# Reading Task Classification Using EEG and Eye-Tracking Data

Nora Hollenstein[*a], Marius Tröndle[b], Martyna Plomecka[b], Samuel Kiegeland[c], Yilmazcan Özyurt[c], Lena A. Jäger[d,e], Nicolas Langer[b]

[a]*Center for Language Technology, University of Copenhagen*
[b]*Department of Psychology, University of Zurich*
[c]*Department of Computer Science, ETH Zurich*
[d]*Department of Computational Linguistics, University of Zurich*
[e]*Department of Computer Sciene, University of Potsdam*

## Abstract

The Zurich Cognitive Language Processing Corpus (ZuCo) provides eye-tracking and EEG signals from two reading paradigms, normal reading and task-specific reading. We analyze whether machine learning methods are able to classify these two tasks using eye-tracking and EEG features. We implement models with aggregated sentence-level features as well as fine-grained word-level features. We test the models in within-subject and cross-subject evaluation scenarios. All models are tested on the ZuCo 1.0 and ZuCo 2.0 data subsets, which are characterized by differing recording procedures and thus allow for different levels of generalizability. Finally, we provide a series of control experiments to analyze the results in more detail.

## 1. Introduction

### 1.1. Motivation & Background

Electroencephalographic (EEG) and eye tracking are considered gold-standard physiological and behavioral measures of cognitive processes involved in reading (Rayner, 1998; Dimigen et al., 2011). Reading task classification, i.e., decoding mental states and detecting specific cognitive processes occurring in the brain during different reading tasks, is an important challenge in cognitive neuroscience as well as in natural language processing.

---

[*]Corresponding author: *nora.hollenstein@hum.ku.dk*

Reading is a complex cognitive process that requires the simultaneous processing of complex visual input, as well as syntactic and semantic integration. Identifying specific reading patterns can improve models of human reading and provide insights into human language understanding and how we perform linguistic tasks. This knowledge can then be applied to machine learning algorithms for natural language processing. Accurate reading task classification can improve the manual labelling process for a variety of NLP tasks, as these processes are closely related to identifying reading intents. Recognizing reading patterns for estimating reading effort has additional applications such as as the diagnosis of reading impairments such as dyslexia (Rello and Ballesteros, 2015; Raatikainen et al., 2021) and attention deficit disorder (Tor et al., 2021).

One of the main bottlenecks in training supervised natural language processing (NLP) and machine learning (ML) applications is that large labeled datasets are required. Generating these labels is often still an expensive and time-consuming manual process. The Zurich Cognitive Language Processing Corpus (ZuCo; Hollenstein et al. 2018, 2020) addresses these challenges. ZuCo is a dataset combining electroencephalography (EEG) and eye-tracking recordings from subjects reading natural English sentences. The EEG and eye-tracking signals lend themselves to train improved ML models for various tasks, in particular for information extraction tasks. One of the advantages of the ZuCo dataset is that it provides ground truth labels for additional machine learning tasks. The availability of labelled data plays a crucial role in all supervised machine learning applications. Physiological data can be used to understand and improve the labelling process (e.g., Tokunaga et al. 2017), and, for instance, to build cost models for active learning scenarios (Tomanek et al., 2010). Is it possible to replace this expensive manual work with models trained on physiological activity data recorded from humans while reading? That is to say, can we find and extract the relevant aspects of text understanding and annotation directly from the source, i.e., eye-tracking and brain activity signals during reading? Using cognitive signals of language processing could be used directly to (pre-)annotate samples to generate training data for ML models. Moreover, identifying reading intents can help to improve the labelling processes by detecting tiredness from brain activity data and eye-tracking data, and subsequently to suggest breaks or task switching.

We leverage the ZuCo dataset with EEG and eye-tracking recordings from reading English sentences for this work. This dataset contains two different reading paradigms: Normal reading (with the only task of reading naturally

for reading comprehension) and task-specific reading (with the purpose of finding specific information in the text). We train machine learning models on eye-tracking and EEG features to solve a binary classification to identify the two reading tasks as accurately as possible. We investigate how two different reading tasks affect both eye movements and brain activity.

Understanding the physiological aspects of the reading process can advance our understanding of human language processing as well as provide benefits for natural language processing. Recent advances in machine learning are providing new methods to approach reading task classification (Haynes and Rees, 2006; Mathur et al., 2021). Moreover, the availability of a neurolinguistic dataset with co-registered EEG and eye-tracking signals such as the ZuCo facilitates this work. The simultaneous recording of EEG and eye-tracking allows us to investigate specific feature sets on different levels of analysis, e.g., sentence level, word level, fixation level. Additionally, the naturalistic setup of the experiments used in this work are crucial for this work and for the ecological validity of experiment in neuroscience in general (Nastase et al., 2020).

### 1.2. Previous Work

The ZuCo dataset is freely available and has recently been used in a variety of applications including leveraging EEG and eye-tracking data to improve natural language processing tasks (Barrett et al., 2018; Mathias et al., 2020; McGuire and Tomuro, 2021), evaluating the cognitive plausibility of computational language models (Hollenstein et al., 2019; Hollenstein and Beinborn, 2021), investigating the neural dynamics of reading (Pfeiffer et al., 2020), developing models of human reading (Bautista and Naval, 2020; Bestgen, 2021). This shows that ZuCo can also be used for other machine learning benchmark tasks. In a recent study, the ZuCo data has been used already for reading task identification (Mathur et al., 2021). The authors propose a complex convolutional network combining eye-tracking and EEG features, which is evaluated on a fixed cross-subject scenario (trained on 12 subjects, validated on 2, and tested on 4 subjects) on the sentences from ZuCo 2.0. The authors achieve 69.79% accuracy on this binary classification task. However, this performance measured on a fixed evaluation setting still leaves room for improvement and open research questions regarding the selection of features.

*1.3. Contributions*

We propose a machine learning approach for reading task classification based on eye-tracking and EEG data. We investigate a large range of features and implement word-level and sentence-level models. We test all models on both ZuCo datasets. Finally, we present a series of control analyses to validate the results. The code for all experiments is available online.[1]

The results show substantially higher performance for sentence-level than for word-level models. Generally, we find that the models achieve high accuracy on the ZuCo 1.0 data, and lower accuracy – but still higher than the baselines – for the ZuCo 2.0 data. Additional analyses show that these differences in performance might be attributed to the session effects present in ZuCo 1.0. Moreover, while the within-subject evaluation yields good results, there is still room for improvement in the cross-subject settings in future work, which is crucial for practical machine learning applications.

## 2. The Zurich Cognitive Language Processing Corpus

In this section, we describe the compilation of the Zurich Cognitive Language Processing Corpus (ZuCo). ZuCo is a dataset combining electroencephalography (EEG) and eye-tracking recordings from subjects reading natural sentences. ZuCo includes high-density EEG and eye-tracking data of 30 healthy adult native English speakers, each reading natural English text for 3–6 hours. We recorded two separate datasets with different participants. The first dataset, ZuCo 1.0, encompasses EEG and eye-tracking data of 21,629 words in 1107 sentences for each of the 12 subjects. The second dataset, ZuCo 2.0, encompasses the same type of recordings of 15,138 words and 739 sentences for each of the 18 subjects. The recordings of ZuCo 1.0 include three reading paradigms. In this work, we consider two paradigms only (present also in ZuCo 2.0): a normal reading experiment and a task-specific reading experiment.

Both datasets, including the raw data and the extracted features, are freely available on the Open Science Framework[2]. Moreover, both datasets have been extensively described in previous publications (ZuCo 1.0 in Hollenstein et al. 2018 and ZuCo 2.0 in Hollenstein et al. 2020). Therefore, in

---

[1]`https://github.com/norahollenstein/reading-task-classification`
[2]ZuCo 1.0: `https://osf.io/q3zws/` and ZuCo 2.0: `https://osf.io/2urht/`.
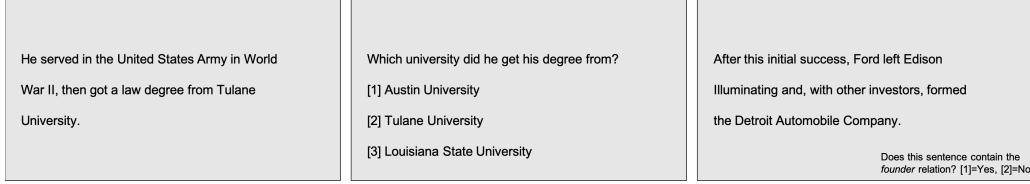
Figure 1: Example sentences on the recording screen: (left) a normal reading sentence, (middle) a control question for a normal reading sentence, and (right) a task-specific annotation sentence.

this article we provide a higher-level general description of the data collection focusing mainly on the reading task paradigms relevant for the benchmark task.

One of the main advantages of the ZuCo dataset is its naturalistic reading setup, defined by the following characteristics: (1) We present full sentences on the screen spanning multiple lines, as opposed to a rapid serial visual paradigm where each word is presented in isolation. (2) There are no time constraints on reading speed. The participants are able to read each sentence in their own pace. (3) The presented stimuli are naturally occurring sentences and not hand-picked or manually constructed for experimental purposes.

*2.1. Reading Materials & Experimental Design*

The reading materials recorded for the ZuCo corpus contain sentences from movie reviews from the Stanford Sentiment Treebank (Socher et al., 2013) and Wikipedia articles from a dataset provided by Culotta et al. (2006). These resources were chosen since they provide ground truth labels for various natural language processing ML tasks. In this work, we focus on the Wikipedia sentences, which were used in the normal reading (NR) and task-specific (TSR) experiment paradigms. Descriptive statistics about the datasets used in this work are presented in Table 1.

For the recording sessions, the sentences were presented one at a time at the same position on the screen. Text was presented in black with font size 20-point Arial on a light grey background resulting in a letter height of 0.8 mm or 0.674° of visual angle. The lines were triple-spaced, and the words double-spaced. A maximum of 80 letters or 13 words were presented per line in both tasks. Long sentences spanned multiple lines (max. 7 lines).

5

|            | ZuCo 1.0 | | ZuCo 2.0 | |
|            | NR | TSR | NR | TSR |
|---|---|---|---|---|
| sentences | 300 | | 349 | 390 |
| sent. length | 21.3 (±10.6) | 20.1. (±10.1) | 19.6 (8.8) | 21.3 (9.5) |
| total words | 6386 | 8164 | 6828 | 8310 |
| word types | 2657 | 2995 | 2412 | 2437 |
| word length | 6.7 (2.7) | 6.7 (2.6) | 4.9 (2.7) | 4.9 (2.7) |
| Flesch score | 51.33 | 51.43 | 55.38 | 50.76 |

Table 1: Descriptive statistics of reading materials (SD = standard deviation), including Flesch readibility scores.

## 2.2. Normal Reading (NR)

In the first task, participants were instructed to read the sentences naturally, without any specific task other than comprehension. The participants were instructed to read one sentence at a time at their own pace and use the control pad to move to trigger the onset of the next sentence. They were informed that a portion of the sentences would be followed by a comprehension question. The task was explained to the subjects orally, followed by instructions on the screen. Figure 1 (left) shows an example sentence as it was depicted on the screen during recording. As shown in Figure 1 (middle), the control condition for this task consisted of single-choice questions about the content of the previous sentence. The questions are presented with three answer options, out of which only one is correct. 12% of randomly selected sentences were followed by such a comprehension question with three answer options on a separate screen. The task was preceded by a practice round.

## 2.3. Task-specific Reading (TSR)

In the second task, the subjects were presented with similar sentences as in the normal reading task, but with specific instructions to search for a specific relation in each sentence they read. The following relation types were contained in the sentences: *award*, *education*, *employer*, *founder*, *job_title*, *nationality*, *political_affiliation*, *visited* and *wife/husband*. This allows us to compare the EEG and eye-tracking signals during normal reading to task-specific reading while searching for a specific relation type.

Instead of comprehension questions, the participants had to decide for each sentence whether it contained the relation or not, i.e. they were actively

annotating each sentence. Figure 1 (right) shows an example screen for this task. 17% of the sentences did not include the relation type and were used as control conditions. All sentences within one recording block involved the same relation type. The blocks started with a practice round, which described the relation and was followed by three sample sentences, so that the participants would be familiar with the respective relation type.

Purposefully, there are some duplicate sentences that appear in both the normal reading and the task-specific reading tasks (48 sentences in ZuCo 1.0 and 63 sentences in ZuCo 2.0.). The intention of these duplicate sentences is to provide a set of sentences read twice by all participants with a different task in mind. Hence, this enables the comparison of eye-tracking and brain activity data when reading normally and when annotating specific relations (see examples in Section 3.1).

## 2.4. Recording Procedure

The main difference and reason for recording ZuCo 2.0 consisted in the experiment procedure, namely, the number of sessions and the order of the reading tasks. For ZuCo 1.0, the normal reading and task-specific reading paradigms were recorded in different sessions on different days. The order of the sessions and sentences within the sessions was identical for all subjects.

Therefore, the recorded data is not fully appropriate as a means of comparison between natural reading and annotation, since the differences in the brain activity data might result mostly from the different sessions due to the sensitivity of EEG and session-specific effects in the eye-tracking signal. This, and extending the dataset with more sentences and more subjects, were the main factors for recording the ZuCo 2.0 dataset.

For ZuCo 2.0, we recorded 14 blocks of approx. 50 sentences for each subject, alternating between tasks: 50 sentences of normal reading, followed by 50 sentences of task-specific reading. The order of blocks and sentences within blocks was identical for all subjects. Each sentence block was preceded by a practice round of three sentences and followed by a short break to ensure a clear separation between the reading tasks.

The differing recording procedures between the two datasets allow us to investigate the impact of possible session biases in the data. As we show in Section 4, these aspects affect the results and are discussed through various control analyses in Section 5.

## 2.5. Participants

For ZuCo 1.0, data were recorded from 12 healthy adults (between 22 and 54 years, all right-handed; 5 female participants). For ZuCo 2.0, data were recorded from 18 healthy adults (between 23 and 52 years old; 2 left-handed; 10 female participants). The native language of all participants is English. See Appendix A for more details on participant demographics and linguistic assessment.

## 2.6. Technical Set-up & Preprocessing

*Eye-Tracking Acquisition.* Eye movements and pupil size were recorded with an infrared video-based eye tracker (EyeLink 1000 Plus, SR Research) at a sampling rate of 500 Hz. The participant was seated at a distance of 68cm from a 24-inch monitor (ASUS ROG, Swift PG248Q, display dimensions 531x299 mm, resolution 800x600 pixels resulting in a display: 400x298.9 mm, a vertical refresh rate of 100 Hz). The eye tracker was calibrated with a 9-point grid at the beginning of the session and re-validated before each block of sentences. The eye-tracker computed eye position data and identified events such as saccades, fixations, and blinks. Saccade onsets were detected using the eye-tracking software default settings: acceleration larger than 8000°/s2, a velocity above 30°/s, and a deflection above 0.1°.

*Eye-Tracking Feature Extraction.* The datasets provide the data as provided by the eye-tracker, consisting of $(x, y)$ gaze location entries for all individual fixations. Coordinates were given in pixels with respect to the monitor coordinates (the upper left corner of the screen was coded as $(0, 0)$ and down/right was positive). Additionally, we provide various engineered reading time features. We extend the features provided in the dataset with a range of fixation and saccade based metrics, both on word level and aggregated on the sentence level (see Table 4).

*EEG Acquisition.* High-density EEG data were recorded simultaneously at a sampling rate of 500 Hz with a bandpass of 0.1 to 100 Hz, using a 128-channel EEG Geodesic Hydrocel system (Electrical Geodesics). To ensure good contact, the impedance of each electrode was checked before recording and was kept below 40 $k\Omega$. In addition, electrode impedance levels were checked approx. every 30 mins and reduced if necessary. The EEG data shared in this project are available as raw data but also preprocessed with Automagic Pedroni et al. (2019) toolbox, a tool for automatic EEG data

cleaning and validation. Before the EEG preprocessing with the Automagic toolbox, data from all 14 blocks (7 NR and 7 TSR) were first merged to avoid high predictive power based on the differences resulting from the preprocessing itself. Afterwards, all subfiles whose average standard deviation exceeded $100\mu V$ were excluded.

*EEG Preprocessing.* First, bad channels were detected by the algorithms implemented in the EEGlab plugin `clean_rawdata`[3]. A channel was defined as a bad electrode when recorded data from that electrode was correlated at less than 0.85 to an estimate based on other channels. Furthermore, a channel was defined as bad if it had more line noise relative to its signal than all other channels (4 standard deviations). Finally, if a channel had a longer flat-line than 5 seconds, it was considered bad. These bad channels were automatically removed and later interpolated using a spherical spline interpolation (EEGLAB function `eeg_interp.m`). The interpolation was performed as a final step before the automatic quality assessment of the EEG files. Next, data were filtered using a 2 Hz high-pass filter and line noise artifacts were removed by applying Zapline de Cheveigné (2020), removing seven power line components. Subsequently, independent component analysis (ICA) was performed. Components reflecting artifactual activity were classified by the pre-trained classifier ICLabel Pion-Tonachini et al. (2019). Components that were classified as any class of artifacts (line noise, channel noise, muscle activity, eye activity, and heart artifacts) with a probability higher than 0.8 were removed from the data. Subsequently, residual bad channels were excluded if their standard deviation exceeded a threshold of $25\mu V$. Very high transient artifacts ($> \pm100\mu V$) were excluded from calculating the standard deviation of each channel. However, if this resulted in a significant loss of channel data ($> 50\%$), the channel was removed from the data. After this, the pipeline automatically assessed the quality of the resulting EEG files based on four criteria: First, a data file was marked as bad-quality EEG and not included in the analysis if the proportion of high-amplitude data points in the signals ($> 30\mu V$) was larger than 0.20. Second, more than 20% of time points showed a variance larger than $15\mu V$ across channels. Third, 30% of the channels showed high variance ($> 15\mu V$). Fourth, the ratio of bad channels was higher than 0.3.

After Automagic preprocessing, 13 of the 128 electrodes in the outermost

---

[3]`http://sccn.ucsd.edu/wiki/Plugin_list_process`

circumference (chin and neck) were excluded from further processing as they capture little brain activity and mainly record muscular activity. Additionally, 10 EOG electrodes were separated from the data and not used for further analysis, yielding a total number of 105 EEG electrodes. Subsequently, the EEG and eye-tracking data were synchronized using the "EYE EEG" extension Dimigen et al. (2011) to enable EEG analyses time-locked to the onsets of fixations and saccades, and subsequently segment the EEG data based on the eye-tracking measures. The synchronization was performed in two steps. First, the eye-tracking data were upsampled by linear interpolation to match the number of EEG sampling points. Subsequently, the algorithm identified the "shared" events. Next, a linear function was fitted to the shared event latencies to refine the start- and end-event latency estimation in the eye tracker recording. Finally, synchronization quality was ensured by comparing the trigger latencies recorded in the EEG and eye-tracker data. All synchronization errors did not exceed 2 ms (one sample). The remaining eye artifacts in data were modelled and removed with Unfold toolbox Ehinger and Dimigen (2019) Finally, the data was referenced to the common average reference.

*EEG Feature Extraction.* To compute oscillatory power measures, we band-pass filtered the continuous EEG signals across an entire reading task for four different frequency bands, resulting in a time-series for each frequency band. The distinct frequency bands were determined as follows: *theta* (4-8 Hz), *alpha* (8.5-13 Hz), *beta* (13.5-30 Hz), and *gamma* (30.5-49.5 Hz). Afterwards, we applied a Hilbert transformation to each of these time-series resulting in a complex time series. The Hilbert phase and amplitude estimation method yields results equivalent to sliding window FFT and wavelet approaches Bruns (2004). We chose specifically the Hilbert transformation to maintain temporal information for the amplitude of the frequency bands to enable the power computation of the different frequencies for time segments defined through fixations from the eye tracking. Finally, for each word in each sentence, the EEG features consist of a vector of 105 dimensions (one value for each EEG channel).

## 3. Method

In this section, we first present a preliminary data analysis to investigate the differences between normal reading and task-specific reading reflected in
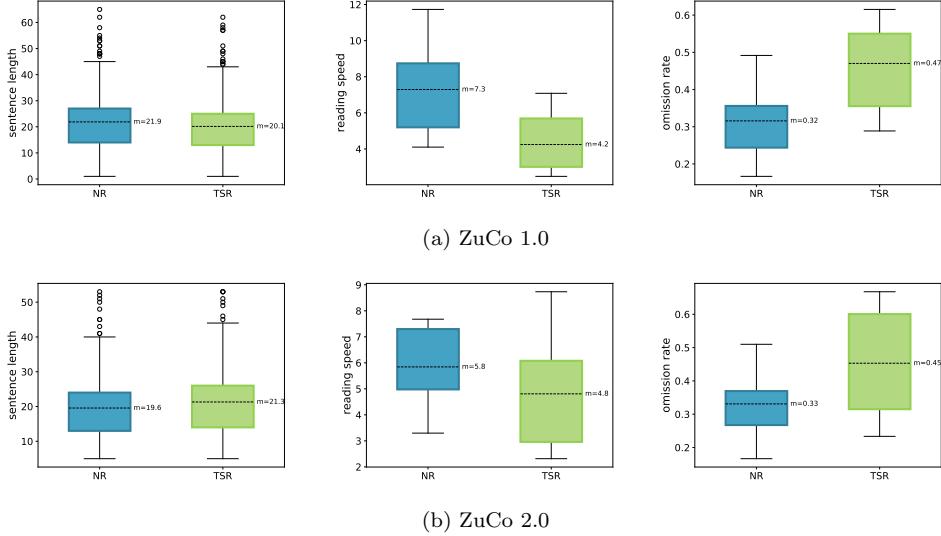
(a) ZuCo 1.0



(b) ZuCo 2.0

Figure 2: Sentence length (words per sentence), reading speed (seconds per sentence) and omission rate (percentage of words not fixated) comparison between normal reading (NR) and task-specific reading (TSR) in ZuCo 1.0 (a) and ZuCo 2.0 (b). The differences in omission rate and reading speed between the two tasks are significant in both datasets.

eye-tracking and EEG data. Following, we describe the machine learning models we developed for the reading task classification. In the first approach (Section 3.2) the models receive word-level features as input, while in the second approach (Section 3.3) the models learn from aggregated sentence-level features.

## 3.1. Preliminary Data Analysis

We analyze the properties of the dataset and compare the two reading tasks. First, we compare the sentence length, reading speed, and omission rate based on the collected eye-tracking data. The sentence length (i.e., the number of words per sentence) was controlled in the selection of reading materials, so that it would not differ significantly between the two tasks (ZuCo 1.0: NR mean 21.9, std 11.1; TSR mean 20.1, std 9.9. ZuCo 2.0: NR mean 19.6, std 8.8; TSR mean 21.3, std 9.5). This is shown in Figure 2, left. Reading speed is defined as the number of seconds spent on a given sentence and the omission rate is the percentage of words that were skipped (i.e., not fixated) per sentence. As expected, in both datasets, the omission rate is higher during task-specific reading (ZuCo 1.0: NR mean 0.32, std 0.09; TSR
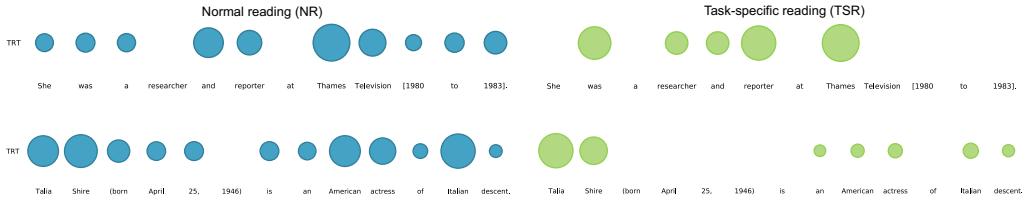
11

Figure 3: Example sentences read in the normal reading paradigm (left) as well as the task-specific experiment paradigm (right) for a single subject. The dots denote the total reading time (TRT) of each word; a larger dot means longer reading time.
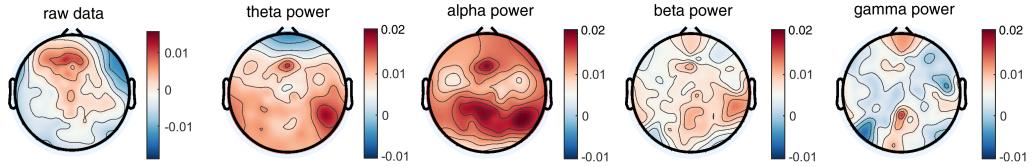


Figure 4: Topography plots showing the EEG activity of the difference between the tasks (NR minus TSR), averaged over all sentences across all subjects from ZuCo 2.0, and averaged across all subjects (scalp viewed from above, nose at the top).

mean 0.47, std 0.11, $p < 0.002$. ZuCo 2.0: NR mean 0.33, std 0.09; TSR mean 0.45, std 0.14, $p < 0.008$), where words are more often skipped, which in turn reduces the reading speed (ZuCo 1.0: NR mean 7.3, std 2.5; TSR mean 4.2, std 1.5, $p < 0.0005$. ZuCo 2.0: NR mean 5.8, std 1.4; TSR mean 4.8, std 2.0, $p < 0.03$). This behavior is visualized in Figure 2, middle and right. Additionally, Figure 3 exemplifies the fixation times of two overlapping sentences between the two reading tasks.

Second, we compare the EEG activity between normal reading and task-specific reading (Figure 4.) These topography plots show higher frontal raw EEG activity in the NR task compared to the TSR task. However, when calculating EEG power in different frequency bands, the lower frequencies (theta and alpha) show widespread increased activity in the NR condition. Contrary, in the high frequency gamma band – typically associated with higher cognitive functions – there is less power in the NR condition compared to TSR.

| Name | Definition | Values |
|------|-----------|--------|
| FIXATION FEATURES | | |
| nFix | number of fixations | 1 |
| FFD | first fixation duration | 1 |
| TRT | total reading time | 1 |
| GD | gaze duration | 1 |
| GPT | go-past time | 1 |
| SACCADE FEATURES | | |
| inSacc_velocity_mean | mean incoming saccade velocity | 1 |
| inSacc_duration_mean | mean incoming saccade duration | 1 |
| inSacc_amplitude_mean | mean incoming saccade amplitude | 1 |
| outSacc_velocity_mean | mean outgoing saccade velocity | 1 |
| outSacc_duration_mean | mean outgoing saccade duration | 1 |
| outSacc_amplitude_mean | mean outgoing saccade amplitude | 1 |
| inSacc_velocity_max | max incoming saccade velocity | 1 |
| inSacc_duration_max | max incoming saccade duration | 1 |
| inSacc_amplitude_max | max incoming saccade amplitude | 1 |
| outSacc_velocity_max | max outgoing saccade velocity | 1 |
| outSacc_duration_max | max outgoing saccade duration | 1 |
| outSacc_amplitude_max | max outgoing saccade amplitude | 1 |

Table 2: Eye-tracking word-level features. *Name* denotes the variable names as used in the dataset.

### 3.2. Word-level Models

First, we describe the reading task classification models leveraging word-level eye-tracking or EEG features. Every sample is assigned a feature vector containing values for each word in the sentence.

*Eye-tracking Features.* Table 2 shows the word-level eye-tracking features we extracted for the reading task classification models. We include standard fixation-based psycholinguistic measures such as the number of fixations and the total reading time. Additionally, we leverage saccade-based features to capture additional reading patterns. For these models, the input for each word is a vector of 5 or 17 feature values, depending on whether saccade features are included or not.

| Name | Definition | Values |
|------|-----------|--------|
| theta | theta frequency band (4-8 Hz) | 105 |
| alpha | alpha frequency band (8.5-13 Hz) | 105 |
| beta | beta frequency band (13.5-30 Hz) | 105 |
| gamma | gamma frequency band (30.5-49.5 Hz) | 105 |
| raw EEG | mean broadband EEG (0-50 Hz) | 105 |

Table 3: EEG word-level features. *Name* denotes the variable names as used in the dataset.

*EEG Features.* Dimigen et al. (2011) demonstrated that EEG indices of semantic processing can be obtained in natural reading and compared to eye movement behavior. The eye tracking data provides millisecond-accurate fixation times for each word. Therefore, the co-registration of both modalities allows the mapping of EEG signals to the visual processing of each word. Hence, by extracting fixation-related potentials (FRPs), we extract the EEG brain activity during the fixation duration of individual words.

Table 3 shows the EEG feature values used for the word-level models. The EEG data for each of these features was extracted for the duration of the total reading time (TRT) of a word. We include a feature composed of the full broadband EEG signals as well as features filtered by frequency bands (theta, alpha, beta, and gamma waves). The input for each word is a vector with 105 electrode values for each frequency band. To capture the total reading time of a word, the EEG signals were averaged over all fixations of a word.

*Model.* Long-Short Term Memory Networks (LSTMs) have been widely adopted for sequence classification tasks (Yu et al., 2019). We implement a bidirectional LSTM to train a reading task classifier based on the word-level features. We train and test the models over five runs with different random seeds, using early stopping to end the training when the validation accuracy is not improving further. Table B.14 presents the values of all hyper-parameters. For the within-subject evaluation, we perform 3-fold cross-validation in each run and average the results across folds and random seeds. We use 10% of the training data used for validation. For the cross-subject evaluation, we train on $n-1$ subjects and test on the data of the left-out subject.
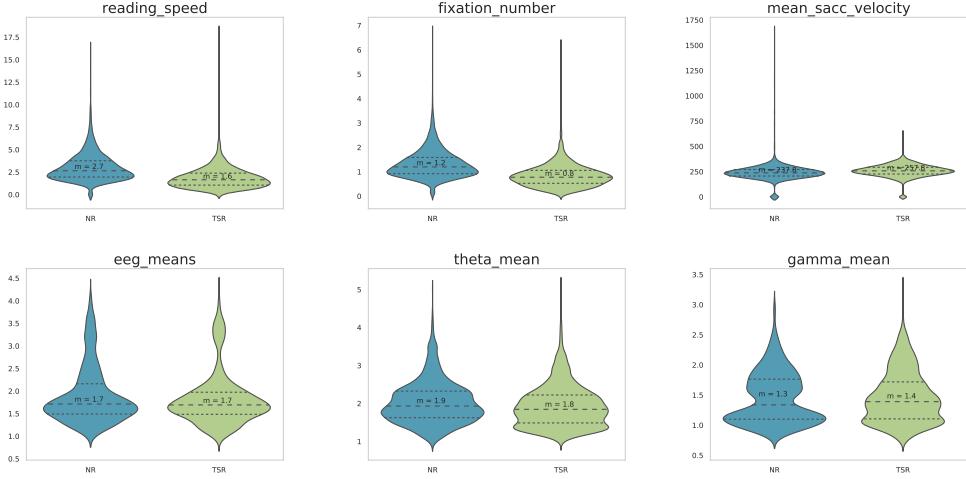
14

Figure 5: Examples of feature distributions across all subjects for ZuCo 1.0.

### 3.3. Sentence-level Models

In addition to word-level features, we investigate the potential of using sentence-level eye tracking and EEG features for the reading task classification. The advantages of sentence-level features consist of the possibility of using simpler machine learning models and reduced training times. Our goal is to investigate if sentence-level features provide the same accuracy as word-level features. Sentence-level features are defined as metrics aggregated over all words in a given sentence. Therefore, every sentence is assigned a single feature value.

*Eye-tracking Features.* We include two types of sentence-level eye-tracking features. The features are summarized in Table 4. First, the fixation-based features - omission rate, number of fixations and reading speed - are aggregated metrics normalized by sentence length, i.e., the number of words in a sentence. Analogous to the word-level models, we also include saccade-based features. These include the mean and maximum duration, velocity and amplitude across all saccades that occurred within the reading time of a give sentence. We test these features individually and combined to investigate the performance increase achieved by adding more features. Examples of these features across all subjects, split by class (normal reading vs. task-specific reading) are shown in Figures and 6, for ZuCo 1.0 and ZuCo 2.0, respectively.
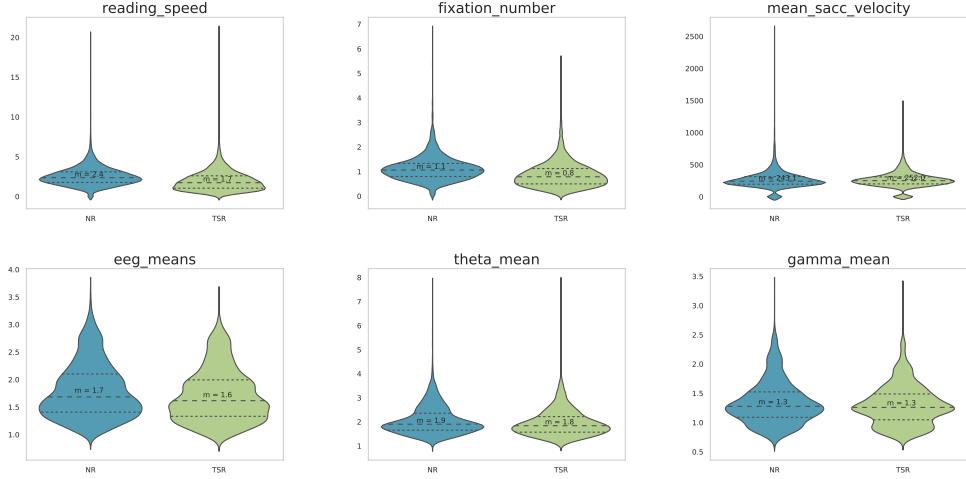
15

Figure 6: Examples of feature distributions across all subjects for ZuCo 2.0.

*EEG Features.* The sentence-level EEG features take into account the EEG activity over the whole sentence duration (even when no words were fixated). We aggregate over the pre-processed EEG signals of the full reading duration of a sentence. The sentence-level EEG features are described in Table 5.

*Model.* The input to the sentence-level model is a single vector representing each sentence. We scale the feature values to a range between $\{-1, 1\}$. We train a support vector machine for classification with a linear kernel. We use the `scikit-learn` SVC implementation[4]. For the within-subject evaluation, we average the results over 50 runs with different random seeds, hence the training and test data are shuffled and split into different sets at every run. The test data is always unseen during training. The models are trained on 90% and tested on 10% of all samples of a single subject during each run. For the cross-subject evaluation, the models are trained on all samples from $n - 1$ subjects and tested on the samples from the left-out subject. Hence, only 1 run is necessary since the samples in the test set remain the same.

*3.4. Baseline Methods*

We compare all models against three baselines: (i) a random baseline, i.e., chance level of 50% for binary classification, (ii) a word embedding baseline,
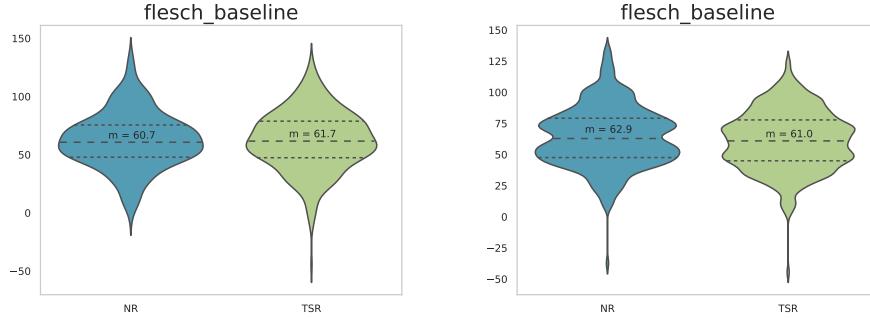
---

[4]https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html

Figure 7: Flesch reading ease (FRE) scores for ZuCo 1.0. (left) and ZuCo 2.0 (right).

and (iii) a text difficulty baseline.

*Word Embedding Baseline.* We compare our models to a textual baseline as a sanity check to ensure the sentences in the data are not easily separable merely by their linguistic characteristics. For this purpose, we use pre-trained textual representations, namely the state-of-the-art contextualized BERT word embeddings (Devlin et al., 2019). We concatenate the embeddings of all words in a sentence feed the into the LSTM model. This word embedding baseline yield a classification accuracy of 58% and 65% for ZuCo 1.0 and ZuCo 2.0, respectively.

*Text Difficulty Baseline.* We also provide a baseline based on text readability. Although the sentences for both reading tasks were chosen to be of similar length and from the same text genre, we want to ensure that both tasks are not separable merely by the difficulty of the sentences. Therefore, we implement a text difficulty baseline, which classifies the sentences into NR and TSR based on their Flesch reading ease score (FRE; Flesch 1948). This score indicates how difficult an English text passage is to understand based on the average number of words in a sentence and the average number of syllables in a word:

$$FRE = x - y\left(\frac{words}{sentences}\right) - z\left(\frac{syllables}{words}\right) \tag{1}$$

where $x$, $y$ and $z$ are language-specific weighting factors (for English $x = 206.835$, $y = 1.015$, $z = 84.6$). We compute FRE scores for each of the

17

English sentences in the ZuCo data. Figure 7 shows the distribution of the FRE across the sentences of ZuCo 1.0 and ZuCo 2.0. This baseline is also above random performance with a classification accuracy of 58% for ZuCo 1.0 and 53% for ZuCo 2.0. Therefore, we also report it in our results.

| Name | Definition | Values |
|---|---|---|
| FIXATION FEATURES | | |
| omission_rate | Percentage of words that is *not* fixated in a sentence | 1 |
| fixation_number | Number of fixations in the sentence divided by the number of words | 1 |
| reading_speed | Sum of the duration of all fixations in the sentence divided by the number of words in the sentence | 1 |
| sent_gaze | Concatenation of the three features described above | 3 |
| SACCADE FEATURES | | |
| mean_sacc_dur | Sum of the duration of all saccades in the sentence divided by the number of words | 1 |
| max_sacc_dur | Maximum saccade duration per sentence | 1 |
| mean_sacc_velocity | Sum of the velocity of all saccades in the sentence divided by the number of saccades | 1 |
| max_sacc_velocity | Maximum saccade velocity per sentence | 1 |
| mean_sacc_amplitude | Sum of the amplitude of all saccades in the sentence divided by the number of saccades | 1 |
| max_sacc_amplitude | Maximum saccade amplitude per sentence | 1 |
| sent_saccade | Concatenation of the six features described above | 6 |
| COMBINED FEATURES | | |
| sent_gaze_sacc | Concatenation of sent_gaze and sent_sacc | 9 |

Table 4: Sentence-level eye-tracking features. *Name* denotes the variable names as used in the dataset.

| Name | Definition | Values |
|---|---|---|
| MEAN FEATURES | | |
| theta_mean | Mean theta band power averaged over all electrodes | 1 |
| alpha_mean | Mean alpha band power averaged over all electrodes | 1 |
| beta_mean | Mean beta band power averaged over all electrodes | 1 |
| gamma_mean | Mean gamma band power averaged over all electrodes | 1 |
| eeg_means | Mean frequency band features averaged over all electrodes, resulting in 1 feature value for each of the 8 frequency bands | 8 |
| ELECTRODE FEATURES | | |
| electrode_features_theta | Mean theta1 and theta2 values of all 104 electrodes | 104 |
| electrode_features_alpha | Mean alpha1 and alpha2 values of all 104 electrodes | 104 |
| electrode_features_beta | mean(sent.mean_b1, sent-mean_b2) | 105 |
| electrode_features_gamma | mean(sent.mean_g1, sent-mean_g2) | 105 |
| electrode_features_all | Concatenation of the four features above | 420 |

Table 5: Sentence-level EEG features. *Name* denotes the variable names as used in the dataset.

| feature set | ZuCo 1.0 | | ZuCo 2.0 | |
| --- | --- | --- | --- | --- |
| | **median** | **MAD** | **median** | **MAD** |
| text difficulty baseline | 0.58 | - | 0.53 | - |
| word embedding baseline | 0.58 | - | 0.65 | - |
| eye_tracking | 0.75 | 0.07 | **0.65** | 0.08 |
| eye_tracking (w/sacc) | **0.75** | 0.05 | **0.65** | 0.04 |
| eeg_raw | 0.88 | 0.06 | 0.64 | 0.05 |
| eeg_theta | 0.98 | 0.01 | 0.62 | 0.07 |
| eeg_alpha | 0.98 | 0.01 | 0.62 | 0.07 |
| eeg_beta | **0.99** | 0.01 | 0.63 | 0.08 |
| eeg_gamma | **0.99** | 0.00 | **0.67** | 0.10 |

Table 6: Summary of all word-level within-subject results. The best results per category are marked in bold.

## 4. Results

In the following, we present the results of both word-level and sentence-level classification models. We present the results on all EEG and eye-tracking feature sets in two settings: within-subject and cross-subjects evaluations. We report all results for both ZuCo datasets. For ZuCo 1.0, we used the data from 11 subjects, and for ZuCo 2.0 from 16 subjects.

### 4.1. Word-level Results

### 4.1.1. Within-subject Evaluation

Table 6 summarizes all word-level within-subject results for both eye-tracking and EEG features. We report the median and median absolute deviation (MAD) of accuracies across all subjects.

Figure 8 show the results of the LSTM models using word-level eye-tracking features including saccade features for ZuCo 1.0 and 2.0. The results without saccade features per subject were almost identical. All within-subject eye-tracking models achieve a performance higher than chance. The median accuracy without saccade features for ZuCo 1.0 is 75% and for ZuCo 2.0 is 65%. For both datasets, the results show that the saccade features do not increase the accuracy, but the variance between multiple runs of individual subjects as well as across all subjects is reduced by including the saccade features.
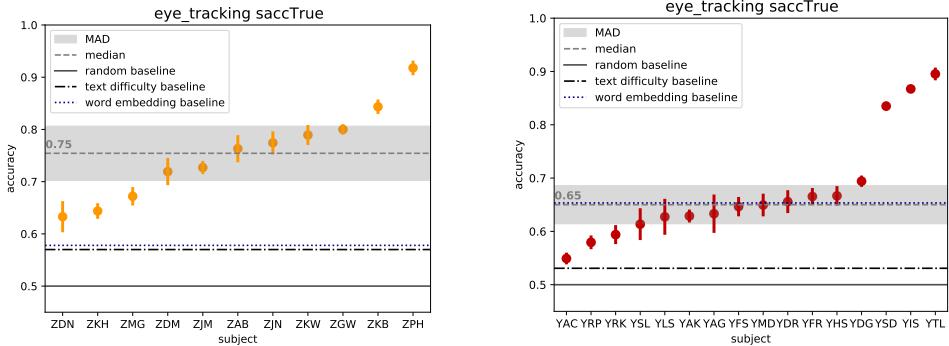
Figure 8: Classification accuracy for each subject on word-level eye-tracking features on ZuCo 1.0 (left) and ZuCo 2.0 (right).
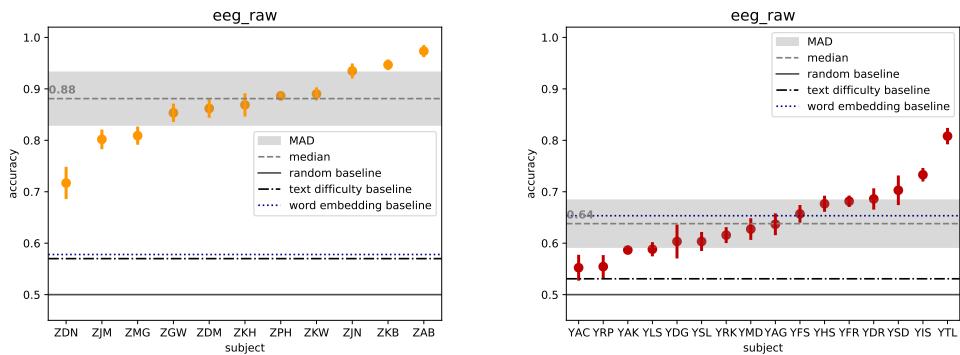


Figure 9: Classification accuracy for each subject on word-level full broadband EEG features on ZuCo 1.0 (left) and ZuCo 2.0 (right).

Figure 9 presents the results achieved by training word-level models on the full broadband EEG signals (*raw_EEG*). The models yield substantially higher classification accuracy when trained and tested on data from ZuCo 1.0 subjects.

The detailed EEG results for all four frequency bands are presented in Figures 10 and 11. We observe near-perfect results (98% accuracy) on all four frequency bands for most subjects of the ZuCo 1.0 dataset, while for ZuCo 2.0 the results show much larger mean absolute deviations and merely 67% accuracy for the best frequency band (i.e., gamma). Generally, all sentence-level feature sets perform better on ZuCo 1.0 than on ZuCo 2.0. This might be due to the session effect, which we discuss in Section 5.
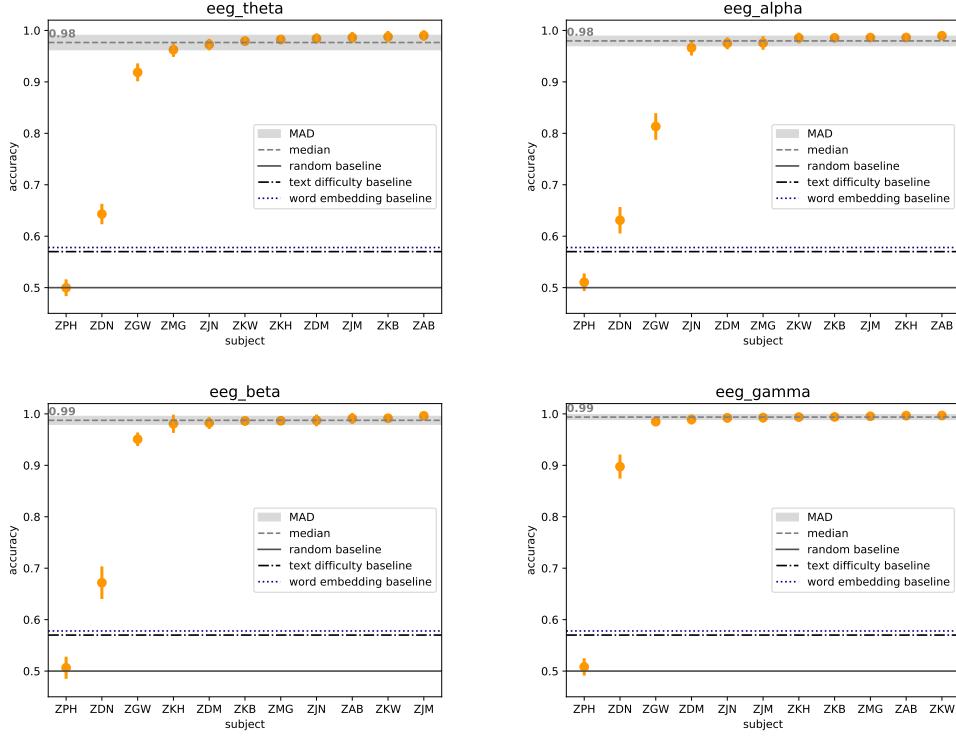
Figure 10: Classification accuracy for within-subject models on word-level EEG frequency band features on ZuCo 1.0.

## 4.1.2. Leave-one-out Cross-subject Evaluation

We also explore the generalization capabilities of the word-level features across subjects in a leave-one-out scenario. The results for the eye-tracking features are shown in Figure 12 and for the EEG features in Figures 13 and 14. The results show that the accuracy of cross-subject models is much lower than that of within-subject models for both datasets. Using the EEG features from ZuCo 2.0, the models of all frequency bands merely reach chance level.
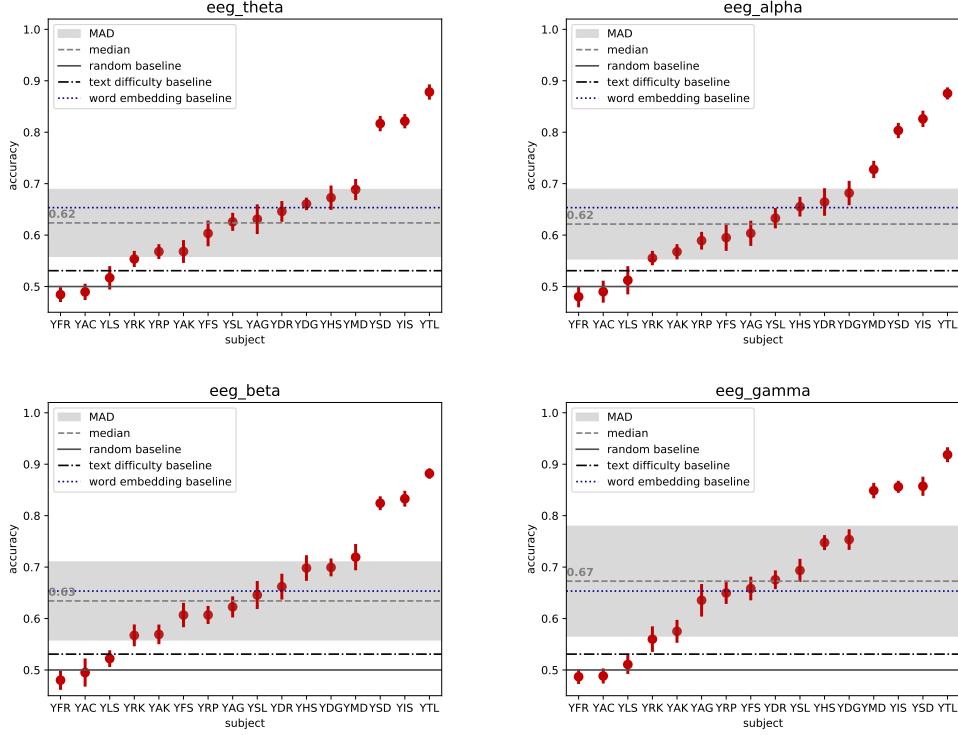
Figure 11: Classification accuracy for each within-subject model on word-level EEG frequency band features on ZuCo 2.0.

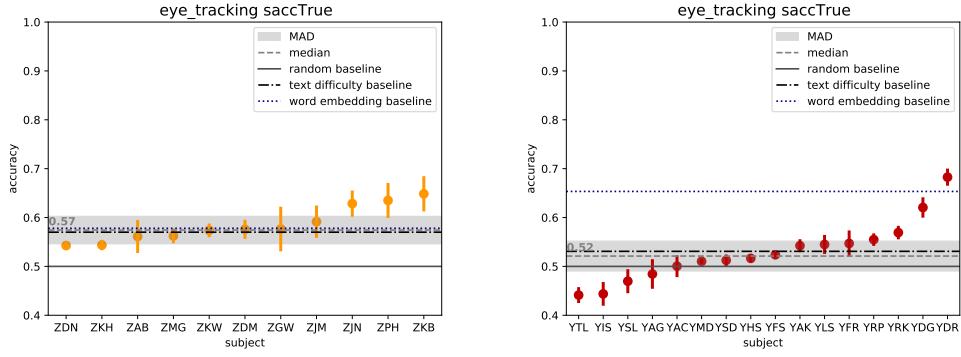| feature set | ZuCo 1.0 | | ZuCo 2.0 | |
|---|---|---|---|---|
| | **median** | **MAD** | **median** | **MAD** |
| eye_tracking | 0.55 | 0.02 | 0.54 | 0.03 |
| eye_tracking sacc | 0.57 | 0.03 | 0.52 | 0.03 |
| eeg_raw | 0.49 | 0.04 | 0.54 | 0.04 |
| eeg_gamma | 0.53 | 0.06 | 0.49 | 0.03 |
| eeg_beta | 0.54 | 0.04 | 0.49 | 0.02 |
| eeg_alpha | 0.54 | 0.02 | 0.50 | 0.02 |
| eeg_theta | 0.56 | 0.06 | 0.49 | 0.02 |

Table 7: Word-level cross-subject result summary.

Figure 12: Accuracy results of the cross-subject eye-tracking word-level classification including saccade features on ZuCo 1.0 (left) and ZuCo 2.0 (right).
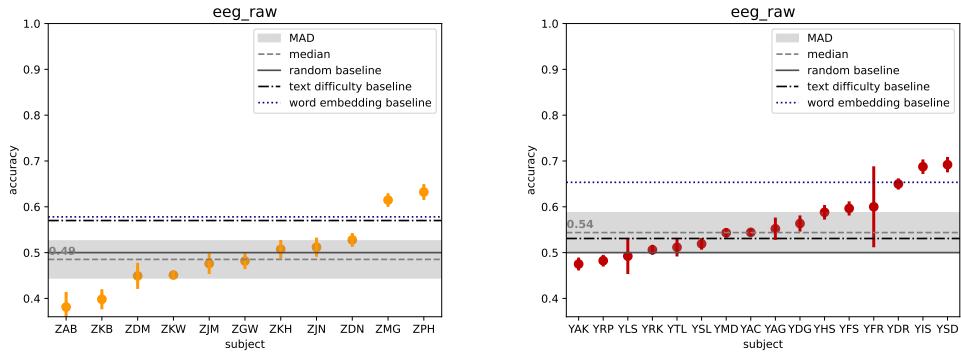


Figure 13: Accuracy results of the cross-subject EEG raw word-level classification including saccade features on ZuCo 1.0 (left) and ZuCo 2.0 (right).
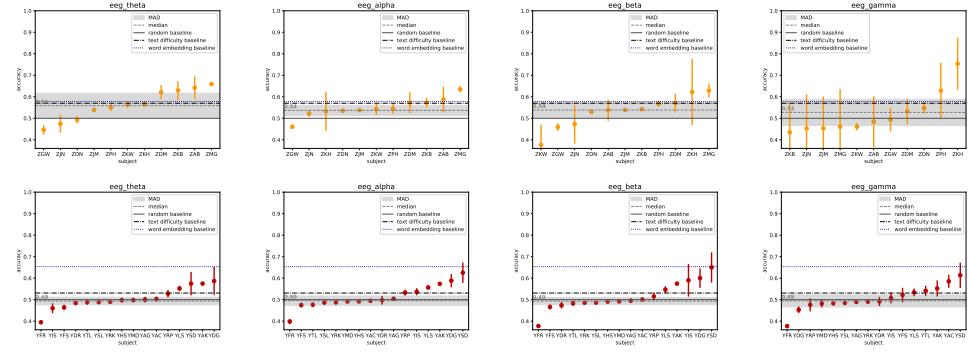


Figure 14: Cross-subject EEG frequency band word-level classification accuracy on ZuCo 1.0 (top) and ZuCo 2.0 (bottom).

*4.2. Sentence-level Results*

*4.2.1. Within-subject Evaluation*

The results for the sentence-level within-subject models are summarized in Table 8. The results on all eye-tracking feature sets are presented in Figure 15 for ZuCo 1.0 and in Figure 16 for ZuCo 2.0. The results show that the models trained on the three fixation-based features perform similarly, but combining them substantially improves the accuracy (*sent_gaze*: ZuCo 1.0 = 76% accuracy, ZuCo 2.0 = 68% accuracy). Moreover, the models trained on the individual saccade features perform worse than the models trained on the fixation features (with only a few subjects considerably higher than random performance), but again the combination of multiple features is helpful (*sent_saccade*: ZuCo 1.0 = 77%, ZuCo 2.0 = 68%). Finally, we observe that combining fixation and saccade features yields the best results using sentence-level eye-tracking features (*sent_gaze_sacc*: ZuCo 1.0 = 82%, ZuCo 2.0 = 70%). Notably, the sentence-level eye-tracking features achieve markedly higher performance on both datasets compared to models based on word-level features.

The results for the EEG feature sets are presented in Figure 17 for ZuCo 1.0 and in Figure 18 for ZuCo 2.0. With the EEG mean features, all frequency bands perform similarly (ZuCo 1.0: lowest - alpha = 62%, highest - gamma = 67%; ZuCo 2.0: lowest - theta = 56%, highest - beta = 58%). Combining all four mean EEG features yields improved results (*eeg_means*: ZuCo 1.0 = 79%, ZuCo 2.0 = 62%). Additionally, the combination of EEG and ET features shows further improvements (*sent_gaze_eeg_means*: ZuCo 1.0 = 88%, ZuCo 2.0 = 72%). However, the most accurate results are achieved using the electrode features. For ZuCo 1.0, most subjects reach 99% accuracy. The gamma frequency band features yield the best performing model. For ZuCo 2, the results of using the theta, alpha, and beta frequency bands range from 70-76% accuracy, whereas using gamma frequency band results in 92% accuracy. This is in line with the quantitative data analysis by Mathur et al. (2021), who find that the EEG frequency ranges are relatively passive during natural reading, although sudden spikes can be observed in task-specific reading.

In general, we observe substantially higher performance on models trained on the ZuCo 1.0 data, i.e., 10% higher accuracy than ZuCo 2.0 for both eye tracking and EEG. Moreover, the gamma frequency band information greatly improves the performance. This is in line with previous research: Adding high frequency band features improves mental task classification

(Zhang et al., 2010). Finally, the sentence-level EEG features yield significantly higher accuracy for ZuCo 2.0 than word-level models. This might be due to the fact that word-level features result in a much high number of parameters than the sentence-level features, which can be detrimental in relation to the low number of training samples available.

| feature set | ZuCo 1.0 median | MAD | ZuCo 2.0 median | MAD |
|---|---|---|---|---|
| text difficulty baseline | 0.58 | - | 0.53 | - |
| word embedding baseline | 0.58 | - | 0.65 | - |
| fixation_number | 0.70 | 0.10 | **0.66** | 0.08 |
| omission_rate | 0.70 | 0.06 | 0.65 | 0.09 |
| reading_speed | **0.72** | 0.10 | 0.65 | 0.09 |
| sent_gaze | **0.76** | 0.08 | **0.68** | 0.08 |
| mean_sacc_dur | 0.59 | 0.06 | 0.55 | 0.05 |
| max_sacc_velocity | 0.61 | 0.06 | 0.57 | 0.04 |
| mean_sacc_velocity | 0.63 | 0.07 | 0.57 | 0.05 |
| max_sacc_dur | 0.58 | 0.04 | 0.54 | 0.04 |
| max_sacc_amp | 0.59 | 0.04 | 0.57 | 0.07 |
| mean_sacc_amp | **0.63** | 0.06 | **0.61** | 0.07 |
| sent_saccade | 0.77 | 0.06 | 0.68 | 0.08 |
| sent_gaze_sacc | **0.82** | 0.06 | **0.70** | 0.07 |
| theta_mean | 0.63 | 0.07 | 0.56 | 0.04 |
| alpha_mean | 0.62 | 0.07 | 0.58 | 0.05 |
| beta_mean | 0.66 | 0.08 | **0.58** | 0.04 |
| gamma_mean | **0.67** | 0.08 | 0.57 | 0.05 |
| eeg_means | 0.79 | 0.09 | 0.62 | 0.05 |
| sent_gaze_eeg_means | **0.88** | 0.06 | **0.72** | 0.06 |
| electrode_features_theta | **1.00** | 0.00 | 0.70 | 0.05 |
| electrode_features_alpha | **1.00** | 0.00 | 0.71 | 0.05 |
| electrode_features_beta | **1.00** | 0.00 | 0.76 | 0.05 |
| electrode_features_gamma | **1.00** | 0.00 | **0.92** | 0.05 |
| electrode_features_all | **1.00** | 0.00 | 0.90 | 0.06 |

Table 8: Summary of all within-subject sentence-level results. The best results per category are marked in bold.
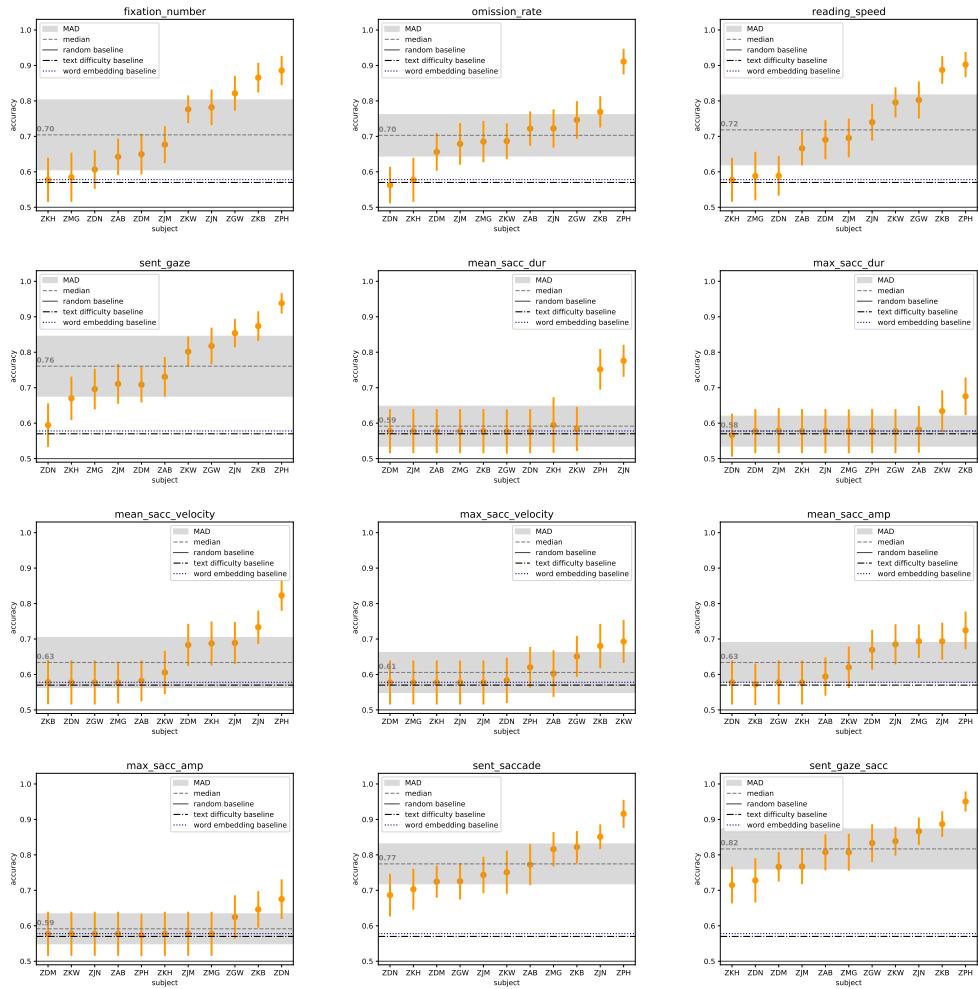
28

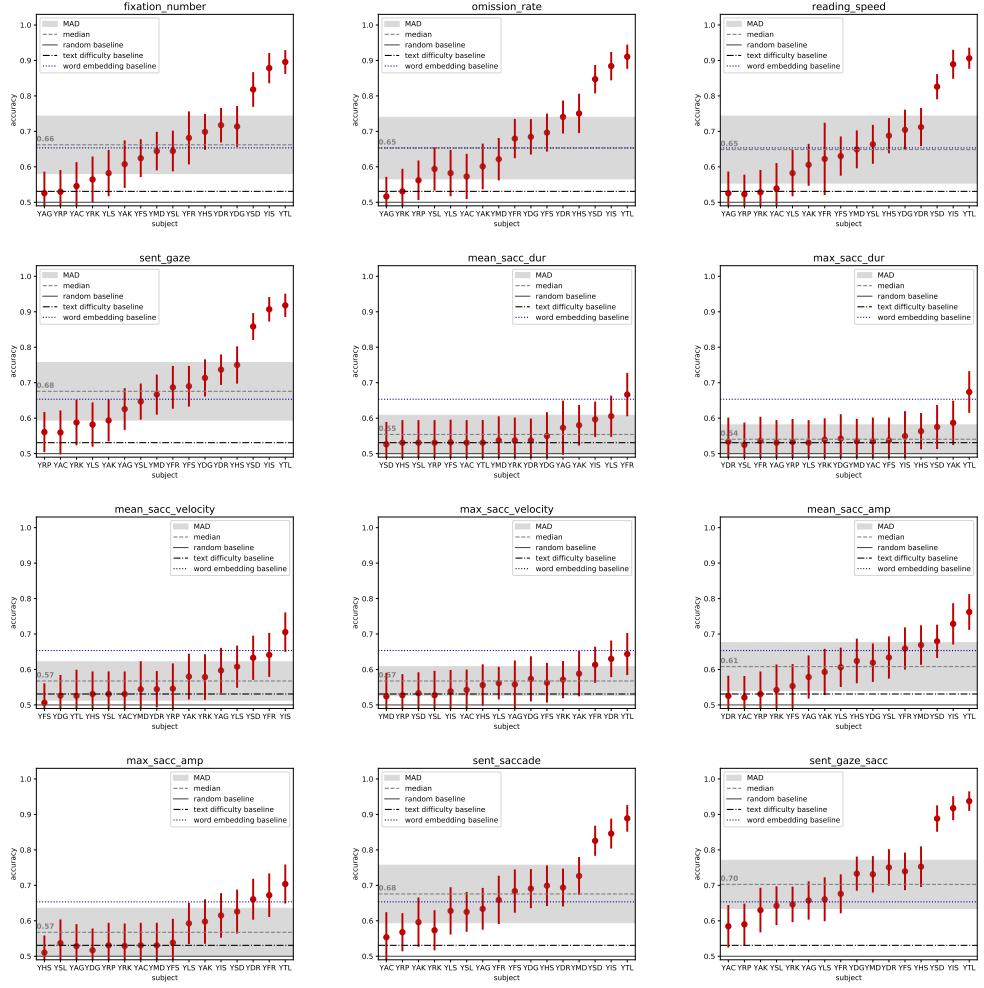Figure 15: Eye-tracking sentence-level classification accuracy on the ZuCo 1.0 data.

Figure 16: Eye-tracking sentence-level classification accuracy on the ZuCo 2.0 data.
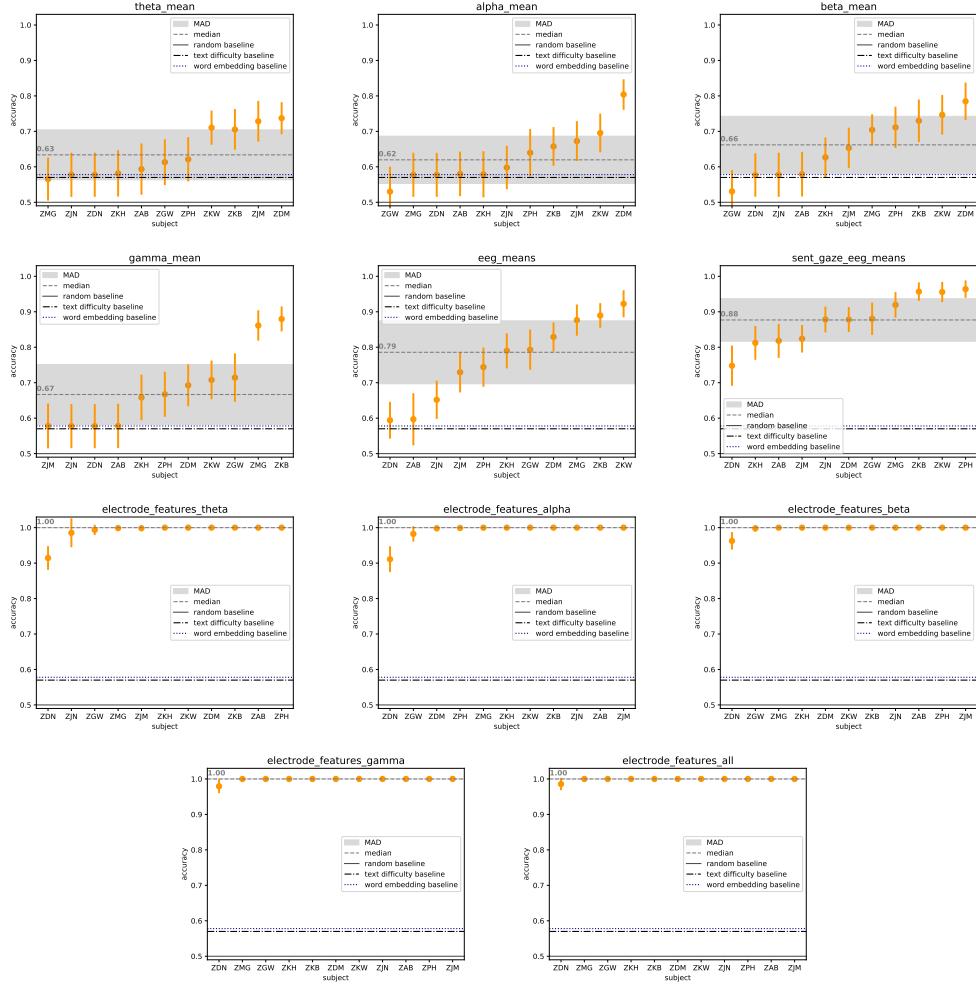
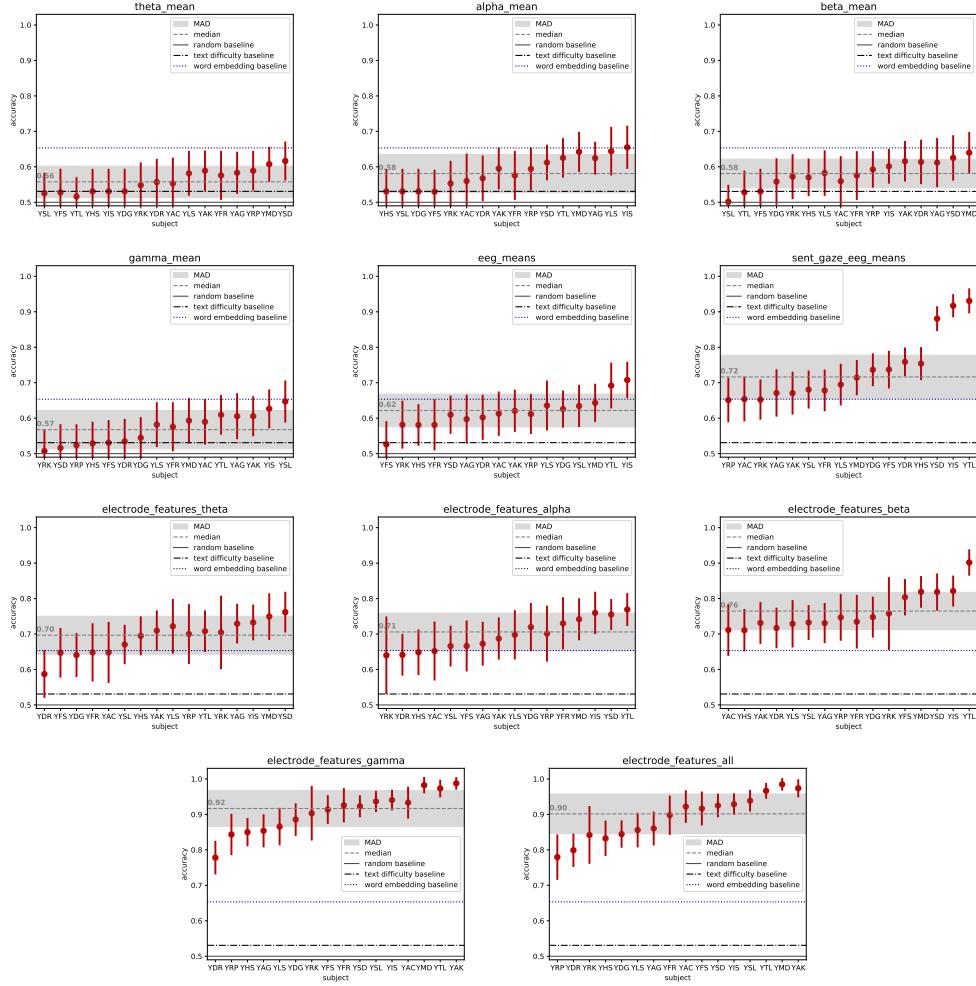Figure 17: EEG sentence-level classification accuracy on the ZuCo 1.0 data.

Figure 18: EEG sentence-level classification accuracy on the ZuCo 2.0 data.

| feature set | ZuCo 1.0 median | MAD | ZuCo 2.0 median | MAD |
|---|---|---|---|---|
| sent_gaze | 0.70 | 0.09 | 0.58 | 0.08 |
| sent_gaze_sacc | 0.60 | 0.03 | 0.60 | 0.08 |
| electrode_features_all | 0.48 | 0.10 | 0.52 | 0.07 |
| electrode_features_gamma | 0.58 | 0.12 | 0.52 | 0.07 |
| electrode_features_beta | 0.50 | 0.09 | 0.53 | 0.04 |
| electrode_features_theta | 0.55 | 0.17 | 0.53 | 0.03 |
| electrode_features_alpha | 0.55 | 0.14 | 0.53 | 0.04 |

Table 9: Sentence-level cross-subject result summary.

*4.2.2. Leave-one-out Cross-subject Evaluation*

To analyze the generalization capabilities of the models, we additionally performed the previously described cross-subject experiments.

The results are presented in Figures 19 - 22. For both ZuCo 1.0 and ZuCo 2.0, when using the fixation-based eye-tracking features (*sent_gaze*), for some subjects we observe accuracies significantly above the text baselines (ZuCo 1.0 median 70%; ZuCo 2.0 median 58%). When including saccade features (*sent_gaze_sacc*), the variability between subjects is reduced, but the accuracy does not improve. Moreover, including saccade features results in similar performance for both ZuCo 1.0 and ZuCo 2.0.

However, using the EEG electrode features, the opposite is the case. ZuCo 1.0 shows high variance between the subjects, the mean accuracy across all models is below random (48%) and only very few subjects achieve performance above the text baseline. The mean performance of the ZuCo 2.0 models is slightly above random (52%). Additionally, we note that in the cross-subject settings, the sentence-level models again perform better than the word-level models.
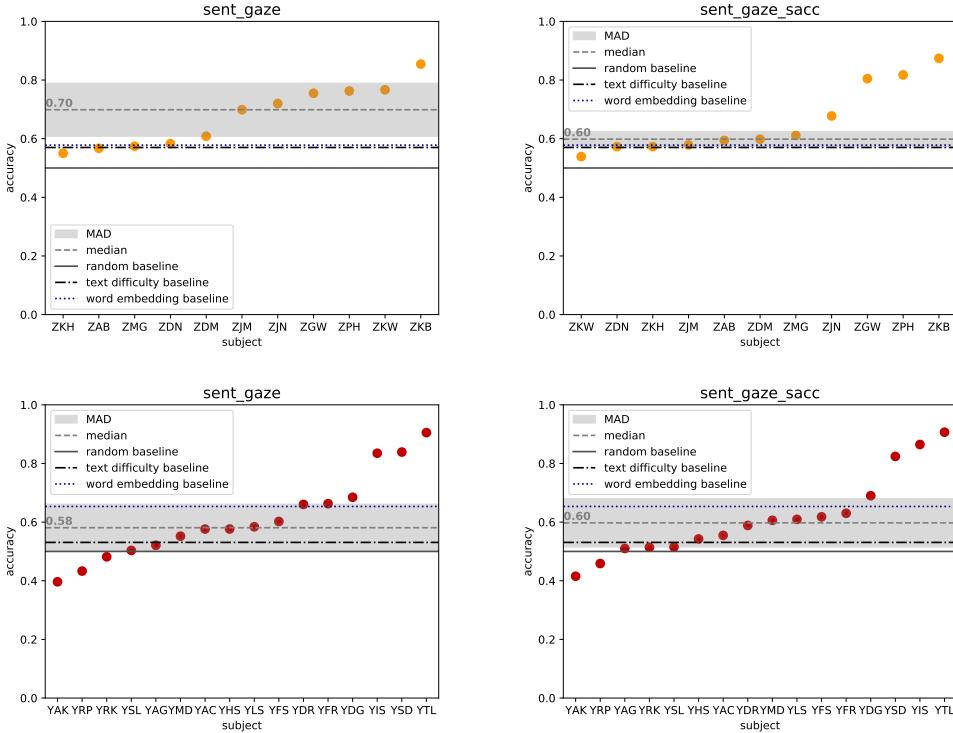
Figure 19: Cross-subject eye-tracking sentence-level classification accuracy on ZuCo 1.0 and ZuCo 2.0, without saccade feature (left) and with saccade features (right).
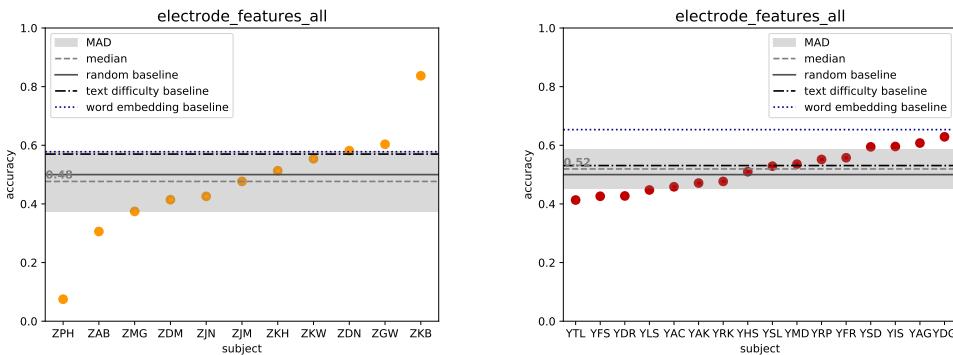


Figure 20: Cross-subject EEG sentence-level classification accuracy on ZuCo 1.0 and ZuCo 2.0 with all electrode features.
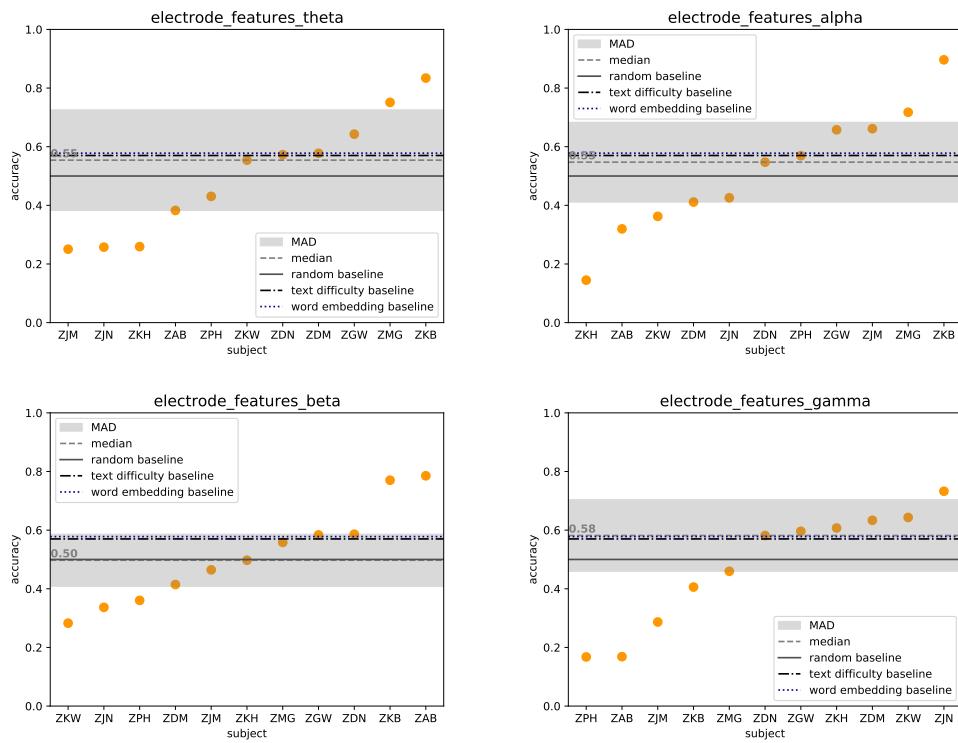
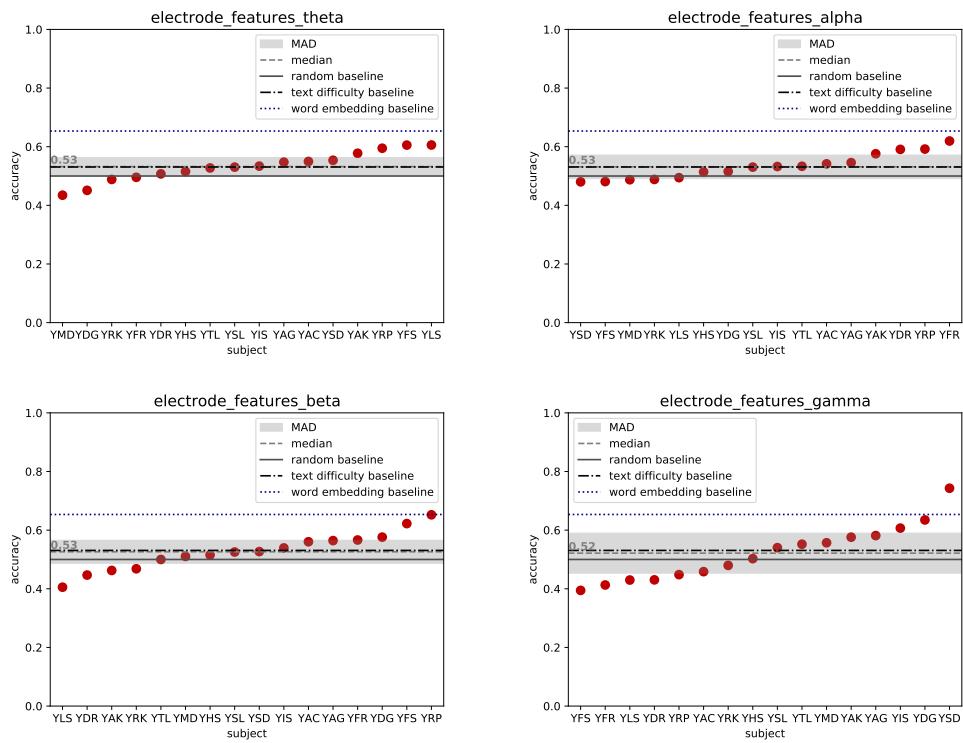Figure 21: Cross-subject EEG sentence-level classification accuracy on ZuCo 1.0.

Figure 22: Cross-subject EEG sentence-level classification accuracy on ZuCo 2.0.

## 5. Control Analyses

In this section, we perform various secondary analyses to validate our results. Since sentence-level models yield higher performance with fewer parameters and lower computational costs, we focus on these results.

### 5.1. Outlier Analysis

First, we detect outlier subjects based on the standard deviation of the feature values and investigate whether the worst- and best-performing subjects are the same as these outliers. We define an outlier as a subject whose feature values deviate from the sample mean by two standard deviations (sample mean±2 std). Table 10 shows the mean feature values before scaling, standard deviation and the resulting outlier subjects for ZuCo 1.0 and ZuCo 2.0. For the EEG mean features, only ZKH from ZuCo 1 and YDG, YSD, YFS from ZuCo 2 are marked as outliers for certain features. For the eye-tracking features, only one subject (ZGW) from ZuCo 1.0 and three subjects (YAG, YRK, YRP) from Zuco 2.0 were marked as outliers. There is no apparent correlation between the model performance and outlier subjects on the within-subject results: On the main EEG feature sets (*electrode_features_all* and *electrode_features_gamma*), the results of all outlier subjects fall within the median absolute deviation (MAD) across all subjects. On the main eye-tracking feature set (*sent_gaze_sacc* including all fixation and saccade features), only the subject YRP from ZuCo 2.0 falls slightly below the MAD.

### 5.2. Correlation Analysis

Second, we explore if the model performance depends on a subject's task performance, level of English language proficiency, or reading speed. Therefore, we investigate whether the performance of the models trained on data from individual subjects correlates with their answer scores in the comprehension questions, their performance in the LexTALE vocabulary test, or their reading speed (see Tables A.12 and A.13).

Table 11 shows the Spearman correlation coefficients for both datasets. We find no significant correlations for the models trained on EEG electrode features. We do see significant negative correlations between the reading speed of the task-specific sentences and the performance of the models trained on eye-tracking features. This is in line with the preliminary data analysis (Section 3.1), showing that participants who fixate on fewer words during task-specific reading have made a clearer mental distinction between the

| | ZuCo 1.0 | | | ZuCo 2.0 | | |
|---|---|---|---|---|---|---|
| **Feature** | Mean | StD | Outliers | Mean | StD | Outliers |
| EYE-TRACKING | | | | | | |
| fixation_number | 1,038 | 0,23 | - | 0,995 | 0,29 | YAG |
| omission_rate | 0,409 | 0,10 | - | 0,393 | 0,11 | - |
| reading_speed | 2,341 | 0,73 | - | 2,240 | 0,63 | - |
| max_sacc_dur | 57,740 | 18,74 | ZGW | 76,834 | 32,84 | YAG,YRK |
| max_sacc_velocity | 576,884 | 104,64 | - | 1023,973 | 446,90 | YRP |
| mean_sacc_dur | 18,512 | 3,03 | ZGW | 22,608 | 5,50 | YRK |
| mean_sacc_velocity | 250,012 | 26,28 | - | 264,913 | 70,07 | YRP |
| EEG | | | | | | |
| theta_mean | 2,786 | 0,63 | - | 2,775 | 0,59 | YDG |
| alpha_mean | 2,585 | 1,16 | ZKH | 2,475 | 0,98 | YSD |
| beta_mean | 2,707 | 1,19 | ZKH | 2,325 | 0,64 | - |
| gamma_mean | 1,785 | 0,39 | - | 1,735 | 0,50 | YFS |

Table 10: Mean feature values before scaling, standard deviation, and the resulting outlier subjects for ZuCo 1.0 (left) and ZuCo 2.0 (right) (sample mean $\pm 2$ std).

two reading tasks. These participants exhibit a different reading strategy and therefore have a higher skipping rate when specifically searching for information in a sentence.

*5.3. Superiority of Gamma Frequency Band Features*

Lastly, to analyze the surprisingly high performance of the gamma frequency band features in all settings, Figure 23 shows topography plots of the SVM coefficients in the within-subject sentence-level reading task classification for the feature electrode_features_gamma for the three best and worst performing subjects of the within-subject models from ZuCo 2.0. To obtain an interpretable topography analysis, the plots were generated by transforming the model into a linear forward model according to the implementation by Haufe et al. (2014).[5] The topography plots show the gamma distribution for the three best performing subjects exhibiting positive SVM coefficients in a

---

[5]https://mne.tools/stable/auto_examples/decoding/linear_model_patterns.html

| | within-subject | | cross-subject | |
| **ZuCo 1.0** | *sent_gaze_sacc* | *EEG elec.* | *sent_gaze_sacc* | *EEG elec.* |
|---|---|---|---|---|
| Score TSR | 0.04 | 0.20 | 0.15 | -0.29 |
| Score NR | 0.33 | 0.05 | -0.03 | -0.13 |
| Speed TSR | -0.18 | 0.30 | **-0.42** | 0.02 |
| Speed NR | **0.58** | **0.50** | 0.11 | 0.27 |
| LexTALE | 0.22 | 0.00 | 0.40 | **-0.32** |
| **ZuCo 2.0** | *sent_gaze_sacc* | *EEG elec.* | *sent_gaze_sacc* | *EEG elec.* |
| Score TSR | 0.54 | -0.12 | 0.26 | 0.01 |
| Score NR | 0.14 | -0.05 | 0.34 | 0.33 |
| Speed TSR | **-0.75*** | 0.07 | **-0.61*** | 0.23 |
| Speed NR | -0.16 | 0.19 | -0.04 | 0.23 |
| LexTALE | 0.27 | **0.22** | 0.15 | **0.43** |

Table 11: Spearman correlation between the within-subject and cross-subject classification accuracies of all subjects in ZuCo 1.0 (top) and ZuCo 2.0 (bottom) for a given feature set and the control scores or reading speed. The feature sets are ET = sent_gaze_sacc, EEG mean = eeg_means, and EEG elec. = electrode_features_all. * marks significant correlations $p < 0.05$.

mid frontal electrode cluster, whereas for the worst performing subjects the positive SVM coefficients are pronounced in the occipital electrodes.

*5.4. Subject Classification*

As a control experiment, we train models to recognize the individual subjects instead of the reading tasks. As shown in Figure 24, the EEG electrode features again yield the highest accuracy for this subject classification for both ZuCo 1.0 and ZuCo 2.0, while the eye-tracking and EEG mean feature sets do not yield such good performance. However, most of the feature sets do reach an accuracy higher than random, defined as the probability of randomly classifying the correct subject in a balanced dataset (i.e, 1/11 for ZuCo 1.0 and 1/16 for ZuCo 2.0. This shows how the EEG electrode features are much more subject-specific than the eye tracking features. This is possibly be due to the loss of information in the eye-tracking features as well as the EEG mean features as they are aggregated on a higher level. This subject classification task opens opportunities for further research exploring privacy preservation in task-specific EEG and eye-tracking signals.
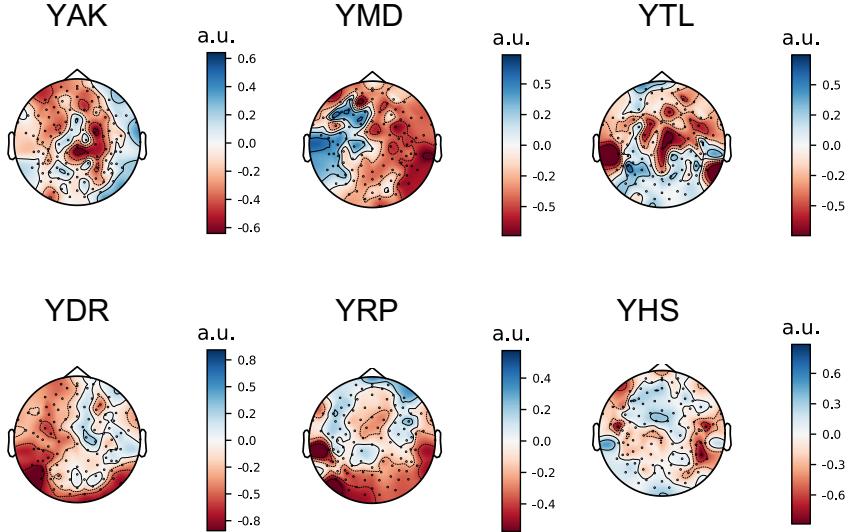
Figure 23: Topography plots of the SVM coefficients for the feature electrode_features_gamma for the three best (top) and three worst (bottom) performing subjects from ZuCo 2.0.

## 5.5. Fixation Ablation

In an additional experiment, we extract the EEG features in their chronological fixation order, instead of the word order within the sentences as in all other experiments in this paper. The input dimension of the feature vectors remains the same. However, instead of averaging the electrode value across all fixations, we average only on the first 10%, 20%, 50% or 75% of fixations. For this control experiment, we use the gamma frequency band features, since these yielded the best results.

We test how the performance of the reading task classification is affected when taking only a proportion of the fixated words. We hypothesize that only a small percentage of the EEG signals is necessary to accurately predict the reading task, speculating that the mental state is different enough between the reading tasks from the beginning.

Figure 25 shows that this hypothesis holds true for the participants of ZuCo 1.0, where only one subject shows a substantial decrease in classification accuracy with a lower percentage of fixations. However, for most subjects even 10% of the EEG signal is enough for an accurate classification (>99% accuracy). This might be due to the fact that the NR and TSR read-
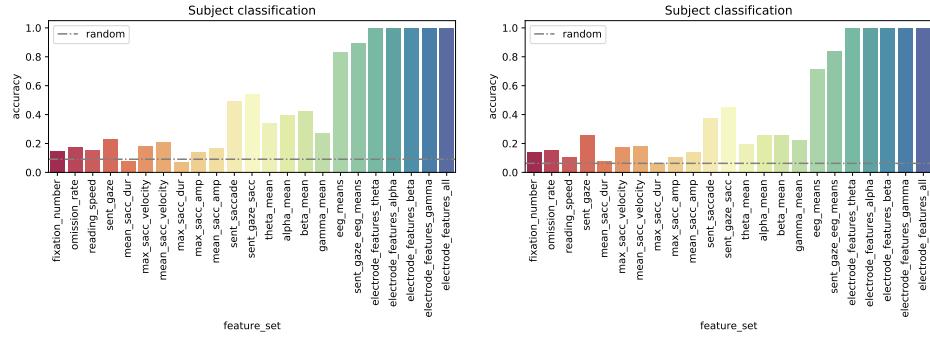
Figure 24: Subject classification for ZuCo 1.0 (left) and ZuCo 2.0 (right). The results are averaged across all subjects.
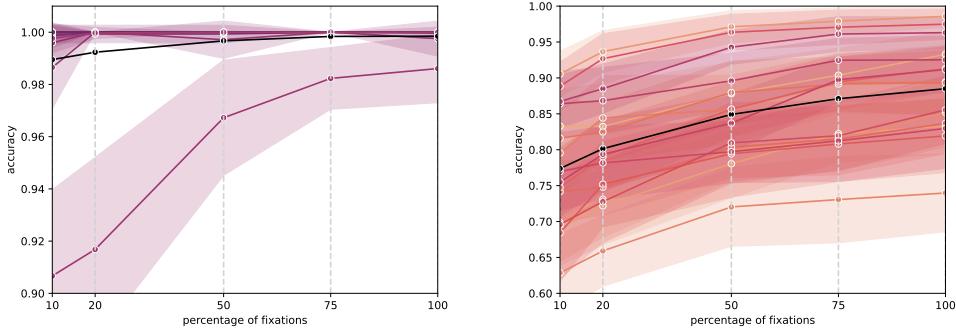


Figure 25: Classification accuracy for each subject for the EEG electrode features of the gamma frequency band in fixation order on sentence level, with varying fixation percentage. ZuCo 1.0 on the left, ZuCo 2.0 on the right.

ing tasks were recorded in separate sessions in ZuCo 1.0. For ZuCo 2.0, the results show that models trained only on the first 10-50% of EEG signals do yield lower accuracy. The highest results are achieved when training on the EEG signals of all fixations in a sentence. These results show the impact of the session bias in ZuCo 1.0: The fact that the reading tasks were recorded in different sessions presents as a systematic differentiation in the brain activity signals. Therefore, the EEG features only from the beginning fixations of every sentence are sufficient for accurate classification. In contrast, for ZuCo 2.0, recorded in a single session, the accuracy decreases with a lower proportion of the signals.

41

## 5.6. Session & Block Classification

*ZuCo 1.0: Session Classification*

Auto-correlations between the trials (or blocks) of an experimental session due to a range of different factors such as familiarization with the task or mental fatigue is a well-known challenge in experimental psychology (Baayen et al., 2017). Moreover, it has been demonstrated that a subject's eye movements show session-specific properties (Bargary et al., 2017). We therefore perform further analyses in an attempt to explore and quantify the session and block effects that emerge from the experiment design. For ZuCo 1.0, we can test the session effect by adding the sentences recorded in a third task, the sentiment reading data (see description in Appendix Appendix A.3). As described above, for ZuCo 1.0, the normal reading and task-specific reading paradigms were recorded in separate sessions. In each of the two sessions, half of the sentiment reading (SR) data was also recorded. Hence, by adding these sentences and training a model to classify in which session each sentence was recorded, we can get an estimate of the session-specific information contained in the EEG signals. By investigating this session bias, we explore whether the model exploits some uncontrolled systematic differences between the sessions to perform the task classification. By adding sentences of the same task recorded in both sessions, we hypothesize that the performance of this session classification task (including NR, TSR, and SR sentences) should be lower than for the reading task classification with only NR and TSR sentences, because the model cannot rely on differences between the recording sessions. This would show that session-specific information is partly a cause of the high accuracy for ZuCo 1.0.

Figure 26 shows how adding the sentence from the additional paradigm recorded in both sessions of ZuCo 1.0, decreases the performance for most features, albeit not by much. Notably, the differences are larger for eye-tracking features and the accuracy is almost not affected for EEG electrode features.

Figure 27 shows the topography plots of the SVM coefficients for the gamma activity of all electrodes for three randomly chosen subjects, the top row for reading task classification and the bottom for this experiment of session classification including the additional sentences. Again, the plots were generated according to the implementation by Haufe et al. (2014). Similar to the analysis based only on the ZuCo 2.0 data set, the highest SVM coefficients also seem to be in a mid frontal electrode cluster, although there is not a
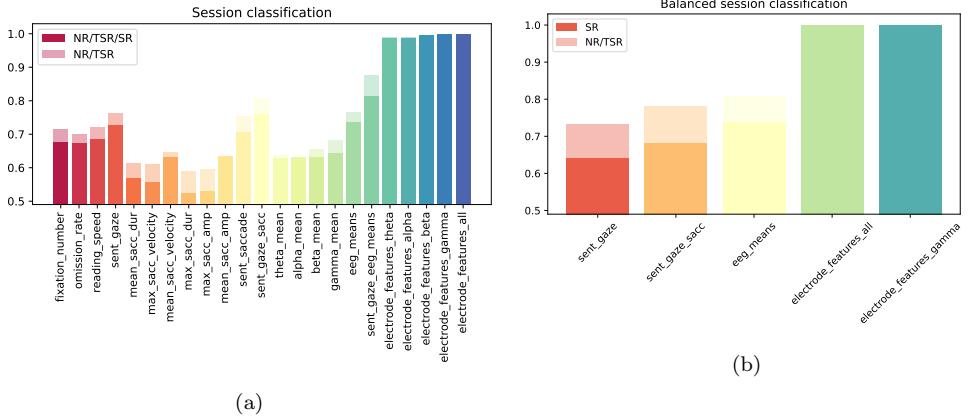
Figure 26: Session classification. **(a)** Classifying the ZuCo 1.0 sentences into the two recording sessions, with and without the additional sentences. In the latter case, the labels are identical to all reading task classification experiments presented previously on ZuCo 1.0. **(b)** Session classification with balanced datasets. Classifying the ZuCo 1.0 sentences into the two recording sessions to quantify the session effect.



Figure 27: Topography plots of the SVM coefficients for all electrodes (gamma) - top for reading task classification and bottom for session classification including the additional sentences.

clear pattern shared across different subjects.

Additionally, for ZuCo 1.0 we perform the session classification on completely balanced datasets, i.e., the same number of samples for the first and second session during training and testing the model. Figure **??** (b) shows
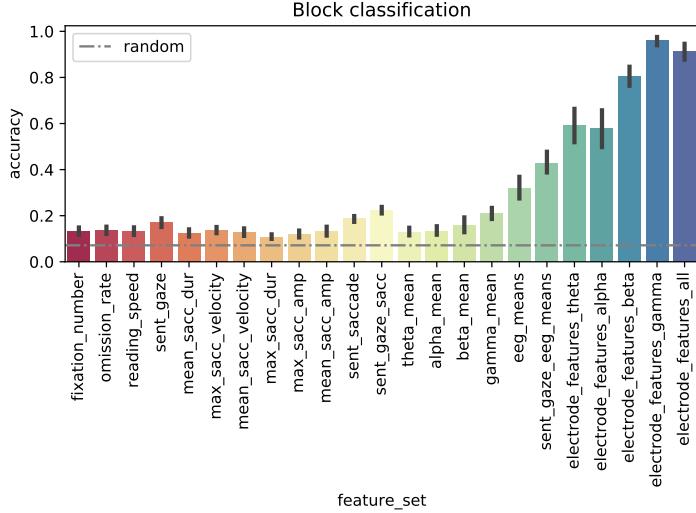
Figure 28: Classifying the ZuCo 2.0 sentences into the 14 recording blocks.

the results for the most relevant eye tracking and EEG feature sets. We compare the performance of classifying the additional sentences of the sentiment reading task (SR) recorded in two sessions, to classifying the same number of sentences from NR and TSR. We see that when classifying these SR sentences, the performance drops between 5% - 10% in accuracy for eye-tracking (with and without saccadic features) as well as for the mean EEG features. However, for the EEG electrode features, the performance does not drop.

These additional session classification experiments show that the session effect does impact the reading task classification for the eye-tracking features and the mean EEG features, but it cannot entirely explain the high performance in the ZuCo 1.0 models, especially for the models with EEG electrode features, where the performance does not decrease at all. From these results, we can conclude that the session bias is more noticeable in aggregated features than in individual electrode features.

*ZuCo 2.0: Block Classification*
For ZuCo 2.0, all sentences were recorded in the same session. Therefore, there are no session-specific biases. However, we can analyze the effect of the order in which the sentences were recorded. As described in Section 2, the sentence blocks of normal reading and task-specific sentences were alternated. Each of the 14 blocks contains approx. 50 sentences of either normal reading
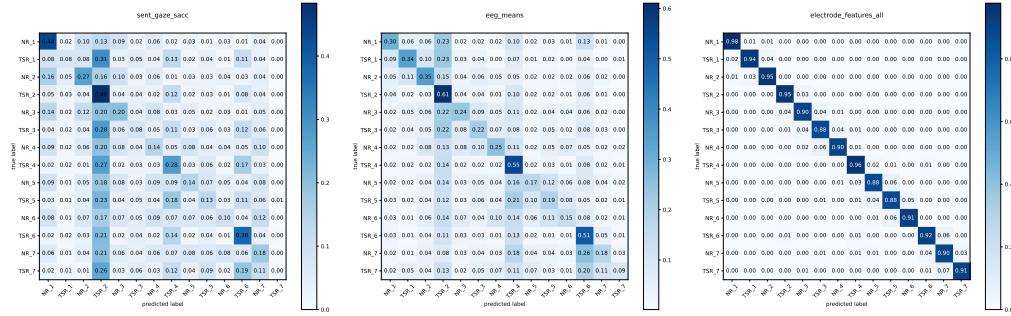
Figure 29: Confusion matrices for block classification of ZuCo 2.0 for the three feature sets *sent_gaze_sacc*, *eeg_means* and *electrode_features_all*.

or task-specific reading. Hence, we also test the performance of classifying the sentences into these 14 blocks with within-subject sentence-level models.

The results are averaged across all subjects. Figure 28 shows how all eye-tracking features as well as the mean EEG features perform only slightly above random for this block classification task, while the electrode features – especially from the gamma frequency band – still achieve high accuracy. A closer look at the confusion matrices presented in Figure 29, shows that with eye-tracking features the blocks are often confounded within the same reading task (left), which shows that effectively the reading patterns between the two tasks are substantially different. When using EEG mean features to classify the sentence into recording blocks, the blocks are more often confounded with neighboring blocks in terms of recording order (middle). Finally, when using the EEG electrode features, the few classification mistakes also occur mostly within neighboring blocks (right). This pattern can be explained by the nature of the EEG experiments, where the electrode impedance was tested and corrected after every 3-4 blocks.

One concern with the recording blocks is that they were recorded in the same order for all participants. Hence, in an attempt to quantify the effect of the variability in the data between the different sequential recording blocks of ZuCo 2.0, we perform an ablation study and train the binary reading task classification using sentence-level models with a decreasing number of blocks. This means that the models are trained on data ranging from one randomly selected block per reading task (i.e., one NR block and one block)
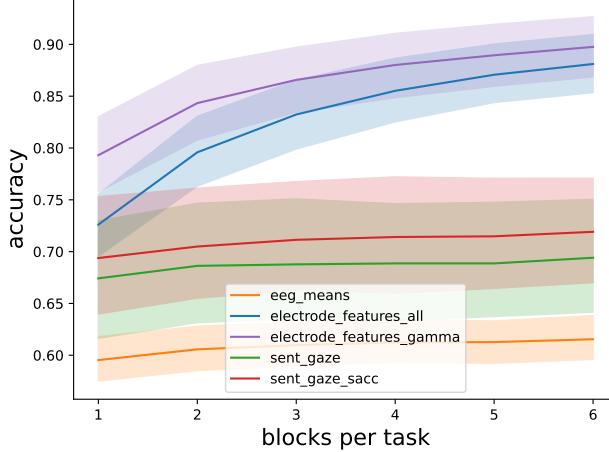
45

Figure 30: Reading task classification on the ZuCo 2.0 data with decreasing number of recording blocks in the training data.

up to six blocks per reading task, and tested on sentences of the remaining blocks. The hypothesis is that if systematic information about the reading tasks is encoded in the order of the blocks, the classification performance will increase when adding more block to the training data.

The results are shown in Figure 30. We observe that the performance of both the eye-tracking features as well as the EEG mean features only increases slightly with additional blocks, while accuracy of the models trained on the EEG electrode features (both all frequency bands and gamma) increase substantially with data from more blocks. This shows that the EEG electrode features are more sensitive to the information about the recording blocks encoded in the brain activity signals and that the reading task classification benefits from this information.

## 6. Conclusion

Reading is a complex cognitive process that requires the simultaneous processing of visual and higher-level linguistic information including syntactic, semantic and discourse integration. Identifying task-specific reading patterns can improve models of human reading and provide insights into human language understanding and how we perform linguistic tasks. This knowledge can then be applied to machine learning algorithms for natural language processing.

Accurate reading task classification can improve the manual labelling process for a variety of NLP tasks, as these processes are closely related to identifying reading intents. Recognizing reading patterns for estimating reading effort has additional applications such as as the diagnosis of reading impairments such as dyslexia (Rello and Ballesteros, 2015) and attention deficit disorder (Tor et al., 2021).

The rise of more accessible behavioral and physiological datasets tailored towards machine learning allows their use to advance natural language processing methods. The data quality validation of both datasets, including a detailed analysis of eye tracking and EEG features and a comparison to previous studies, can be found in the original publications (Hollenstein et al., 2018, 2020).

In this work, we presented various machine learning methods for reading task classification. The ML models learn to distinguish between a normal reading task and task-specific information-searching reading. We develop models that learn from eye movements or EEG signals. While the within-subject evaluation yields high classification accuracy, further research is required to improve the cross-subject generalization capabilities of the models. This is crucial for any potential applications of this classification task. Furthermore, we extensively analyze the differences between ZuCo 1.0 and ZuCo 2.0, and how the session bias of ZuCo 1.0 impacts the classification results. We address some of the open challenges in building robust ML models for task classification based on eye movement and brain activity data, including signal-to-noise ratio, high inter-subject variability, and generalization across sessions and datasets.

### References

Keith Rayner. Eye movements in reading and information processing: 20 years of research. *Psychological bulletin*, 124(3):372, 1998.

Olaf Dimigen, Werner Sommer, Annette Hohlfeld, Arthur M Jacobs, and Reinhold Kliegl. Coregistration of eye movements and EEG in natural reading: analyses and review. *Journal of Experimental Psychology: General*, 140(4):552, 2011.

Luz Rello and Miguel Ballesteros. Detecting readers with dyslexia using machine learning with eye tracking measures. In *Proceedings of the 12th International Web for All Conference*, pages 1–8, 2015.

Peter Raatikainen, Jarkko Hautala, Otto Loberg, Tommi Kärkkäinen, Paavo Leppänen, and Paavo Nieminen. Detection of developmental dyslexia with machine learning using eye movement data. *Array*, 12:100087, 2021.

Hui Tian Tor, Chui Ping Ooi, Nikki SJ Lim-Ashworth, Joel Koh En Wei, V Jahmunah, Shu Lih Oh, U Rajendra Acharya, and Daniel Shuen Sheng Fung. Automated detection of conduct disorder and attention deficit hyperactivity disorder using decomposition and nonlinear techniques with eeg signals. *Computer Methods and Programs in Biomedicine*, 200:105941, 2021.

Nora Hollenstein, Jonathan Rotsztejn, Marius Troendle, Andreas Pedroni, Ce Zhang, and Nicolas Langer. ZuCo, a simultaneous EEG and eye-tracking resource for natural sentence reading. *Scientific Data*, 2018.

Nora Hollenstein, Marius Troendle, Ce Zhang, and Nicolas Langer. ZuCo 2.0: A dataset of physiological recordings during natural reading and annotation. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 138–146, 2020.

Takenobu Tokunaga, Hitoshi Nishikawa, and Tomoya Iwakura. An eye-tracking study of named entity annotation. *Proceedings of the International Conference Recent Advances in Natural Language Processing*, pages 758–764, 2017.

Katrin Tomanek, Udo Hahn, Steffen Lohmann, and Jürgen Ziegler. A cognitive cost model of annotations based on eye-tracking data. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 1158–1167, 2010.

John-Dylan Haynes and Geraint Rees. Decoding mental states from brain activity in humans. *Nature Reviews Neuroscience*, 7(7):523–534, 2006.

Puneet Mathur, Trisha Mittal, and Dinesh Manocha. Dynamic graph modeling of simultaneous eeg and eye-tracking data for reading task identification. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1250–1254. IEEE, 2021.

Samuel A Nastase, Ariel Goldstein, and Uri Hasson. Keep it real: rethinking the primacy of experimental control in cognitive neuroscience. *NeuroImage*, 222:117254, 2020.

Maria Barrett, Joachim Bingel, Nora Hollenstein, Marek Rei, and Anders Søgaard. Sequence classification with human attention. In *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pages 302–312, 2018.

Sandeep Mathias, Diptesh Kanojia, Abhijit Mishra, and Pushpak Bhattacharyya. A survey on using gaze behaviour for natural language processing. *Proceedings of IJCAI*, 2020.

Erik McGuire and Noriko Tomuro. Relation classification with cognitive attention supervision. In *Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics*, pages 222–232, 2021.

Nora Hollenstein, Antonio de la Torre, Nicolas Langer, and Ce Zhang. CogniVal: A framework for cognitive word embedding evaluation. In *Proceedings of the 23nd Conference on Computational Natural Language Learning*, 2019.

Nora Hollenstein and Lisa Beinborn. Relative importance in sentence processing. *Association for Computational Linguistics*, 2021.

Christian Pfeiffer, Nora Hollenstein, Ce Zhang, and Nicolas Langer. Neural dynamics of sentiment processing during naturalistic sentence reading. *NeuroImage*, page 116934, 2020.

Louise Gillian Bautista and Prospero Naval. Towards learning to read like humans. In *International Conference on Computational Collective Intelligence*, pages 779–791. Springer, 2020.

Yves Bestgen. LAST at CMCL 2021 Shared Task: Predicting Gaze Data During Reading with a Gradient Boosting Decision Tree Approach. In *Proceedings of the NAACL Workshop on Cognitive Modeling and Computational Linguistics*, 2021.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, 2013.

Aron Culotta, Andrew McCallum, and Jonathan Betz. Integrating probabilistic extraction models and data mining to discover relations and patterns in text. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics*, pages 296–303, 2006.

Andreas Pedroni, Amirreza Bahreini, and Nicolas Langer. Automagic: Standardized preprocessing of big EEG data. *NeuroImage*, 2019.

Alain de Cheveigné. Zapline: A simple and effective method to remove power line artifacts. *NeuroImage*, 207:116356, 2020.

Luca Pion-Tonachini, Ken Kreutz-Delgado, and Scott Makeig. ICLabel: An automated electroencephalographic independent component classifier, dataset, and website. *NeuroImage*, 198:181–197, 2019.

Benedikt V Ehinger and Olaf Dimigen. Unfold: an integrated toolbox for overlap correction, non-linear modeling, and regression-based EEG analysis. *PeerJ*, 7:e7838, 2019.

Andreas Bruns. Fourier-, Hilbert- and wavelet-based signal analysis: Are they really different approaches? *Journal of neuroscience methods*, 137 (2):321–332, 2004.

Yong Yu, Xiaosheng Si, Changhua Hu, and Jianxun Zhang. A review of recurrent neural networks: Lstm cells and network architectures. *Neural computation*, 31(7):1235–1270, 2019.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.

Rudolph Flesch. A new readability yardstick. *Journal of applied psychology*, 32(3):221, 1948.

Li Zhang, Wei He, Chuanhong He, and Ping Wang. Improving mental task classification by adding high frequency band information. *Journal of medical systems*, 34(1):51–60, 2010.

Stefan Haufe, Frank Meinecke, Kai Görgen, Sven Dähne, John-Dylan Haynes, Benjamin Blankertz, and Felix Bießmann. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *Neuroimage*, 87: 96–110, 2014.

Harald Baayen, Shravan Vasishth, Reinhold Kliegl, and Douglas Bates. The cave of shadows: Addressing the human factor with generalized additive mixed models. *Journal of Memory and Language*, 94:206–234, 2017.

Gary Bargary, Jenny M. Bosten, Patrick T. Goodbourn, Adam J. Lawrance-Owen, Ruth E. Hogg, and JD Mollon. Individual differences in human eye movements: An oculomotor signature? *Vision Research*, 141:157–169, 2017.

Kristin Lemhöfer and Mirjam Broersma. Introducing LexTALE: A quick and valid lexical test for advanced learners of English. *Behavior Research Methods*, 44(2):325–343, 2012.

## Appendix A. Data Collection

Please refer to Hollenstein et al. (2018) and Hollenstein et al. (2020) for detailed descriptions of the data acquisition, preprocessing and feature extraction methods for ZuCo 1.0 and ZuCo 2.0, respectively.

### Appendix A.1. Participants

For ZuCo 1.0, data were recorded from 12 healthy adults (between 22 and 54 years old; all right-handed; 5 female subjects). For ZuCo 2.0, data were recorded from 18 healthy adults (between 23 and 52 years old; 2 left-handed; 10 female subjects). The native language of all participants is English, originating from Australia, Canada, UK, USA or South Africa. In addition, all subjects completed the standardized LexTALE test to assess their vocabulary and language proficiency (Lexical Test for Advanced Learners of English; Lemhöfer and Broersma, 2012). All participants gave written consent for their participation and the re-use of the data prior to the start of the experiments. The study was approved by the Ethics Commission of the University of Zurich.

| ID | LexTALE | Score NR | Score TSR | Speed NR | Speed TSR |
|---|---|---|---|---|---|
| ZKW | 96.25 | 91.67 | 94.84 | 11.73 | 6.14 |
| ZDN | 97.5 | 86.11 | 92.87 | 4.10 | 2.93 |
| ZPH | 97.5 | 94.44 | 97.05 | 7.55 | 2.71 |
| ZMG | 100 | 88.89 | 95.82 | 5.33 | 3.73 |
| ZAB | 100 | 86.11 | 90.42 | 5.14 | 3.32 |
| ZJN | 97.5 | 83.33 | 79.12 | 11.30 | 7.10 |
| ZKH | 81.25 | 83.33 | 93.12 | 6.43 | 5.57 |
| ZGW | 91.25 | 86.11 | 92.14 | 8.06 | 4.17 |
| ZJS | 97.5 | 91.67 | 93.86 | 4.18 | 2.88 |
| ZKB | 100 | 86.11 | 95.33 | 8.43 | 2.48 |
| ZDM | 100 | 80.56 | 96.81 | 5.13 | 3.32 |
| ZJM | 77.5 | 97.22 | 96.56 | 8.73 | 6.30 |
| **mean** | **94.69** | **87.96** | **93.16** | **7.18** | **4.22** |

Table A.12: ZuCo 1.0 Subject demographics, LexTALE scores, and control scores and reading speed (i.e. seconds per sentence) for each task. The * next to the subject ID marks a bilingual subject.

| ID | LexTALE | Score NR | Score TSR | Speed NR | Speed TSR |
|---|---|---|---|---|---|
| YAC | 76.25% | 82.61% | 83.85% | 5.27 | 4.96 |
| YAG | 93.75% | 91.30% | 56.92% | 7.64 | 8.73 |
| YAK | 100.00% | 74.07% | 96.41% | 3.83 | 5.89 |
| YDG | 100.00% | 91.30% | 96.67% | 4.97 | 3.93 |
| YDR | 85.00% | 78.26% | 96.92% | 4.32 | 2.32 |
| YFR | 85.00% | 89.13% | 94.36% | 6.48 | 4.79 |
| YFS | 90.00% | 91.30% | 96.15% | 3.96 | 2.85 |
| YHS | 90.00% | 78.26% | 97.69% | 3.30 | 2.40 |
| YIS | 97.50% | 89.13% | 98.46% | 5.82 | 2.58 |
| YLS | 93.75% | 91.30% | 92.31% | 5.57 | 5.85 |
| YMD | 100.00% | 86.96% | 95.64% | 7.50 | 6.24 |
| YMS | 86.25% | 89.13% | 95.38% | 7.68 | 3.35 |
| YRH | 81.25% | 86.96% | 95.64% | 5.14 | 4.32 |
| YRK | 85.00% | 97.83% | 96.15% | 7.35 | 7.70 |
| YRP | 82.50% | 78.26% | 90.00% | 7.14 | 8.37 |
| YSD | 95.00% | 93.48% | 94.36% | 5.01 | 2.87 |
| YSL | 71.25% | 84.78% | 83.85% | 6.73 | 6.14 |
| YTL* | 81.25% | 80.43% | 94.10% | 7.48 | 3.23 |
| **mean** | **88.54%** | **86.36%** | **91.94%** | **5.84** | **4.81** |

Table A.13: ZuCo 2.0 Subject demographics, LexTALE scores, and control scores and reading speed (i.e. seconds per sentence) for each task. The * next to the subject ID marks a bilingual subject.

*Appendix A.2. EEG Data*

In this section, we present the EEG data extracted from the ZuCo corpus for this work. We describe the acquisition and preprocessing procedures as well as the feature extraction.

*Appendix A.3. Sentiment Reading Task*

ZuCo 1.0 includes a third reading task. We only use this data for the control analyses. Therefore, we describe it here.

For this task, the subjects were presented with positive, negative, or neutral sentences from the Stanford Sentiment Treebank (Socher et al., 2013). The Stanford Sentiment Treebank contains single sentences extracted from movie reviews with manually annotated sentiment labels. We randomly selected 400 very positive, very negative, or neutral sentences (4% of the full

treebank). The 400 selected sentences are comprised of 123 neutral, 137 negative and 140 positive sentences. The sentences were split into two blocks, one for each recording session.

The participants were asked to read the sentences normally. The objective was to analyze the elicitation of emotions and opinions during reading. As a control condition, the subjects had to rate the quality of the described movies in 47 of the 400 sentences. The average response accuracy compared to the original labels of the Stanford Sentiment Treebank is 79.53%.

## Appendix B. Model Parameters

Table B.14 shows the hyper-parameters used in the word-level models presented in Section 3.2.

| Parameter | Value |
|---|---|
| Learning rate | 0.001 |
| LSTM dimension | 64 |
| Dense dimension | 64 |
| Batch size | 40 |
| Epochs | 200 |
| Patience | 104 |
| Min. delta | 0.0000001 |

Table B.14: Hyper-parameters of word-level LSTM model.