader focuses on specific information (keywords). TSR is usually employed as a pre-reading strategy, in which the user can decide whether a text provides relevant information for the task in question, or after reading, to locate segments of interest. Therefore, recognizing when students engage in one strategy of reading could guide the neurotutor, leading it to adjust to contents that promote deeper comprehension, all while trying to maintain the engagement and motivation of the user.

MATERIALS AND METHODS

1. Dataset description and preprocessing

The Zurich Cognitive Language Processing Corpus 2.0 (ZuCo 2.0) is a dataset of two physiological signals, EEG and eye-tracker, of 18 English native-speaking subjects recorded during two different read-

ing tasks. In the experiment, subjects were asked to read sentences from an annotated Wikipedia corpus (in English), in which each sentence is associated with a specific semantic relation (i.e. Political affiliation, education, founder, wife/husband, job title, nationality, and employer). Participants were asked to read sentences with two different purposes, (1) to find an implicit relation in a text (Task Specific Reading (TSR), 390 sentences) and (2) to fully comprehend the meaning of the sentence (Normal Reading (NR), 349 sentences). To encourage reading comprehension during the NR task, some control questions were randomly presented after some instances. Additionally, each participant was required to do a linguistic assessment (Lexical Test for Advanced Learners of English) to measure their language proficiency [6].

In this work, only EEG signals were employed for classification, seeking to reduce the amount of data needed for the classification. Moreover, EEG signals are currently being researched to derive implicit information about the user's state (memory load, emotions, etc.), this information could aid in the development of a neurotutor. The signals were recorded using a 128-electrode Geodesic Hydrocel System with a sampling frequency of 500 Hz. EEG signals were preprocessed by [6], using Matlab's Automagic and the Multiple Artifact Rejection Algorithm (MARA). Twenty-three electrodes were removed during this stage because of their predominant muscular and ocular information. Furthermore, to reduce data size and computational load, signals were downsampled to 256 Hz.

2. Epoch extraction

Since natural and untimed reading was encouraged during both task reading, the duration of reading epochs was variable. To standardize the length of the signals of all participants, epochs of 1, 2, and 3 seconds were selected to analyze the significance of epoch length in classification performance, full-length signals were also analyzed to establish a reference. Records of at least 3 seconds were selected for this analysis, which created a class imbalance; thus a random selection of 100 epochs (samples) per class was applied. Participant YDR was excluded from this study because of insufficient epochs per class. Epochs were extracted from the center region of each sentence, which means the 1-second epoch is contained in the other two epochs, and the 2-second epoch is contained in the 3-second epoch. The whole epoch was also analyzed to obtain a benchmark and evaluate the significance of using a shorter epoch for analysis.

3. Channel selection

After the preprocessing stage, a set of 105 channels were kept. Given the inherent characteristics of EEG, certain channels exhibit redundant information; for this reason and aiming to reduce classi-

fication times, our experiments were focused on 31 electrodes taken from a standard 32-electrode setup.

4. Feature extraction

Wavelet-based methods have been proven to be accurate techniques to characterize and process biomedical signals [9], which due to their complexity and variability tend to be hard to analyze. In this work, three wavelet-based techniques were analyzed to find the best characterization of the EEG signals related to reading tasks.

4.1 Discrete Wavelet Transform (DWT)

This method decomposes the signal using a series of filters. The filtering process is limited by the sampling frequency and the length of the signal. DWT provides n -levels of decomposition, by dividing the signal into a high-frequency component (detail coefficients) and a low-frequency component (approximation coefficient). A second-order Daubechies is used as the mother wavelet, with 6 levels of decomposition. A 2nd-order Daubechies is chosen because of the similarity between the wavelet and EEG patterns; moreover, it has been used successfully to classify EEG signals for seizure detection in epilepsy [9]. The EEG signal, located within 0.5 to 50 Hz, has often been characterized in terms of five brain rhythms: delta (0.5 - 4 Hz), theta (4 - 8 Hz), alpha (8 - 14 Hz), beta (14 - 30 Hz), and gamma waves (higher than 30 Hz). DWT analyzes the signal in the time-frequency domain by decomposing it into sub-bands. Given a 256 Hz sampling frequency, DWT efficiently matches these frequency bands, enabling the extraction of characteristics pertinent to cognitive tasks. Then for each level of decomposition, eleven features were calculated: mean, root-mean square (RMS), kurtosis, median, maximum and minimum amplitude, standard deviation, energy, Instantaneous Wavelet Energy (IWE), Teager Wavelet Energy (TWE), and Hierarchical Wavelet Energy (HWE). In total, 77 features were computed for each channel, resulting in 2387 features per epoch.

4.2 Continuous Wavelet Transform (CWT)

Continuous Wavelet Transform highlights the intricate relationship among the frequency, time, and energy of a signal, through a visual representation known as the 'scalogram'. In this study, CWT of each EEG channel was computed using an Analytic Morlet wavelet as the mother wavelet. The scalogram was then divided into the EEG bands described in Section 4.1: delta, theta, alpha, beta, and gamma. From each band, a comprehensive set of 29 characteristics was extracted. These features include the flux at 0, 45, and 90 degrees, as well as the energy of the scalogram, which reflects amplitude variations across the frequency and time axes of the scalogram. Additionally the RMS, mean, standard deviation, skewness, kurtosis, maximum value, entropy, and three key percentiles (75th, 50th, and 25th) were computed. Furthermore, an entropy filter

from Matlab's Image Processing Toolbox was applied to each EEG band, this filter computes the entropy across the image, highlighting the dynamic changes and complexity of the entropy measure within the CWT coefficients. For each entropy-filtered scalogram, features such as the mean, standard deviation, RMS, skewness, kurtosis, and the three percentiles were computed. Subsequently, the average waveform was derived by calculating the mean at each time point across the frequency spectrum of the segmented scalogram. Seven features were then extracted from this averaged signal: mean, median, standard deviation, kurtosis, skewness, RMS, and sample entropy. As a result, 145 characteristics were computed for each scalogram, given that 31 channels were employed within the study, a total of 4495 features were computed per sample.

4.3 Wavelet Scattering Transform

Wavelet Scattering Transform (WST) is a novel wavelet-based method used for the analysis of time series that exhibit non-linear and non-stationary characteristics, such as EEG signals. This advanced mathematical technique yields sparse representations that are invariant to translations and stable to deformations.

In the first level of WST, a decomposition produces a series of coefficients at different scales. A modulus operation is then applied to these coefficients to capture the signal's energy across various frequencies. The resulting modulus wavelet coefficients are subsequently averaged, yielding translation-invariant features of the signal. This operation is recursive yielding higher order coefficients. Typically, first and second-order coefficients capture the majority of relevant frequency information of naturally occurring phenomena. These features are unique to the scattering transform's framework and serve as the foundation for its powerful signal analysis capabilities.

In this work, the scattering time-invariant first-order coefficients are divided into five segments related to brain rhythms (delta, theta, alpha, beta, and gamma). Similar to CWT, WST yields a visual representation often referred to as 'scattergram', which relates time, frequency, and power information of the EEG signal. For each EEG-Band derived coefficients, the 29 descriptors described in Section 4.2 were computed, resulting in a total of 4995 features per epoch.

All three wavelet-based methods were applied to the three epoch lengths (1, 2, and 3 seconds) and for the complete signal.

5. Classification

This study aimed to compare the wavelet-based methods for the classification of EEG signals obtained during two types of reading. Given that EEG signals are highly variably across subjects and even across sessions, an intra-subject approach was pursued, to validate the discrimination power of the

proposed method. Nevertheless, training individualized models requires gathering extensive data, which can be time-consuming, so an evaluation of an inter-subject classification scheme was attempted.

5.1 Intra-subject Classification

For intra-subject classification (IAC), a model was trained for each subject in the dataset. The model was evaluated using a 5-fold validation, with 40 samples (20 of each class) per fold for the testing stage. Three classification algorithms were tested: Support Vector Machine (SVM with a quadratic kernel), K-Nearest Neighbors (KNN with 5 neighbors), and Random Forest (RF with 100 trees). A total of 36 classifiers were trained per subject due to the lengths of the four epochs, three feature-extraction methods, and three classification algorithms were compared.

5.2 Inter-subject Classification

In addition to IAC, an inter-subject classifier (IEC) was trained, using a leave-one-subject-out cross-validation. Prior knowledge of the best epoch duration and feature extraction method was inferred from IAC. 2-second and 3-second epochs were analyzed, using DWT-based features and RF.

RESULTS AND DISCUSSION

The experiments were carried out to evaluate whether machine learning algorithms, trained on time-frequency representations of EEG signals and with different epoch lengths, could detect differences in brain patterns from subjects engaging in two types of reading: TSR and NR. Furthermore, this rationale was analyzed in two scenarios of classification: intra-subject (personalized models) and inter-subject (generalized models).

1. Intra-subject experiments

Figure 1 shows an analysis of epoch length and its impact on classification performance. In this Figure the global average accuracy for all subjects is taken, regardless of the classification algorithm used, primarily to determine if the length of the signal affects classification outcomes. A trend is observed across all wavelet-based methods; as the epoch size is increased, the performance of the classifier is enhanced. This does not hold when analyzing full-length signals in DWT-derived features. Last, for this and the remaining figures the chance level (50% for two balanced classes) is shown as a dashed line. A non-parametric, Kruskal-Wallis test, with a post-hoc follow-up Dunn's test was performed for each wavelet feature group. Significant differences were found between the one-second epoch and the 3-second epochs in all characterizations. Given the trend observed, and the reduced computational costs in epoch analysis, the two-second and 3-second epochs were further analyzed as promising for TSR and NR classification. A comparative analysis of machine learning algorithms (KNN, SVM and RF) was conducted for

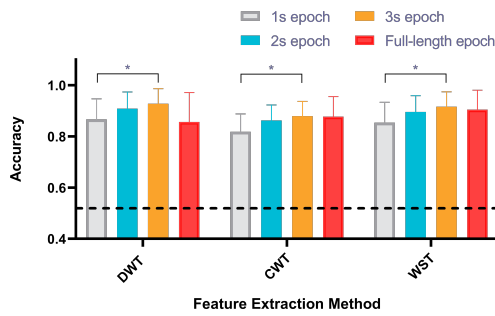


Figure 1: Comparing how epoch length affects classification using wavelet-based features, average accuracy from IAC is shown across KNN, SVM, and RF classifiers. Regardless of classifier type, accuracies are averaged to emphasize epoch length’s impact. Significant differences between epoch lengths are marked with asterisks ($p < 0.05$, Dunn’s test).

all feature extraction methods, using the 3-second epoch, which achieved the highest averaged performance. Figure 2 presents a comparison of the classification algorithm and its effect on global accuracy across all subjects. Significant differences were found in CWT-derived methods, for KNN and RF, as well as for KNN and SVM algorithms. Using WST, differences were found between KNN and RF. On the other hand, the best performances were obtained for DWT and WST regardless of the machine learning algorithm used.

For both wavelet-based methods, the best classification performance was yielded using RF with an accuracy of $94.40 \pm 5.10\%$, DWT and $94.80 \pm 5.10\%$, WST). Furthermore, DWT-based descriptors exhibited consistent performance across all three machine learning algorithms. DWT has positioned itself as a valuable tool for EEG classification since it provides a time-frequency analysis without information loss or alteration while reducing computational costs [9]. A baseline classifier was trained using the characteristics proposed by [7], and evaluated through the same classification scheme. Baseline characteristics were obtained by filtering each epoch into the relevant EEG bands and obtaining the mean amplitude from each EEG component. To ensure comparable results, the 31 channels selected in this study were also used for benchmark classification.

Even though no significant differences were found within the proposed features and the benchmark, our approach utilizes epochs of 2 and 3 seconds, meanwhile, the average signal length from the original recordings is 5.84 seconds for NR and 4.81 seconds for TSR[6]. Short epochs reduce computational costs and would be more suitable for online applications. Figure 3 shows accuracies obtained for each subject, employing both the DWT-derived features with RF classifier and 3-second epochs. The best accuracy (i.e. 99.50 ± 1.11) was achieved by subjects YAK and YMS. Besides that, all subject accuracies were

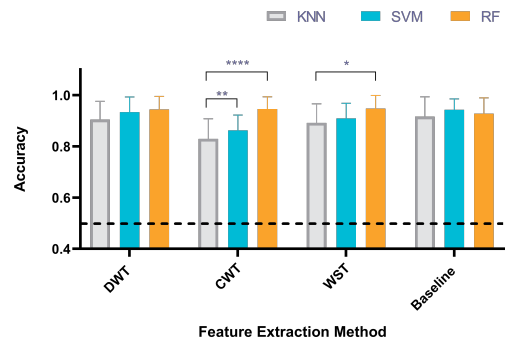


Figure 2: Intra-subject performances compared between wavelet-based and benchmark features using three-second epochs (optimal length). Asterisks indicate significant differences ($p < 0.05$, Dunn’s test).

greater than 85%. Additionally, the scores provided in [6] for NR control questions (NR scores), and correct semantic text identification (TSR scores) were analyzed to see if task classification was correlated to individual performance in each task. No correlation was found between classification performance and NR/TSR scores (Spearman test). However, a low performance across classifiers and epochs was consistently observed for YAG, achieving one of the lowest performances (i.e. 87.70 ± 4.67). This outcome could be explained because YAG also exhibited a low performance on semantic identification. On the other hand, subjects such as YRK, YLS, YMD, YMS, exhibited great performance in both control tests and similarly an accurate classification in the proposed methodology. Although YAK received the lowest score for the set of random questions in NR, the algorithm demonstrated good performance. Since questions were randomly presented, for NR scores it is difficult to assess if scores truly reflect the quality of the task being performed by the user.

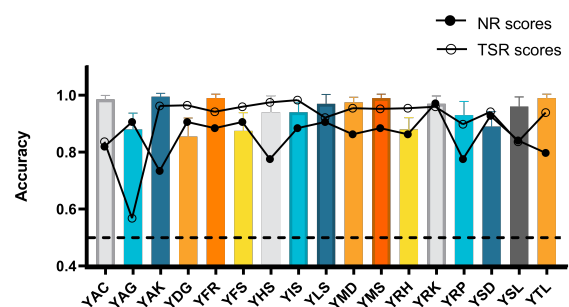


Figure 3: Intra-subject accuracies by RF after 5-fold cross-validation using DWT-based features and 3-second epochs. Also, NR and TSR scores are shown [6]. NR scores show the accuracy of responses to randomly posed comprehension questions, while TSR scores refer to the correct identification of semantic relations within the text.

2. Inter-subject experiments

Despite constructing a universal model to interpret highly variable EEG signals being a complex challenge, an experiment for inter-subject classification was also performed on the data, using a leave-one-subject-out validation. For this experiment, only two and 3-second epochs were analyzed, using the most stable characterization (DWT). Additionally, a deep-learning architecture specialized in EEG data (EEGNet) was trained for the identification of NR and TSR, for the 2 and 3 second epochs. The hyperparameters for the network were selected based on the recommendations of the EEGNet developers (128 kernel length) [10]. Also, some parameters were taken from [4] such as Adam optimizer and a batch size of 16 instances. Since that work analyzed EEG signals related to the reading process and obtained promising results. Likewise, 100 iterations were calculated.

A baseline classifier was trained using the characteristics used in [7] and evaluated through the same classification algorithms. Benchmark features were originally calculated using full-length signals. All the models were trained using data from the 31 selected channels. Table 1 presents the performance for all trained models.

The best classification performance, 54.4 ± 6.7 , was achieved by DWT-derived features, using a 2-s epoch. Despite baseline features achieving a global classification accuracy lower than the chance level, all proposed methods and EEGNet slightly outperformed the chance level for the two classes. This could imply different cognitive processes are undertaken in both reading tasks that could be generalized efficiently across subjects through the proposed methodology and EEGNet. Despite no classifier got a higher average accuracy than the empirical chance level for 100 trials per class (58%) [11]; this threshold was overcome for eight subjects.

Table 1: Performance comparison between DWT-based models, EEGNet and baseline descriptors for inter-subject classification.

Method	Average	Median	Max - Min
2-s epoch	54.4 ± 6.7	53.5	68.0 - 45.0
3-s epoch	51.7 ± 7.2	52.0	66.0 - 37.0
EEGNet 2-s	52.9 ± 8.8	52.0	74.5 - 38.5
EEGNet 3-s	48.06 ± 11.54	50.0	66.5 - 27.0
Baseline	49.5 ± 12.0	46.5	80.0 - 31.5

Figure 4 shows the individual test accuracy obtained after leave-one-subject-out validation. Both the proposed methodology and EEGNet implementation generally achieved accuracies surpassing the theoretical random classifier. Specifically, EEGNet showed an accuracy above the 50% threshold for approximately 8 subjects, whereas the proposed method achieved this for 12 out of 17 subjects. For the empirical random classifier, 5 subjects surpassed the

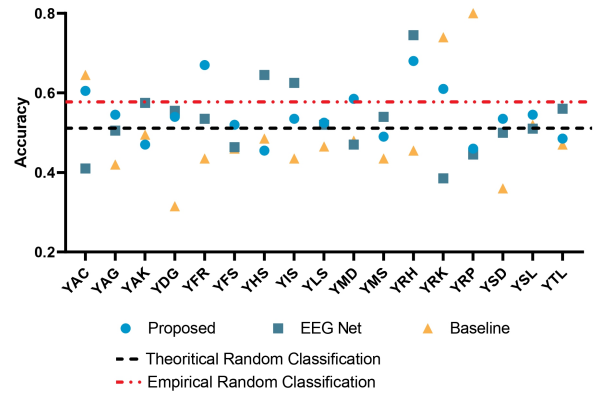


Figure 4: Test accuracy for each subject after leave-one-subject-out validation is shown for three methods: (1) Proposed method using a 2-second epoch, DWT-derived features and RF; (2) EEGNet classifier, and (3) Baseline method trained with the features from [7] and RF. The theoretical and empirical [11] chance levels are indicated by black and red dashed lines, respectively.

58% threshold. This suggests time-frequency features, along with RF might be useful to discern between reading tasks. Moreover, when compared to the baseline classifier leveraging features proposed by [7], this study showed the majority of the participants (9) obtained the lowest classification accuracy, while 13 subjects did not surpass the random classifier. Therefore, baseline characteristics do not seem suitable for the task (even though the overall best accuracy was obtained using them, through subject YRP).

Interestingly, subjects who under-perform in the intra-subject approach (YAG, YDG, and YFS), have similar low results for the inter-subject classifier, obtaining results near random classification (Proposed, EEGNet), or below it (Baseline). Similarly, subjects YAC, YFR, YRK, YMD, overperform both in intra-subject and intersubject analysis.

Reading is a complex task, that requires the activation of various brain sub-processes; beyond language processing and visual decoding, reading evokes responses from attention, working memory load, abstract reasoning, and memory pathways. Consequently, we hypothesized NR and TSR could be differentiated through EEG patterns since cognitive and attention demands are different in each reading strategy. Within this study, DWT-based features allowed the distinction of the two types of reading using EEG data. These differences, represented through the time-frequency domain, result from specific brain processes, such as attention or cognitive load. Additionally, inter-subject classification rates above the random classifier reflect subtle patterns that are generalized across subjects. Research by Hollenstein et al. furtherly support this idea, using eye-tracking. In their study, results indicated subjects focused uniformly on each word when engaging

in NR, on the other hand, in TSR subjects focused on words that determined a specific semantic relation and skimmed through the rest. Reading times were also reduced while engaging in TSR [7]. Despite this, more studies are required to identify features that allow classification, based on cognitive differences.

An ideal neurotutor, would benefit from the detection of different reading states to provide active feedback about the psycho-emotional state of the user. Cognitive overload and academic stress can impair students' well-being and decrease academic performance. Likewise, professors could use this as a tool for monitoring students' progress or evaluating contents in concordance with students' comprehension. Although scanning or TSR is a common reading strategy, useful for navigating through large amounts of information, their extensive use could result in surface level comprehension [12]. A neurotutor could adapt the contents' readability and encourage reading comprehension through different activities, thus helping to reduce the use of scanning and pursuing deeper comprehension in the user.

CONCLUSIONS

The first approach to NR and TSR classification involved an IAC in which an analysis of epoch length, feature extraction method, and classification algorithm were assessed. Best accuracy (i.e. 94.40 ± 5.10) was achieved using 2s and 3s epochs, DWT-derived features, and RF classifier. All subjects performed beyond 85% accuracy. Furthermore, for some cases, a relation was found between classifier performance and control test scores, which could imply that diminished performance in task completion reduces classification outcomes.

IEC proved to be a complex task; nevertheless, both the proposed method (2s epoch, DWT, RF) and EEGNet performed above the random and benchmark classifier. From this, it could be inferred there are brain patterns shared across subjects when performing reading tasks, namely NR and TSR.

Future works will explore the relationship between classification performance and control questions, to determine if removing low-performing subjects could increase classifier performance. Further, techniques of data augmentation could be employed to gain more insight into differences in the reading patterns. Likewise, changes in the inclusion criteria for epoch selection could be performed, even if it results in unbalanced classes. Inter-subject complexity derives from the intricate nature of EEG signals, since they are variable across subjects and even across sessions. A generalized model, enhanced by a limited number of training samples from the new user, could improve classification accuracy while maintaining the benefits of limited training time.

ACKNOWLEDGMENTS

The authors, H. R-M and J. M-MdO, wish to express their gratitude for the graduate scholarships granted by CONAHCYT, Mexico. These grants have enabled them to conduct the research presented in this work.

REFERENCES

- [1] Aricò P, Borghini G, Di Flumeri G, Sciaraffa N, Babiloni F. Passive BCI beyond the lab: Current trends and future directions. *Physiological measurement*. 2018;39(8):08TR02.
- [2] Lin CC, Huang AY, Lu OH. Artificial intelligence in intelligent tutoring systems toward sustainable education: A systematic review. *Smart Learning Environments*. 2023;10(1):41.
- [3] Müller-Putz G *et al*. The future in brain/neural computer interaction: *Horizon 2020*. 2015.
- [4] Torres-García AA, Martínez-Santiago F, Montejo-Ráez A, Ureña-López LA. Toward an educative EEG-based neuroIIR system for adapting contents. *International Journal of Human-Computer Interaction*. 2023:1–15.
- [5] Ye Z *et al*. Towards a better understanding of human reading comprehension with brain signals. In: *ACM Web Conference*. 2022, 380–391.
- [6] Hollenstein N, Troendle M, Zhang C, Langer N. ZuCo 2.0: A dataset of physiological recordings during natural reading and annotation. In: *12th LREC. ELRA: Marseille, France, May 2020*, 138–146.
- [7] Hollenstein N, Tröndle M, Plomecka M, Jäger LA, Langer N. The zuco benchmark on cross-subject reading task classification with eeg and eye-tracking data. *Frontiers in Psychology*. 2023;13:1028824.
- [8] Liu X, Cao Z. Enhance reading comprehension from eeg-based brain-computer interface. In: *Australasian Joint Conference on Artificial Intelligence*. 2023, 545–555.
- [9] Chen D, Wan S, Xiang J, Bao FS. A high-performance seizure detection algorithm based on Discrete Wavelet Transform and EEG. *PloS one*. 2017;12(3):e0173138.
- [10] Lawhern VJ, Solon AJ, Waytowich NR, Gordon SM, Hung CP, Lance BJ. Eegnet: A compact convolutional neural network for eeg-based brain-computer interfaces. *Journal of neural engineering*. 2018;15(5):056013.
- [11] Combrisson E, Jerbi K. Exceeding chance level by chance: The caveat of theoretical chance levels in brain signal classification and statistical assessment of decoding accuracy. *Journal of neuroscience methods*. 2015;250:126–136.
- [12] Elleman AM, Oslund EL. Reading comprehension research: Implications for practice and policy. *Policy Insights from the Behavioral and Brain Sciences*. 2019;6(1):3–11.