

CASES: A Cognition-Aware Smart Eyewear System for Understanding How People Read

ANONYMOUS AUTHORS

The process of reading has attracted decades of scientific research. Work in this field primarily focuses on using eye gaze patterns to reveal cognitive processes while reading. However, eye gaze patterns suffer from limited resolution, jitter noise, and cognitive biases, resulting in limited accuracy in tracking cognitive reading states. Moreover, using sequential eye gaze data alone neglects the linguistic structure of text, undermining attempts to provide semantic explanations for cognitive states during reading. Motivated by the impact of the semantic context of text on the human cognitive reading process, this work uses both the semantic context of text and visual attention during reading to more accurately predict the temporal sequence of cognitive states. To this end, we present a Cognition-Aware Smart Eyewear System (CASES), which fuses semantic context and visual attention patterns during reading. The two feature modalities are time-aligned and fed to a temporal convolutional network based multi-task classification deep model to automatically estimate and further semantically explain the reading state timeseries. CASES is implemented in eyewear and its use does not interrupt the reading process, thus reducing subjective bias. Furthermore, the real-time association between visual and semantic information enables the interactions between visual attention and semantic context to be better interpreted and explained. Ablation studies with 25 subjects demonstrate that CASES improves multi-label reading state estimation accuracy by 20.90% for sentence compared to eye tracking alone. Using CASES, we develop an interactive reading assistance system. Three and a half months deployment with 13 in-field studies enables several observations relevant to the study of reading. In particular, observed how individual visual history interacts with the semantic context at different text granularities. Furthermore, CASES enables just-in-time intervention when readers encounter processing difficulties, thus promoting self-awareness of the cognitive process involved in reading and helping to develop more effective reading habits.

CCS Concepts: • Human-centered computing → Mobile devices.

Additional Key Words and Phrases: Smart eyewear, reading, cognition-aware, eye-tracking, visual attention

ACM Reference Format:

Anonymous Authors. 2022. CASES: A Cognition-Aware Smart Eyewear System for Understanding How People Read. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 0, 0, Article 0 (2022), 30 pages. <https://doi.org/00.0000/00.0000>

1 INTRODUCTION

Reading is a fundamental approach to learning, through which people can expand their vocabulary, gain knowledge, and develop skills. Research has shown a positive relationship between reading and learning; for example, the more people read, the more effectively they improve vocabulary, knowledge levels, and cognitive skills [15]. In fact, reading has long been considered the most important path to lifelong learning, and lifelong readers are generally more successful, both personally and professionally [24, 76].

The science of reading has attracted decades of interest in human-computer interaction (HCI) [27, 81], cognitive science [40, 44], psychology [67], educational psychology [11, 77], cognition and neuroscience [82], pedagogy [38],

Author's address: Anonymous Authors.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

2474-9567/2022/0-ART0 \$15.00

<https://doi.org/00.0000/00.0000>

and brain science [2, 51]. Reading is a cognitive process and understanding it benefits numerous research communities. Studying how people understand the semantics and syntax of text can aid in understanding natural language representation and processing, which are key functionalities of human-level intelligence [35]. Understanding the reading process can also advance the theory of human behavior, thus **benefiting** the domains of applied psychology, pedagogy, and educational psychology. For instance, we can scrutinize human cognitive abilities [36] such as verbal **working**, memory capacity, inhibitory control ability, perceptual speed, and immediate and delayed effects on reading processes. Furthermore, understanding how people read sheds light on reading patterns and strategies, potentially helping readers achieve metacognitive awareness and read more efficiently [21, 58, 84]. In particular, HCI researchers have studied enhancing human reading efficiency [80], reading proficiency [52], reading skills [55], reading comprehension performance [34, 43], and reading outcomes [28].

Reading is a multi-level interactive eye-mind cognitive process. In the short term, readers visually perceive each word, encode it, and mentally assign semantics. In the long term, readers visually perceive a sentence and mentally associate it with context and domain knowledge [39]. Reading can be viewed as a sequence of numerous time-varying states. For instance, some studies explored the state of mind wandering, to detect whether a reader is cognitively engaged or decoupled from the current reading task [19, 54]. Furthermore, some researchers studied the state of having difficulty processing unfamiliar words [33, 72]. However, we note that processing difficulties can present at multiple granularities, e.g., readers may encounter difficulties at the level of a single word, a sentence, or a paragraph. Since it is hard to enumerate all reading states, we focus on the problem of probing the reading cognitive process to detect and explain multiple states at word and sentence levels. Specifically, we investigate whether a reader's mind is wandering, whether the reader is positively engaged, and when comprehension is delayed due to word- or sentence-level processing difficulties.

Eye movements are good indicators to infer the cognitive process [1, 64, 74, 83]. This is based on the eye-mind hypothesis [39], which states that there is a close relationship between where the eyes look and where the mind is engaged. Owing to the fast development of eye-tracking technologies, we can easily access eye-tracking data [3, 50] to explore eye-mind relationships. Numerous researchers have extended the relationship between eye movements and cognitive processes [65, 67]. Also, numerous prevalent methods design eye-tracking reading systems to automatically track the participants' eye movements in a non-intrusive way [16, 38, 72, 73]. These works have summarized some hand-engineered eye movement features to probe the reading cognitive process [72, 73].

However, eye-tracking technologies suffer from a number of shortcomings. The error of commercially available eye-tracking technologies typically ranges from 1 to 4 degrees [46, 60]. Under reading scenarios, this angular accuracy translates to a spatial tracking resolution of about 1.4–2.6 cm. Considering a computerized-reading task where the distance from eye to screen is 40–50 cm, this means that the resolution of the eye tracker is about 3 to 4 lines for a single-spaced document and about 1 to 3 words in the horizontal direction. Such low spatial resolution makes it infeasible to track reading states during word-by-word and line-by-line reading because we cannot locate the words and lines accurately. Previous studies tackle this problem by using an unrealistic setting with a very wide line spacing (e.g., triple-spaced [16]), leaving them unsuitable for use with normally spaced text. In addition, eye-tracking techniques are subject to the inherent transient jitter [9] of human gaze and vertical drift, which require constant calibration [7]. Eye-tracking techniques suited to real-world scenarios have the potential to advance the study of reading.

Furthermore, existing methods ignore contextual influences from text, resulting in less accurate reading state estimation and undermining semantic explanation for these states. Given the same reading context and motivations, the factors influencing reading states mainly pertain to the reading material's and subject's domain knowledge about the content. For example, a good reader may cross-reference previously read text to assist in understanding new and unfamiliar text [33]. In such cases, the high reading frequencies of the earlier text do not necessarily imply that they are difficult. To correctly estimate the current reading state, it

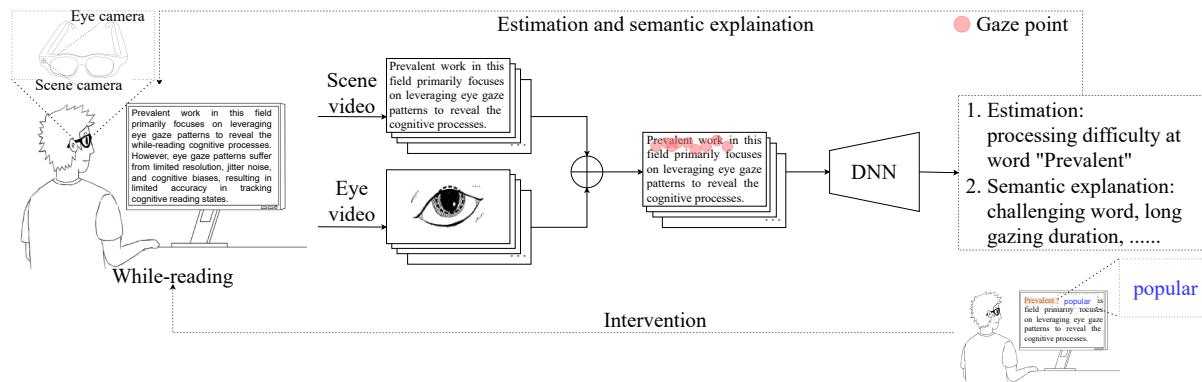


Fig. 1. The proposed CASES smart eyewear system.

is important to be aware of the semantic meaning of the current text, the cross-referenced text, their semantic correlations, and real-time eye gaze patterns. However, it is a non-trivial task to properly fuse the semantics of reading text and eye movements and learn from them in progressive reading scenarios, and it is more challenging to infer semantic explanations for reading state timeseries.

This work aims to provide accurate estimations and semantic explanations for reading state timeseries to support research and outreach efforts in the field of reading science. To this end, we pose the following two research questions (RQ) and posit the corresponding hypotheses.

RQ1: Do readers in the same reading states show different visual attention distributions on the reading text?

Hypothesis 1: Readers in the same reading **state** will show varying visual attention histories (detailed in Section 3), e.g., different total fixation duration, reading times, scanning paths, etc. That is, the visual attention histories of readers in the same reading state differ from each other.

RQ2: When readers are in the same reading states, e.g., encountering difficulty progressing, how does reader visual attention interact with semantic cues in the text?

Hypothesis 2: As indicated by previous studies [19, 73], readers' cognitive effort in processing text is positively related to the difficulty of the text. However, in contrast with previous studies, we further hypothesize that readers can overcome reading difficulties by fetching contextual semantic cues from the surrounding text. When progress is blocked, easy text that is semantically related to difficult text also receives more visual attention and cognitive effort.

The motivation for this work is that the semantic context of text has a direct impact on the multi-level interactive eye-brain cognitive reading process. Leveraging the rich semantic information about reading materials, which can be extracted by advanced natural language processing (NLP) techniques [61, 87], can improve estimation accuracy and provide semantic interpretation of reading states. The semantic information is high-resolution because NLP models can provide semantics at the word level [61, 87]. The inherent hierarchical structure of the semantic information can also be inferred by summarizing the semantics of words to a sentence level. The high-resolution semantic information can compensate for the low-resolution eye movements for more accurate reading state tracking. More importantly, the real-time interaction of eye movements and semantic context can provide semantic explanations for the ongoing reading states.

To this end, we present a **Cognition-Aware Smart Eyewear System (CASES)** capable of measuring reading (**cognitive**) state timeseries. Figure 1 illustrates the workflow of the proposed system. At the heart of CASES is a bi-modal multi-task network named CASES-Net, which takes the bi-modal data, i.e., the eye-tracking and reading

109 text data, as inputs and estimates cognitive reading states in real-time at two granularities: word and sentence
 110 level. To collect high-quality bi-modal data, CASES uses two cameras to record the two required modalities
 111 automatically: **an outward-facing scene camera to capture text** and an inward-facing camera to track gaze points
 112 during reading. CASES is implemented in the form of eyewear to avoid interfering with the reading process
 113 when collecting data. Surveys are deferred until after a reading task is completed, also to avoid interference.

114 CASES-Net employs a four-layer temporal convolutional network (TCN) based module to fuse the two
 115 types of sequential modalities, one of which is informed by semantic information extracted from a pre-trained
 116 NLP model [6, 18]. We treat estimations at two granularities as two distinct but related tasks and propose a
 117 shared convolutional filter mechanism within the TCN to learn the characteristics of the two tasks and their
 118 commonalities. Moreover, we design a multi-task and hierarchical loss function to guide reading state estimation.
 119 To evaluate CASES, we first collect and construct a dataset and then demonstrate that CASES has higher reading
 120 state estimation accuracy than baseline methods. To sum up, CASES-Net combines gaze and semantic information
 121 to better estimate reading states. More importantly, it provides semantic explanations for these reading states.
 122 The well-trained deep model can automatically detect when users encounter reading difficulty without requiring
 123 further inputs (e.g., feedback) from them, thereby limiting potential subjective biases.

124 This work makes the following contributions.

- 125 • We present a cognition-aware smart eyewear system (CASES) to probe and explain human cognitive
 126 processes while reading. CASES aims to support the study of reading and learning to read, as well as
 127 supporting HCI and educational applications investigations on improving reading productivity. The CASES
 128 system is equipped with a deep neural network, CASES-Net, that extracts features pertaining to the visual
 129 attention history and text semantic content. It fuses the two types of features via a shared convolutional
 130 filter mechanism based on TCN to enable accurate reading state estimation at various granularities.
- 131 • CASES is evaluated in real-world contexts. We conduct an ablation study involving 25 participants, in
 132 which CASES delivered superior reading state detection to baseline methods. Specifically, encoding text
 133 semantic content facilitates learning from context cues and improves reading state estimation accuracy.
 134 **Compared with the conventional eye-tracking-only method, we improve accuracy by 20.90% for sentence.**
 135 Furthermore, the text semantic context enables quantitative explanations of reading (cognitive) states.
- 136 • We integrate CASES into a novel interactive reading assistant system. **Three and a half months** of deployment
 137 with 13 in-field studies demonstrate that the integrated system can enable helpful interventions for readers,
 138 thus improving self-awareness in the reading process and helping readers adopt more effective reading
 139 habits.

140 The rest of this paper is organized as follows. Section 2 surveys related work. Section 3 clarifies the key concepts
 141 used in this work. Section 4 details the proposed network and our built real-time reading state detection and
 142 intervention system. Section 5 presents the experimental setups and results. Section 6 presents our findings when
 143 using CASES in practice, general discussion, and future direction. Finally, Section 7 concludes this work.

144 2 RELATED WORK

145 This work is mostly relevant to three broad areas: reading science, eye-tracking in reading, and natural language
 146 processing.

147 2.1 Science of Reading

148 Reading science has attracted decades of interest in various research communities, e.g., HCI, pedagogy, and
 149 educational psychology. These studies primarily deal with the outcomes of reading [28] and reading compre-
 150 hension [43]. **Recently, researchers have studied reading patterns and strategies that improve the efficiency of**
 151 **reading** [21, 27, 58, 84], e.g., interactive reading systems that detect mind wandering during reading [19, 54].

152 They mitigated the negative effect of mind wandering on reading comprehension using just-in-time interventions [19, 54]. Other methods detect words readers do not know automatically [27] and provide appropriate
 153 help [33, 72]. In psychology, applied psychology, and educational psychology, researchers primarily focused
 154 on studying how texts are read and comprehended [11, 62, 67, 77, 83]. For example, Perfetti et al. delivered a
 155 blueprint of reading, consisting of the visual process, representation process that converts visual perception into
 156 a linguistic representation, and operation process on the representation [62]. In cognition science, neuroscience,
 157 and brain science, extensive reading studies focus on developing computational theories of cognition [47]. One
 158 important branch studies the representations and processing of natural languages by the human brain [35]. For
 159 example, Lewis et al. contributed a theoretical framework to explain how verbal working memory supports
 160 sentence processing [47]. Kamide et al. studied how the global and local information in texts impact sentence
 161 processing [40]. Cognitive scientists usually jointly consider language representation and processing [23] based
 162 on the belief that discovering language representation can help answer questions about computation, and vice
 163 versa. Schrimpf et al. provided computationally explicit evidence that language comprehension mechanisms in
 164 human brains are fundamentally shaped by predictive processing through an integrative modeling approach [69].

165 In summary, previous works on the science of reading primarily focus on leveraging eye-tracking during
 166 reading to study the reading process and outcomes. However, they focus less on how individual readers perceive
 167 and process the text in real-time. This study introduces context information from texts to the study of reading
 168 cognitive processes.

170 2.2 Eye-Tracking in Reading

171 Eye-tracking technology can acquire real-time eye movements in a non-intrusive manner [8]. It is natural to
 172 utilize eye movement data to probe the reading process, as the reading process initiates visual input and operates
 173 as an interactive eye-mind cognition process [39]. Over the past decades, numerous studies have focused on
 174 analyzing eye movement data obtained during reading to understand the reading cognitive process and provide
 175 reading assistance [4, 19, 25, 32, 54, 73]. For example, Hyrskykari proposed a gaze-aware reading assistance system
 176 to provide help at the right time without interrupting the reader's thoughts [32]. Cheng et al. proposed a social
 177 reading system, in which they demonstrated that sharing eye gaze annotations generated by experts promoted
 178 reading comprehension for non-experts [10]. Bottos and Balasingam presented an approach to accurately track
 179 the horizontal eye-gaze points in reading scenarios [4]. In addition, there are also many studies focused on
 180 detecting reading behaviors, such as mind wandering [19, 54] or encountering difficulties in comprehending
 181 unfamiliar words [33, 72].

182 In general, these relevant methods have demonstrated that eye movement data helps understand the reading
 183 cognitive process. However, the semantic information of the text, which is closely related to the reading process,
 184 is rarely used in previous studies. This study jointly considers text semantic information and eye movement data
 185 can facilitate understanding the reading process and how readers comprehend texts.

186 2.3 Nature Language Processing

187 Natural language processing (NLP) uses computational techniques to represent and analyze human languages [12]
 188 (see [42] for a comprehensive review). NLP can usually be classified into two categories: natural language under-
 189 standing and natural language generation. As discussed above, this work uses natural language understanding
 190 techniques to obtain semantic contextual information from texts. Successful natural language understanding
 191 techniques can provide generic models for NLP downstream tasks, such as analyzing the association among
 192 text components [18], extracting keywords [6, 71], and analyzing syntax [48]. For example, Linzen et al. pointed
 193 out that, given targeted syntax supervision, a long short-term memory (LSTM) network can learn syntax infor-
 194 mation [49]. Later, they further stated that linguists and neural network researchers might contribute to each

other's areas [48]. Furthermore, NLP neural networks can provide good representations of text; for example, the bidirectional encoder representations from transformers (BERT) model [18], which is based on transformers [79], can obtain state-of-the-art results on several NLP tasks by providing high-quality language representations. Considering the dependency between the masked positions and the discrepancy from pretrain-finetune that BERT neglects, Yang et al. proposed a generalized autoregressive pretraining method to overcome the limitations of BERT [86]. Their pre-trained model, XLNet, outperforms BERT on various tasks. This work builds on recent progress in NLP by using pre-trained NLP models to help understand the reading cognitive process.

202 3 PROBLEM FORMULATION

203 This section clarifies three important concepts used in this work: eye movements, visual attention, and semantic
204 attention.

205 *Eye Movements:* Eye movement patterns can reveal reading strategies and are vital to understanding the reading
206 cognitive process. As shown in existing studies [16, 53], reading generally consists of a series of pauses and rapid
207 shifts in gaze locations. The pauses are called fixation and the shifts are called saccades. These patterns reflect
208 the low-level oculomotor characteristics during reading, typically determined by the physical properties of text,
209 such as the positions or lengths of words.

210 By exploring eye movement patterns, researchers establish connections between low-level eye movement
211 behaviors and higher-level cognitive processes during reading [74]. First, research shows that the direction and
212 duration of eye fixation reveal how the cognitive process unfolds over time [72, 73]. More specifically, fixation
213 locations indicate the attended content, while fixation duration suggests the level of cognitive effort invested by
214 the reader, i.e., longer fixation suggests more effort. Second, the processing time-course of eye movement patterns
215 is widely used to reveal the temporally continuous reading process, which is often linked with comprehending or
216 memorizing. For example, one common temporal reading activity is to move the gaze backward to review the
217 already-read content. In this case, the informative eye movement patterns might be the reading and regression
218 durations, which is also called the second pass [31]. Finally, to alleviate the potential inter-person variations, recent
219 work also designs global features or statistical features based on eye movement patterns to access the reading
220 process, such as the number of saccades, saccade frequencies, and variations in fixation duration [16]. Given
221 the potential ability of eye movement patterns in revealing reading cognitive processes, this work also employs
222 these hand-engineered features as valuable indicators. However, to better suit our case, we first distinguish the
223 representing eye movement patterns at two granularities, and then we re-design them at word and sentence
224 levels, respectively. More details can be found in Section 4.1.3.

225 *Visual Attention:* Although no previous work explicitly defines visual attention in reading scenarios, substantial
226 studies demonstrate a strong correlation between eye movement patterns and attentional processing during
227 reading. For instance, the E-Z reader model [66] posits that attention during reading moves from word to
228 word continuously. The serial-processing assumption states that attention is linked to changes of focus in text
229 processing [26, 56, 85]. Following these studies, our work describes visual attention during reading by establishing
230 the connection between eye movement patterns and the corresponding while-reading text components, such as
231 words and sentences. More specifically, we define the visual attention state to be the collection of eye movement
232 features on each text component. For example, when reading the sentence “*They race to maturity, with the shortest*
233 *generation time of any vertebrate*”, the visual attention for the word “*vertebrate*” consists of fixation duration,
234 reading times, number of fixations, etc. At the sentence level, the visual attention state is defined based on the
235 total dwell time, saccade times, etc.

236 *Semantic Attention:* We are interested in exploring how the semantic meaning from text assists in estimating the
237 time-series reading states and how they explain these states. From this perspective, it is necessary to have a holistic
238 semantic understanding of while-reading texts. Furthermore, such understanding should cover the semantic

meaning of different grain sizes of texts, ranging from single words and sentences to passage levels. This works terms this semantics collection at various granularities as *semantic attention*. For example, semantic attention can hint at whether the while-reading text components are difficult. These difficult components may be unfamiliar or ambiguous words or sentences with complex syntax, which often delay reading. In this case, appropriately using such semantic meaning regarding the difficult score can provide additional evidence in revealing the current reading state and deliver a reasonable interpretation regarding why the current text components block the reading. More details regarding semantic attention can be found in Section 4.1.2.

4 SYSTEM DESIGN

This section describes the CASES system design. We first detail the CASES network (CASES-Net), a deep neural network for detecting and interpreting ongoing reading states. Then, we describe a real-time reading state estimation and intervention system aiming to boost reading comprehension performance.

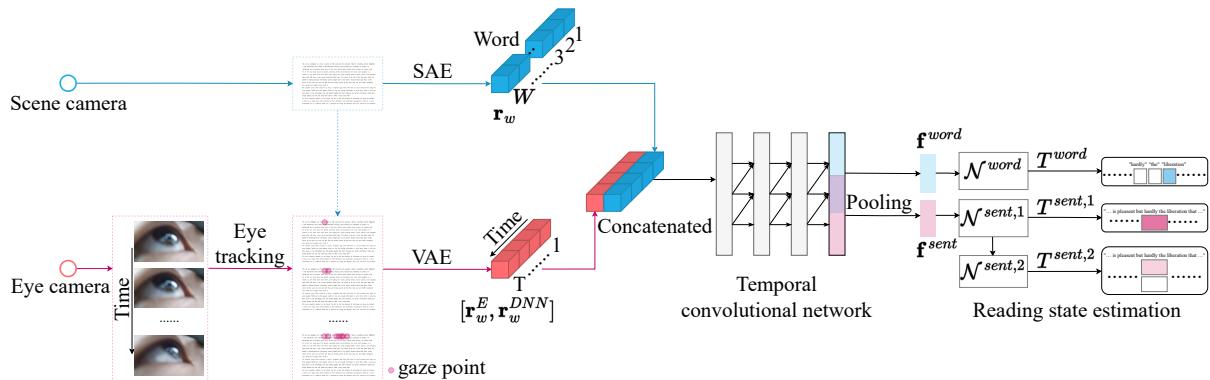


Fig. 2. Framework of the CASES.

4.1 CASES Network

4.1.1 Overall Pipeline. Figure 2 depicts the general framework of CASES-Net. It consists of four modules: semantic attention extraction (SAE), visual attention extraction (VAE), cross-attention extraction (CAE), and reading state estimation/explanation.

The first step in the CASES pipeline provides a comprehensive semantic understanding of the text before the reading begins. This semantic meaning information compensates for the low-resolution eye-tracking data, thus enabling accurate reading state estimation. Semantic meaning also enables explanations during reading state detection tasks in later pipeline stages. To extract semantic meaning, the system turns on the outward-facing scene camera to obtain the text to be read. The SAE module then runs once on the text. It utilizes NLP techniques to extract the high-resolution semantic features and the inherent linguistic structure from the text, thus facilitating subsequent tasks.

Texts contain rich semantic information, but for better individual reading state estimation, personalized visual attention data are also necessary. To capture it, the VAE module is triggered to obtain the online visual attention features corresponding with text components (e.g., while-gazing words or sentences). More specifically, the CASES system senses reader eye images to predict gaze sequences using continuous eye-tracking [46, 60]. Then, the VAE module extracts visual attention features from the sequential gaze data. In parallel, the scene camera records time-aligned scene images to help track gaze positions.

267 Since the obtained semantic meaning of the text and visual attention features are at different spatial resolutions,
 268 we propose the CAE module to properly align them. Here, we use words to segment the visual attention features,
 269 because words are the minimal units considered in this work, upon which sentences and global context depend.

270 The TCN-based network estimates the reading states at word and sentence levels, aiming to explore the task-
 271 specific features for the assistance of the multi-task output. One feature represents the binary determination of
 272 whether a reader has difficulty processing a word; we call this the “word-level task”. The second task is hierarchical
 273 multi-label classification at the sentence level, which includes (Task I) estimating whether a reader is having
 274 sentence-level processing difficulty and if so, (Task II) estimating whether the reader is facing comprehension
 275 challenges, the reader’s mind is wandering, or both. A multi-task and hierarchical loss function for training
 276 guides CASES-Net. We can qualitatively understand the reasons for the predicted reading states by visualizing
 277 the learned semantic attention and visual attention features.

278 The rest of this section explains the technical details of each of the proposed modules.

279 4.1.2 *Semantic Attention Extraction Module.* The aim of the SAE module is to understand the high-resolution
 280 semantic meaning of the document \mathbf{R} , ranging from the word level to the document level. There are two primary
 281 prerequisites for extracting accurate semantic features: obtaining the while-gazing locations and text contents.
 282 The former, i.e., while-gazing locations, can be obtained by using eye tracking and represented as Points of Gaze
 283 (PoG) timeseries. Each PoG corresponds to a two-dimensional coordinate in the scene image recorded by the
 284 scene camera. Given the locations of PoG, we can easily load the while-gazing text contents because the reading
 285 system has already stored all the reading materials in advance. After that, we propose to extract the following
 286 three types of semantic features by utilizing various advanced pre-trained NLP models.

- 287 (1) Each word in \mathbf{R} is encoded as a 768-dimensional vector by XLNet model [86], which can learn the semantic
 288 meaning of the document by processing the whole text passage once. To lower the potential adverse effect
 289 incurred by the high dimensionality, we reduce the XLNet features to 64 dimensions via a fully-connected
 290 (FC) layer and denote them as $\mathbf{r}^B = \{\mathbf{r}_w^B\}_{w=1}^W$, where $\mathbf{r}_w^B \in \mathbb{R}^{64}$ and W is the total number of words.
- 291 (2) To understand the keyword information in the document, we calculate the probability of each word
 292 describing the whole document via the YAKE model [6]. The keyword features are denoted as $\mathbf{r}^K = \{r_w^K\}_{w=1}^W$,
 293 where $r_w^K \in \mathbb{R}$.
- 294 (3) We use word difficulty to assist in the final task of identifying the reading state. Following Franklin et
 295 al. [22], we describe the word difficulty using the length of the word, number of syllables, and familiarity
 296 scored by the MRC psycholinguistic database [14]. We denote the difficulty of words by $\mathbf{r}^D = \{\mathbf{r}_w^D\}_{w=1}^W$,
 297 where $\mathbf{r}_w^D = [l_w, s_w, f_w] \in \mathbb{R}^3$.

298 Finally, each word in the document is represented by the concatenation of the three feature vectors; that is
 299 $\mathbf{r}_w = [\mathbf{r}_w^B, r_w^K, \mathbf{r}_w^D] \in \mathbb{R}^{68}$ ($w = 1, 2, \dots, W$). Note that the semantics regarding more coarse levels (e.g., sentence- and
 300 passage- level) can be generalized from that of the word level, as words are inherently structured and semantically
 301 connected — a passage consists of multiple sentences and a sentence of multiple words.

302 4.1.3 *Visual Attention Extraction Module.* A reliable gaze sequence is the foundation for accurate visual attention
 303 feature extraction. However, the raw gaze points are noisy due to difficult-to-avoid human motion and limited
 304 eye-tracking resolution. To alleviate this issue, we design a filtering algorithm to smooth the raw gaze points,
 305 leveraging their sequential characteristics. More specifically, we first employ an existing eye-tracking technology
 306 to estimate the PoGs and record the PoGs sequences as $\mathbf{E} = \{\mathbf{e}_t\}_{t=1}^T$, where T is the total number of timestamps
 307 considered. The designed filtering method first uses median filtering to discard outliers due to gaze jitter. Then,
 308 we use mean filtering to stabilize the fluctuations of sequential PoGs due to the limited eye-tracking resolution.
 309 After filtering, we obtain the smoothed PoGs $\mathbf{E}^* = \{\mathbf{e}_t^*\}_{t=1}^T$. We segment each word and sentence using \mathbf{E}^* and
 310 then send them to the next step for visual attention extraction.

311 The number of PoGs will increase rapidly during reading. To reduce the size of PoGs, experts have engineered
 312 a large number of representative features reflecting how people comprehend characters during reading [16, 31,
 313 72, 73] or whether they are disengaged from reading [19, 54]. In this work, we propose to further enrich the
 314 engineered visual features. The following features are widely used to describe word-level processing state while
 315 reading: fixation duration, number of fixations, and number of repeated word readings. However, we observe that
 316 these three features vary not only person-to-person but also during reading. Such variation significantly affects
 317 estimation performance. The personal variation is usually removed by normalizing personal data [29]; however,
 318 the latter while-reading variation is rarely considered. This work introduces local information to tackle the latter
 319 problem: every τ seconds, we add the statistical features to describe the mean and the variance of each engineered
 320 feature, for $E^* = \{e_t^*\}_{t=1}^\tau$, to describe the visual attention for each word. In total, we obtain a 9-dimensional feature
 321 for each word. Moreover, we normalize the four sentence-level representative visual features, including dwell
 322 time [17], saccade times [20], forward saccade times [59], and backward saccade times [59], using the sentence
 323 length, so these features better describe the local variation. Given that we segment M words during τ , each word
 324 is represented using $r_w^E \in \mathbb{R}^{(9+4)}$ ($w = 1, 2, \dots, M$). There are nine word-level features and four sentence-level
 325 features that are identical to the words in the same sentence.

326 Lastly, we propose to use the deep features of the sequential gaze data, as recent studies have demonstrated
 327 the effectiveness of deep neural networks (DNN) on eye movement classification. We adopt the existing feature
 328 extractor based on the 1D-CNN with BLSTM backbone [75] (denoted as N_{eye} to extract 8-dimensional deep
 329 features during time duration τ , i.e., $r_w^{DNN} \in \mathbb{R}^8$ ($w = 1, 2, \dots, M$), and discard the classifier.

330 4.1.4 *Cross-Attention Extraction Module.* To facilitate downstream multi-task learning, the cross-attention
 331 extraction (CAE) module first fuses the two modalities, then explores the commonalities to predict at different
 332 granularities and the distinct task-specific information.

333 Before fusing the two modalities, we use the following strategy to time synchronize them. Specifically, for
 334 each smoothed PoGs sequence e_t^* , we identify the M words being processed at time t , and concatenate the two
 335 features vectors to obtain $f_t^w = [r_w, r_w^E, r_w^{DNN}] \in \mathbb{R}^{(68+13+8)}$ as the overall representation of the two modalities.
 336 For all other words w' that have not been visually processed till time t , we pad the semantic attention feature
 337 vector $r_{w'}$ with a zero vector, i.e., $f_t^{w'} = [r_{w'}, \mathbf{0} \in \mathbb{R}^{21}]$. In this way, the word being processed at time t can
 338 be properly described semantically with its corresponding visual attention features. In contrast, the unread words
 339 padded with zeros are given less attention.

340 The CAE module uses a Temporal Convolutional Network (TCN) model, which can capture temporal de-
 341 pendencies. Specifically, the module uses temporal convolutional filters/kernels to process input sequences.
 342 Each filter calculates a weighted average in the time domain, and the parameters of the filters are learned to
 343 optimize the objective function. Each TCN layer consists of temporal convolutions, a non-linear Relu activation
 344 function, and a max pooling function or an upsampling function. The CAE module has four TCN layers. To learn
 345 different tasks more efficiently, the filters of the last layer are divided into task-specific filters, namely word-level
 346 filters/sentence-level filters, and task-shared filters, namely common filters. The features extracted by word-level
 347 filters and common filters are used for word-level tasks. The features extracted by sentence-level filters and
 348 common filters are used for sentence-level tasks.

349 4.1.5 *Reading States Estimation and Explanations.* After obtaining the cross-attention features, we are ready to
 350 detect the reading state of “processing difficulty”. We have the following three tasks. (1) Word-level binary-class
 351 classification task T^{word} : The word-level features are fed to a fully connected layer to predict whether the reader
 352 finds the word being processed difficult. Sentence-level and word-level tasks differ. Since we know that mind
 353 wandering may co-occur with reading difficulty for a sentence, we formulate the task at the sentence level in the
 354 following hierarchical fashion. (2) Sentence-level binary-class classification task $T^{sent,1}$: With the sentence-level

355 features, we first determine whether the reader is in a normal reading state without any processing difficulties
 356 using a binary classifier. (3) Sentence-level multi-label classification task $T^{sent,2}$: If the reader enters into an
 357 abnormal state, the reader can be either mind wandering or processing difficulty, or both; This is a multi-label
 358 classification task, where multi labels can be assigned simultaneously; label 1 is mind wandering and label 2 is
 359 processing difficulty.

360 Finally, to train the network, we propose the following loss function reflecting the performances of all tasks:

$$\mathcal{L} = \mathcal{L}(T^{word}) + \mathcal{L}(T^{sent,1}) + \mathcal{L}(T^{sent,2}). \quad (1)$$

Binary Cross Entropy (BCE) loss is used for T^{word} and $\mathcal{L}(T^{word})$ is illustrated as follows

$$\mathcal{L}(T^{word}) = -\frac{1}{W} \sum_{w=1}^W \left(y_w^{word} \log p_w^{word} + (1 - y_w^{word}) \log(1 - p_w^{word}) \right),$$

where W denotes the number of word; y_w^{word} denotes the label of word w , $y_w^{word} = 0$ indicates the reader finds the word w easy, $y_w^{word} = 1$ indicates the reader finds the word w difficult; p_w^{word} is the word-level estimation results given by the network N^{word} .

BCE loss is also used for $T^{sent,1}$ and $\mathcal{L}(T^{sent,1})$ is illustrated as follows

$$\mathcal{L}(T^{sent,1}) = -\frac{1}{S} \sum_{s=1}^S \left(y_s^{sent,1} \log p_s^{sent,1} + (1 - y_s^{sent,1}) \log(1 - p_s^{sent,1}) \right),$$

where S denotes the number of sentences; $y_s^{sent,1}$ denotes the binary classification label of the s th sentence, $y_s^{sent,1} = 0$ indicates the reader is in a normal reading state without any processing difficulties for sentence s , $y_s^{sent,1} = 1$ indicates the reader is in an abnormal reading state; $p_s^{sent,1}$ is the sentence-level binary classification estimation results given by the network $N^{sent,1}$.

For sentences with $y_s^{sent,1} = 1$, to solve the multi-label problem, BCE loss is used for each label separately and the loss of $T^{sent,2}$ is illustrated as follows

$$\mathcal{L}(T^{sent,2}) = -\frac{1}{\sum_{s=1}^S \mathbf{1}(y_s^{sent,1} = 1)} \sum_{s=1}^S \sum_{l=1}^L \mathbf{1}(y_s^{sent,1} = 1) \left(y_{s,l}^{sent,2} \log p_{s,l}^{sent,2} + (1 - y_{s,l}^{sent,2}) \log(1 - p_{s,l}^{sent,2}) \right),$$

361 where $L = 2$ denotes the number of labels, i.e. label 1 as mind wandering and label 2 as processing difficulty;
 362 $y_{s,l}^{sent,1}$ denotes the supervised information of the l th label for sentence s , $y_{s,l}^{sent,1} = 1$ indicates sentence s has the
 363 l th label, $y_{s,l}^{sent,1} = 0$ indicates sentence s does not have the l th label; $\mathbf{1}(\cdot)$ is an indicator function, $\mathbf{1}(y_s^{sent,1} = 1) = 1$
 364 when $y_s^{sent,1} = 1$, $\mathbf{1}(y_s^{sent,1} = 1) = 0$ when $y_s^{sent,1} = 0$; $p_{s,l}^{sent,2}$ is the sentence-level multi-label estimation results
 365 given by the network $N^{sent,2}$.

366 4.2 EYEReader: A Real-Time Reading State Detection and Intervention System

367 Our goal is to determine reading state series that influence reading fluency and mitigate the negative effects of
 368 reading processing difficulties. To this end, we build a real-time reading state detection and intervention system
 369 (called EYEReader) for English language. For the convenience of readers, EYEReader is implemented in the form
 370 of a website, enabling cross-platform compatibility.

371 This section first gives a concrete example to show the key features of EYEReader and how to use it. Then it
 372 details the system architecture, along with its operation pipeline. At last, it describes the hardware prototype.

When I speak of ‘defining’ reference and the various propositional attitudes, I am not, of course, thinking of defining an “analytic” definition, one which analyzes the ‘concept’ or ‘meaning’ of ‘refers to’, ‘believes’, ‘desires’, and so on. Nor am I longer being concerned with the question of whether there are ‘semantically’ or ‘consciously’ reducible to physicalistic (or computational) predicates. The question is whether there are semantic and propositional-attitude properties and relations which are ‘reducible’ to physical-cum-computational properties and relations in the way in which (to use a familiar example) the temperature of an ideal gas is reducible to mean molecular kinetic energy.

When we say that temperature has [] an molecular kinetic energy, we are claiming more than just that ‘temperature’ [] ‘equal’ [] ‘mean molecular kinetic energy’. These two magnitudes could be [] ‘coextensive’ [] ‘equal’ [] ‘mean molecular kinetic energy’ differed only in a class of physically possible cases which were not ever produced by investigators and which did not occur spontaneously in nature.

(a)

A problem in the radical interpretation of a language spoken by beings more sophisticated than we. How would this be determined using beliefs of our culture such that ‘everyone knows them, everyone knows that everyone knows them, and so on’?

These future speakers might believe that snow is white on the basis of a quantum mechanical calculation.

Imagine they think differently than us, and need more sensory input or time to process it before reaching a mental state.

white ‘with co-extensibility’ will not apply to [] ‘coextensive’ [] ‘equal’ [] ‘mean molecular kinetic energy’.

beliefs in a different way than we, but that the sensory stimulations that suffice for us do not suffice for most of them, at least without lengthy calculation, they do not go immediately from the sensory stimulations to the ‘mental state’.

Imagine they think differently than us, and need more sensory input or time to process it before reaching a mental state.

(b)

What we have just seen to hold in the case of my (former) account also holds for Lewis's account. His account too says that to have psychological states is just to be a ‘model’ of a certain theory. He too needs to restrict the notion of ‘model’ to those cases where the theory in question is ‘metaphysically’ true. Lewis does this by requiring the readers to have the metaphysical property he calls ‘naturalness’ or else. But this means that his theory requires that each organism to which folk psychology applies be a model for folk psychology in just the sense that my account in ‘Philosophy and Physical States’ does. That is, that each organism possesses one physical state per propositional attitude.

You just mind wandered; please reveal the missed content.

I can exploit my arguments to make a further point. Folk psychology cannot play the explanatory role Lewis wishes it to play unless it is true that each organism possesses one physical state per propositional attitude which is ‘natural’ in the sense that it is not ruled out by any belief that the organism has about its own future.

That is, that such an organism possesses one physical state per propositional attitude.

I can exploit my arguments to make a further point. Folk psychology cannot play the explanatory role Lewis wishes it to play unless it is true that each organism possesses one physical state per propositional attitude which is ‘natural’ in the sense that it is not ruled out by any belief that the organism has about its own future.

That is, that such an organism possesses one physical state per propositional attitude.

(c)

Fig. 3. Screenshots of three intervention examples. Detection and Interventions: at the word level (left), (b) at the sentence level (middle), and (c) on mind wandering (right).

373 4.2.1 ***Key Features and Operation Process of EYEReader.*** We first give some key features of EYEReader, and we
374 then use a concrete case to show the automatic detection and intervention process.

375 **Key feature 1: text materials selection.** The text materials should contain various topics, as the intervention
376 is anticipated to be text-agnostic. We select 36 reading comprehension materials with diverse topics from an
377 English qualification test to match the participants’ reading comprehension ability. Each article has around 450
378 words on average. Users can log in to the system, select their preferred articles from existing materials, and start
379 reading by simply clicking a button.

380 **Key feature 2: friendly reading interface.** Because we have overcome the limited resolution issues when
381 eye-tracking is used during reading scenarios, the interface of text presentation of EYEReader is similar to
382 common computerized reading settings. More specifically, articles are automatically divided into several different
383 pages (around 240 words per page) with a regular line height, approximately single-spaced. We adopt an 18-point
384 default font typeface.

385 **Key feature 3: intervention design.** The interventions are designed to help users overcome the three
386 while-reading processing difficulties, i.e. mind wandering, challenging words and complex sentences, that may
387 lead to a negative impact on their reading comprehension performance. Three interventions are designed for
388 the three difficulties respectively: 1) providing an immediate reminder once mind wandering is detected, which
389 reminds readers to focus on the current reading; 2) simplifying the challenging words; and 3) streamlining the
390 complex sentences. We provide the following three examples to further clarify how the interventions support
391 reading.

- 392 (1) **Simplifying challenges words.** Once the system detects that a reader is facing a challenging word, it
393 highlights the word in blue and provides a more comprehensible one in the pop-up window. For example,
394 when a user struggles with “coextensive”, the pop-up window offers a more straightforward and easy-
395 to-understand one, “equal”. Figure 3 (a) provides a screenshot of this intervention. After receiving the
396 interventions, users can click on the highlighted words to hide the reminders and continue reading.
- 397 (2) **Streamlining complex sentences.** The procedure of streamlining complex sentences is similar to that of
398 simplifying challenging words. Differently, the system highlights the complex sentences in red and provides
399 simpler sentences in the pop-up window. For example, for a long and complex sentence, “Imagine that they
400 not only come to the belief in a different way than we, but that the sensory stimulations that suffice for us
401 do not suffice for most of them, at least without lengthy calculation, they do not go immediately from the
402 sensory stimulations to the ‘mental state.’” the system provides a relatively more straightforward version:
403 “Imagine they think differently than us, and need more sensory input or time to process it before reaching
404 a mental state.” A screenshot of this type of intervention is depicted in Figure 3 (b). Users can also click on
405 the highlighted sentences to hide the pop-up window and continue reading.

406 (3) **Giving mind wandering reminder.** When the system detects that the reader is distracted while reading, it
 407 highlights the missed content in yellow and displays a pop-up message in the centre of the screen, showing
 408 that “You just mind wandered; please reread the missed content.” A screenshot of this type of intervention
 409 is depicted in Figure 3 (c). The pop-up message automatically fades out after one second.

410 Participants' eye gazes are calibrated prior to their reading in order to correlate the two cameras equipped in
 411 the eyewear. The calibration method follows Pupil Capture¹ [41]. Specifically, during the calibration phase, the
 412 participants wear the eyewear and sit in front of the computer, a pupil calibration marker appears on the screen
 413 with fixed locations. The participant is instructed to gaze at the marker for approximately two seconds. The same
 414 procedure is executed for the other four calibration markers on the screen. In this way, the system would record
 415 these positions to correlate the two cameras.

416 During the reading process, readers wear the prototype eyeglass and sit in front of the computer to read. The
 417 pre-trained CASES-Net model is always-on to automatically detect potential abnormal reading states, i.e., whether
 418 the user is struggling with difficult words or complex sentences, or their mind is wandering. When abnormal
 419 events that affect reading are detected, the system triggers interventions automatically. The text components will
 420 be highlighted, and the corresponding treatments will be shown on the right-top of the text content automatically
 421 in a pop-up window.

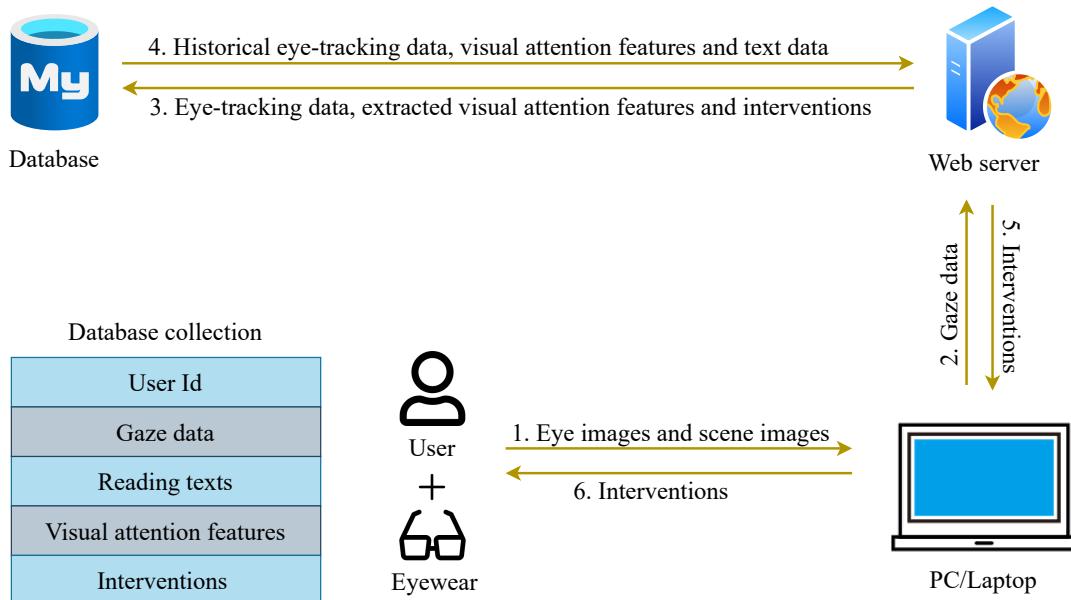


Fig. 4. The architecture of the reading state detection and intervention system.

422 4.2.2 *System Architecture.* Figure 4 illustrates the overall architecture of EYEReader. We use the Vue.js framework
 423 to develop the front-end website, while we choose Django for the back-end of the website, as it is a widely-used
 424 Python web framework [5]. Django offers a variety of third-party tools for building communication between the
 425 front-end and back-end efficiently following the REST API specification. To store and manage the data on the
 426 server, we adopt one of the widely-used open-source database management systems – MySQL [57]. **The eyewear**

¹<https://docs.pupil-labs.com/core/software/pupil-capture/#calibration>

427 and the PC used for reading are in the same LAN(Local Area Network). The eyewear is running on Android
 428 12. We developed a service app without user interfaces to read the real-time video stream recorded by the two
 429 cameras and push the video stream to the PC via the RTSP protocol². On the PC edge, we receive the coming
 430 video stream from the eyewear using the RTSP protocol. The received video is then handled by Pupil Capture
 431 and Pupil Service provided by Pupil Labs³.

432 Next, we describe the overall operation workflow of the built intervention system and show how it provides
 433 just-in-time interventions for users encountering reading processing difficulties. There are mainly six steps
 434 described below.

- 435 • Step 1: During system operation, EYEReader loads the pre-trained CASES-Net from the server when
 receiving the requests from the front end.
- 436 • Step 2: The recorded eye/scene images captured by eyewear are pushed to the user's PC for eye-tracking
 using the Pupil Capture [41].
- 437 • Step 3: The tracked gaze points are sent to the server for further visual attention feature extraction.
- 438 • Step 4: The server loads the historical eye-tracking data, visual attention features, and texts to decide
 whether it is the right time to intervene.
- 439 • Step 5: Once processing difficulties are detected, the estimation results are returned to the front end for
 triggering interventions. The corresponding treatment is shown at the front end to facilitate the current
 reading.
- 440 • Step 6: After that, the current interventions and all other data are saved in Database.

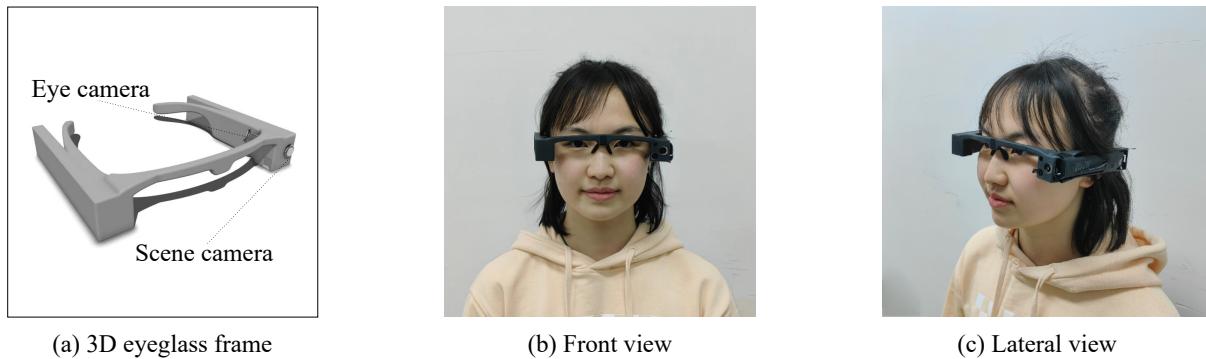


Fig. 5. Hardware prototype of CASES eyewear.

446 4.2.3 *Hardware Design*. We design prototype eyewear and integrate CASES-Net into the eyewear, as eyewear is
 447 a natural way to be used in various reading scenarios.

448 We presume that the eyewear will be well-migrated to various reading scenarios. Therefore, we adopt a
 449 stand-alone scheme to integrate the computing components and power supply into the headset frame. Figure 5
 450 shows the eyewear hardware prototype.

451 The eye-tracker follows the Pupil³, and we make slight adjustments to suit our case. More specifically, we use
 452 Qualcomm Snapdragon 865 platform directly integrated into the left leg of the eyewear. The eye camera and
 453 scene camera modules are replaced with 20 MegaPixels (MP) Samsung S5K3T2 and 64 MP Samsung S5KGW1,

²https://en.wikipedia.org/wiki/Real_Time_Streaming_Protocol

³<https://docs.pupil-labs.com/core/diy/>

454 respectively. The eye camera is used to record eye videos to perform eye tracking. The scene camera senses scene
 455 videos to capture the text being read. We design the 3D eyeglass frame to fit the two cameras into the left leg
 456 of the mounting frame. To balance the weight of the headset, the battery is integrated into the right leg of the
 457 eyewear.

458 5 EVALUATIONS

459 This section describes experiments to evaluate CASES, the cognition-aware eyewear system for estimating
 460 reading states. We first detail the experimental setup, data collection, and evaluation measures. We present
 461 results and quantify the technical capabilities of CASES. All experimental procedures are approved by the ethical
 462 committee at our University.

463 5.1 Evaluation Methodology

464 5.1.1 *Experimental Setup.* We recruited 25 participants by posting a questionnaire at our university campus.
 465 Specifically, we distributed informed consent forms for the participants before the experiment started. In the
 466 consent form, we informed them about the purpose of our study and the procedure of the experiments and gave
 467 them the option to withdraw at any time during the experiments. All participants signed the informed consent
 468 form. After completing their sessions, the subjects received either local currency equivalent to 14 dollars or a
 469 thank-you gift worth approximately 14 dollars for their participation. A summary of the participant demographics
 470 follows.

- 471 • **Age:** 22–28 years old with an average age of 23.5,
- 472 • **Gender ratio:** 19 males (76.0%) and 6 females (24.0%).
- 473 • **Native/non-native speaker:** 5 native speakers (20.0%) and 20 non-native ones (80.0%).

474 As shown in Figure 6 (a), the participant wears eyeglasses and sits in front of the computer to read. While reading,
 475 we record videos using the eye camera and time-aligned videos using the scene camera.

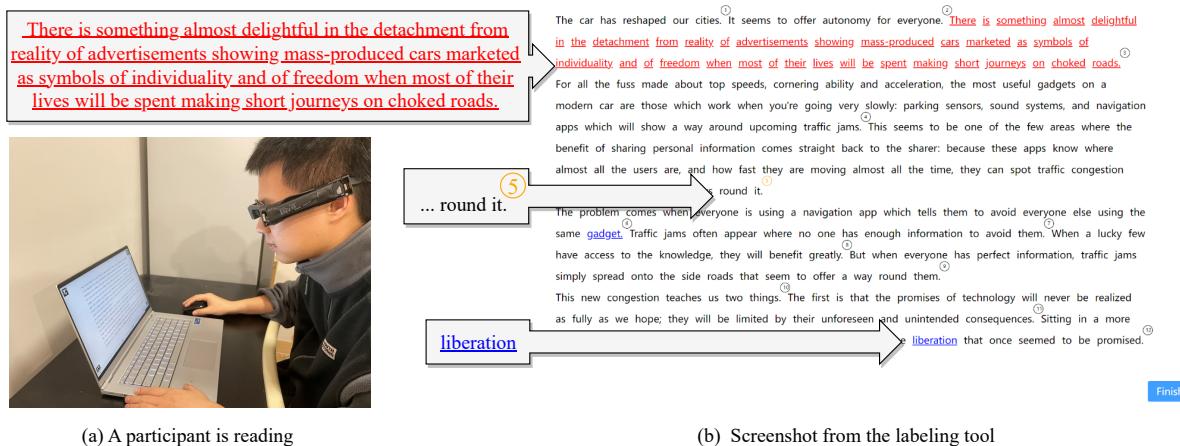


Fig. 6. The in-lab setting of CASES experimental study.

476 5.1.2 *Text Material Selection.* Texts should cover a wide range of subjects so readers can enter multiple reading
 477 states. Moreover, each text should be short, allowing participants to read several texts. This study selects 36
 478 articles with the following three subjects:

479 **Subject matter 1** One-minute BBC world news⁴: 10 articles with approximately 300 words per article on
 480 average.

481 **Subject matter 2** English qualification tests⁵ : 16 articles containing reading comprehension materials
 482 with approximately 450 words per article on average.

483 **Subject matter 3** Philosophy related [63]: 10 articles with approximately 500 words per article on average.

484 The first two of these provide challenging words and sentences, respectively. The third may lead to mind
 485 wandering. We anticipate that most participants are unfamiliar with the third subject matter, and it is hard to
 486 understand the content without prior knowledge. The idea of mundane subject selection to introduce mind
 487 wandering follows a recent work [54].

488 Considering the diverse backgrounds and prior knowledge of various participants, texts should also cover a
 489 wide range of subject classes. According to Dewey Decimal Classification (DDC) method [70], we categorize
 490 the selected articles into ten subject classes, including “social science”, “religion”, and eight other subjects. Prior
 491 to data collection, we select an approximately equal number of articles from each topic class, except for the
 492 philosophy articles.

493 5.1.3 *Dataset*.

494 (1) *Data Collection*. The CASES requires time-aligned eye gaze data and text data (i.e., the words or sentences
 495 being read) to detect reading states. In addition, the synchronized data should capture continuous reading, during
 496 which users may encounter various reading states. To the best of our knowledge, there are no publicly available
 497 datasets suitable for our problem. Therefore, we first develop an online system to collect data meeting our
 498 requirements. To facilitate research, we release the collected dataset, which is online available.⁶

499 Prior to data collection, we select an approximately equal number of articles from each topic class, except
 500 for the philosophy articles. Then, these articles are randomly assigned to each participant. Then, each article
 501 is divided into pages. There are around 240 words per page in single-spaced 18-point typeface. After that, we
 502 randomly select articles from each topic for the participants to ensure that they cover all three subject matters.
 503 This design allows most participants to encounter numerous reading states. Each article is read by an average of
 504 five participants. We verbally instruct participants on how to use the data-collection system, such as navigating
 505 to the next/previous page. Finally, each participant reads the texts. Reading one article takes approximately six
 506 minutes.

507 (2) *Ground-Truth Labeling*. After completing an article, the participant is immediately instructed to label their
 508 reading states. We developed a labeling tool with a GUI window to accelerate labeling. Participants can review
 509 each page of the article. On each page, they use a single click to label the words they cannot comprehend and
 510 use a double-click to label sentences they do not comprehend. We also provide a button at the right top of each
 511 sentence for users to mark whether their minds wandered when reading it. The annotated words and sentences
 512 are highlighted in different colors so users can quickly double-check their annotations, as shown in Figure 6 (b).
 513 Annotating one article takes around three minutes. In total, the data collection process, including the annotation
 514 collection, took us approximately fourteen days.

515 (3) *Dataset Statistics*. The collected dataset is randomly split into training (80%) and test (20%) sets per par-
 516 ticipant/article. The total numbers of labels for “word-level processing difficulties”/“sentence-level processing
 517 difficulties”/“mind wandering” are 1005/244/200.

⁴<https://www.bbc.com/news>

⁵<https://cet.neea.edu.cn/>

⁶<https://drive.google.com/drive/folders/1AZmL1YhUU49ZOmCJKxqWsQUVFIW5nedo?usp=sharing>

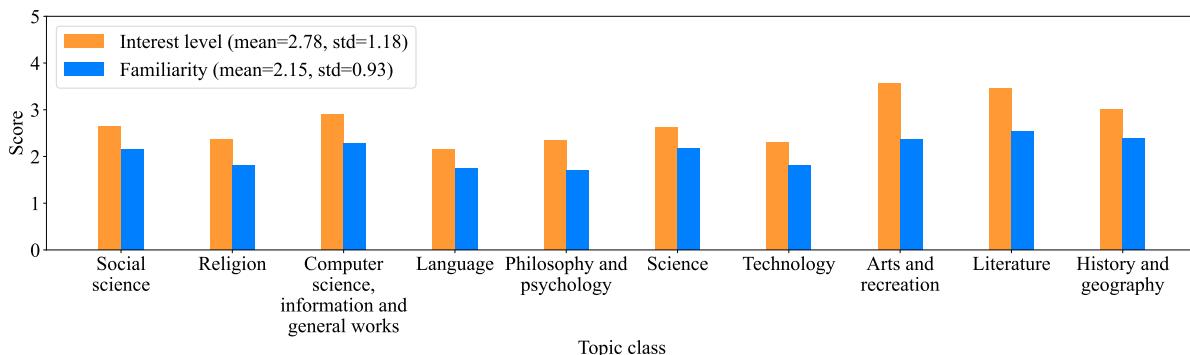


Fig. 7. Score of interest level and familiarity for each topic class.

518 We survey the participants' interest level and familiarity with the ten topics to further verify the fairness of the
 519 selected topics, i.e. we expect that the interest level and familiarity are evenly distributed across all topics. Using
 520 the Likert scale⁷, we asked participants to score their interest level and familiarity with the articles they read.
 521 The scale ranges from 1 to 5; a higher score indicates that a participant is more familiar with or more interested
 522 in the article. As shown in Figure 7, participants gave roughly similar interest scores (mean = 2.78, std = 1.18) and
 523 familiarity (mean = 2.15, std = 0.93) on the ten topic classes, indicating that the ten topic classes have covered
 524 individual participants evenly. The means of interest level and familiarity of all topics are around 2.5, suggesting
 525 that the topics are intermediate to participants.

526 We also visualize the distribution of the average number of labels per article at various levels of interest and
 527 familiarity in Figure 8 (left). It is clear that readers give approximately the same number of labels per article
 528 under each interest level. Also, Figure 8 (right) shows that the number of labels per article decreases as familiarity
 529 increases. This is in line with our intuition, as participants often give more labels for their unfamiliar articles.

530 **5.1.4 Evaluation Metrics.** Because our framework is hierarchical and multi-task, we need to adopt appropriate
 531 measures to evaluate each task. The first task is binary classification of whether a reader is facing difficulty
 532 processing a word. We evaluate its performance using accuracy and the receiver operating characteristics (ROC)
 533 curve. The second task is hierarchical multi-label classification at the sentence level, which includes sentence-level
 534 Task I and Task II. Task II is multi-labeled, following previous work [88]. We therefore use the multilabel-based
 535 macro-averaging metric, i.e., averaged-accuracy and ROC curve, to evaluate it.

536 **5.1.5 Baseline methods.** We conduct ablation studies to evaluate CASES, as there is no prior work solving the
 537 problem addressed in this work, thus making direct comparisons with prior work infeasible. We use the following
 538 three baseline methods for evaluation.

539 (1) Visual: Previous studies have demonstrated that some reading states, such as mind wandering, can be
 540 identified using gaze-relevant features [19, 54], which are closely related to our work. To validate whether
 541 the eye-relevant features are sufficient for reading state recognition at multiple text element granularities
 542 (words and sentences), this work uses baseline method leveraging 13 eye-relevant features (9 word-level
 543 features and 4 sentence-level features described in Section 4.1.3) to identify the state while reading. We
 544 use the support vector machine (SVM) method to conduct the three classification tasks: word-level task,
 545 sentence-level Task I, and sentence-level Task II. This work adopts SVM as it has been successfully applied

⁷https://en.wikipedia.org/wiki/Likert_scale

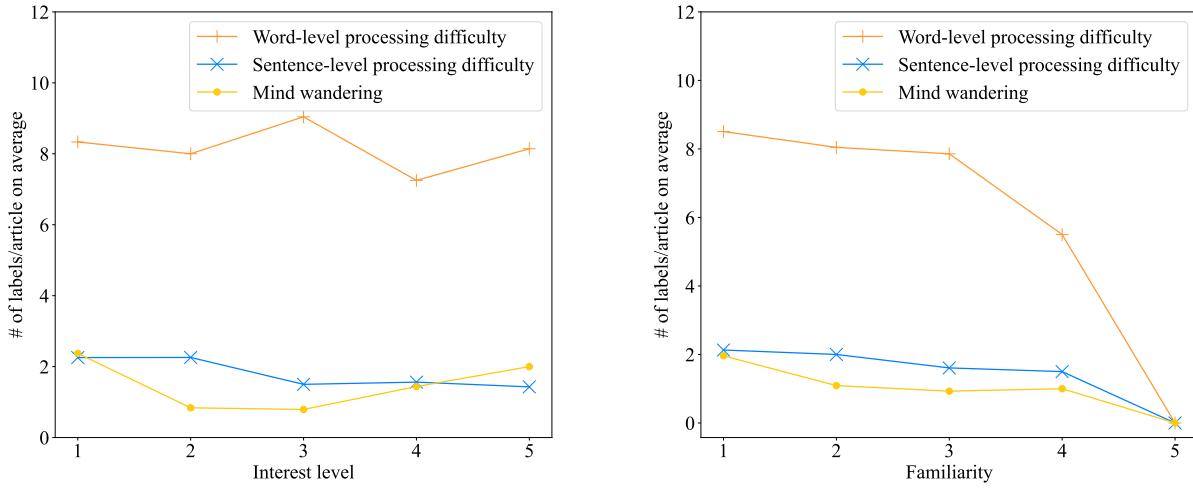


Fig. 8. Number of labels per article on average under each interest level and familiarity.

to various classification tasks [78], and is one of the widely used methods in similar tasks [20, 54]. For simplicity, we refer to this method as *Visual*.

- (2) *Visual+*: Eye movement patterns are good indicators for reading state recognition. Inspired by prior work [75] that leverages deep neural network (DNN) to achieve accurate eye movement pattern identification, we use the 8-dimensional deep features extracted from a deep neural network (1D-CNN with BLSTM [75]) to improve the accuracy of reading state estimation. To make a fair comparison, the extracted deep features are concatenated with the 13 expert-designed features and sent to the CAE module to estimate reading state. This baseline method is an improved version of the *Visual* method called *Visual+*.
- (3) *NLP*: *Visual* and *Visual+* identify reading states based solely on visual attention features. To verify the classification performance based on the semantic content of texts, we designed this baseline method, dubbed *NLP*. As in the *Visual+* method, we first extract semantic features using the SAE module and then send the extracted features to the CAE module to infer reading states.

5.2 Results

5.2.1 Overall Performance. Figure 9 shows the reading state recognition performance of our methods and three baseline methods. The proposed method achieves the best performance. Compared with the *Visual* method, i.e., conventional eye-tracking only, CASES improve the accuracy by 6.85%, 8.55%, 20.90% for the word-level task and the sentence-level Task I and Task II. Furthermore, compared with the baseline method *Visual+* and *NLP*, CASES has superior reading state estimation. For example, the sentence-level Task II detection accuracy of CASES is 86.64% while it is 79.15% or lower for the baseline methods. We conclude that using context derived from text improves reading state estimation.

We plot the Receiver Operating Characteristic (ROC) of different methods. Figures 10a, 10b, and 10c demonstrate that CASES outperforms the baseline methods in Area Under the Curve (AUC), which is one of the most widely used performance measures in classification or retrieval problems.

The next section further explains why CASES outperforms the baseline methods and how it offers semantic explanations of the predicted reading states.

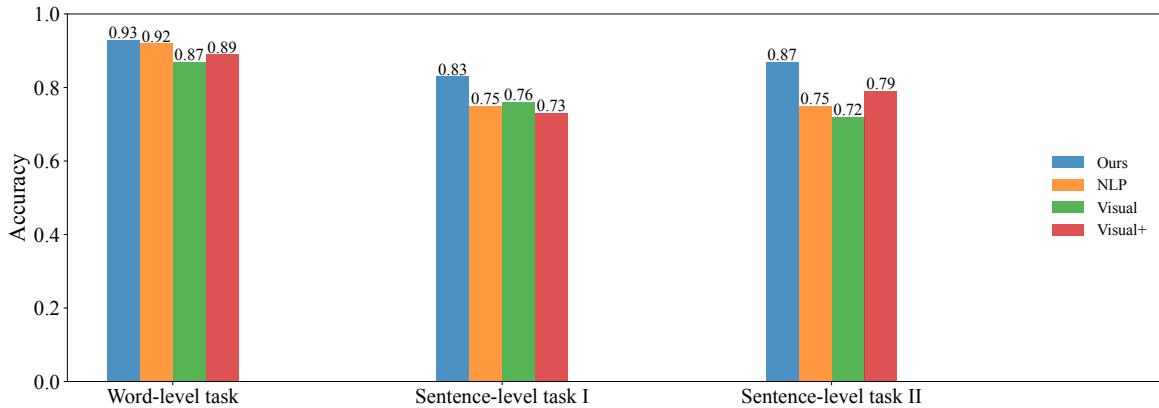


Fig. 9. Reading state classification accuracy for CASES and the baseline methods.

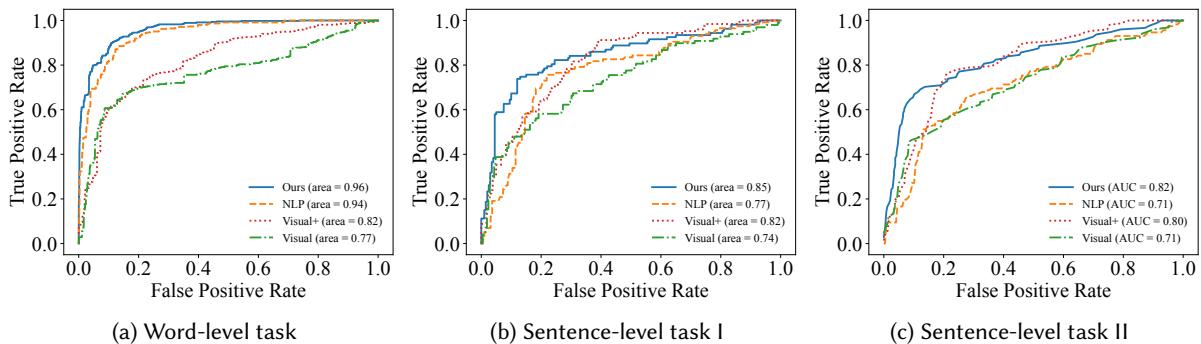


Fig. 10. ROC for CASES and the baseline methods.

571 6 PILOT STUDY

572 This work aims to study progression through cognitive states while reading to assist our understanding of the
 573 reading process. To this end, we have conducted in-field pilot studies using CASES, the proposed system, **for**
 574 **totally three and a half** months. This section first summarizes the initial findings around our the designed two
 575 RQ and hypotheses using CASES. Then, it demonstrates the capability of EYEReader to make helpful real-time
 576 interventions when reading difficulties are encountered. Finally, it describes the limitations of our system and
 577 indicates possible extensions of this work.

578 6.1 The Procedure of the Pilot Study

579 We recruited **thirteen** volunteers to participate in our pilot study from our University. The average age is **23.9**
 580 **years (SD=1.6, min=22, max=28)**, with **n=4 (30.8%) female and n=9 (69.2%) males**. There are **10 non-native readers**
 581 **(76.9%) and 3 native ones (23.1%)**. The non-native participants reported that they have passed the college English
 582 test and the native readers are college-level students at our university.

583 During the pilot study, participants wore the prototype eyeglasses, sat in front of the computer, and logged in to
 584 the development website to read. The participants can either take the eyeglasses with them and use the eyeglasses
 585 whenever they would like to do the experiments, or come to our laboratory for the experiments. Participants are
 586 encouraged to use the system whenever they read, as reasonable observations require the prolonged engagement
 587 of participants.

588 The pilot study lasts three and a half months and consisted of two stages. During the first stage, we require
 589 participants to label the words and sentences they encountered difficulty processing, and these labels are treated
 590 as ground truth. Based on the qualitative evaluation [68], we examine the labelled data point by point at different
 591 granularities around the designed RQ and hypotheses. Then, we make several findings on how people read at
 592 different granularities, i.e., single words and sentences, and summarized the following six patterns to discuss. The
 593 second stage focuses on applying EYEReader in practice. At the end of the pilot study, each participant completes
 594 a survey of their opinions on the usability and value of EYEReader. Finally, we confirm the proposed hypotheses.

595 6.2 Key Observations

596 6.2.1 *Observations at the Word Level.* In this section, we present three observations on how users read at the
 597 single-word level.

598 *Observation I: Users comprehend the lexical meanings of words by directing their gazes more frequently toward*
 599 *material they find difficult to process.* When users encounter difficulty processing a word, they usually gaze at
 600 it longer, and more times than typical. This observation is consistent with prior evidence about the process of
 601 comprehending single words during reading [16, 53]. Figure 11 illustrates one example of this observation, where
 602 participant P6 has difficulty comprehending the meaning of “mitigate” and “debris”. P6 fixates “mitigate” (fixation
 603 label 13) and “debris” (with fixation label 19) for a long time and reads them more than two times. In particular,
 604 P6 has the longest fixation duration on the word “debris” and he has the most reading times on the word “debris”
 and “mitigate”.

Officials at the White House announced a new space policy focused on managing the increasing number of satellites that companies and governments are launching into space. Space Policy Directive-3 lays out general guidelines for the United States to mitigate the effects of space debris and track and manage traffic in space. This policy sets the stage for the Department of Commerce to take over the management of traffic in space.

Officials at the White House announced a new space policy focused on managing the increasing number of satellites that companies and governments are launching into space. Space Policy Directive-3 lays out general guidelines for the United States to mitigate the effects of space debris and track and manage traffic in space. This policy sets the stage for the Department of Commerce to take over the management of traffic in space.

Fig. 11. Visualization of visual attention for P6 reading a sentence. Each circle represents a fixation point. The larger the area of the circle, the longer the fixation duration. The circle number denotes the timestamp of the fixation time. *Top:* Raw text; *Bottom:* Text with filtered point of gazes

605
 606 Figure 12 provides another example of this observation for a native participant. Participant P12* (* indicates
 607 native user here in after) is facing challenging words “counterbalanced” (with fixation label 10) and “sketch”

608 (fixation labels 19 and 21). Under the text context presented in Figure 12, we observe that P12* has the most
 609 prolonged fixation duration on words “sketch” and “counterbalanced”.

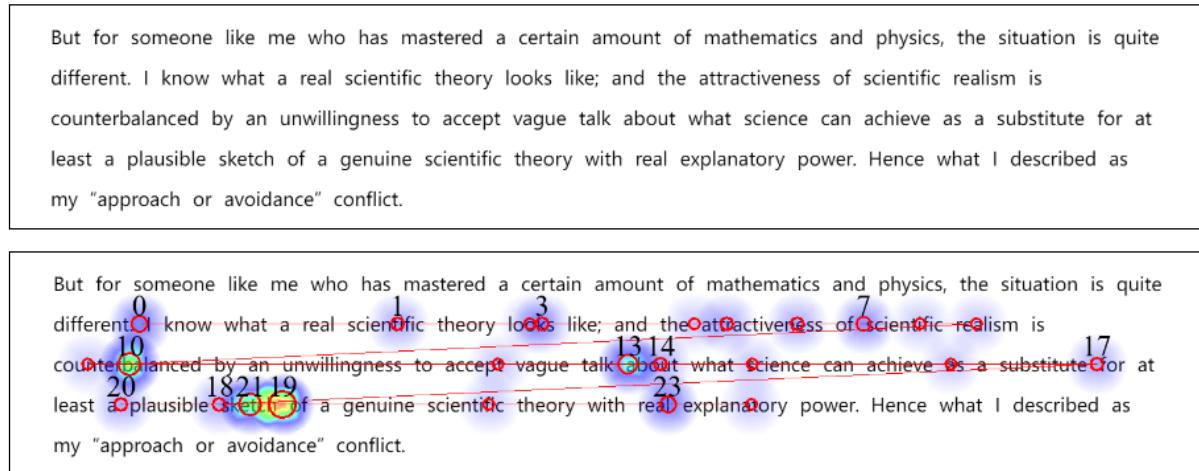


Fig. 12. Visualization of visual attention for P12* reading a sentence. *Top:* Raw text; *Bottom:* Text with filtered point of gazes.

610 *Observation II: When a user encounters difficulty processing a word, the user first directs their gaze to the word and*
 611 *then to other words to examine the semantic context.* Readers generally avoid breaking their chain of thinking by
 612 stopping when a difficult word is encountered, especially when the word does not affect their understanding of
 613 the text. However, when readers consider a difficult word to be highly topic-relevant or meaningful for subsequent
 614 text comprehension, they tend to interrupt their reading and attempt to deduce the semantic meaning of the
 615 word from its semantic context. This observation differs from a previous study [16] and our next observation
 616 complements it.

617 *Observation III: When users examine the semantic context of a difficult-to-process word, they gather semantic*
 618 *clues by shifting their gazes to different locations even when considering the same difficult word, from the same text,*
 619 *under similar reading conditions.* Readers typically attempt to find an appropriate location in the text to help
 620 comprehend the current difficult-to-process word. The text at the location should reveal the relevant information
 621 about the difficult word. Also, that location varies from person to person, depending on their current cognitive
 622 states about the context.

623 Figure 13 shows the proportion of the three above observations for each participant by summarizing their past
 624 experienced processing difficult words. We observe that the ten non-native participants experience Observation I
 625 in most cases (around 87.47% cases on average), and they fall into Observation II & Observation III in fewer times,
 626 i.e., around 12.53% on average. In contrast, the native participants experience Observation II & Observation III in
 627 most cases (around 60.67% cases on average), and they fall into Observation I fewer times, i.e., around 39.33% on
 628 average. This aligns with our intuition, as we anticipate that native readers are more adept at leveraging the
 629 context cues from texts to help their reading comprehension.

630 Figure 14 shows an exemplary case to provide further insights on Observation II and Observation III. Here two
 631 readers, P2 and P5, face the same reading difficulty in comprehending the word “liberation” when they read the
 632 same sentence from the same article. We can see that the two participants first direct their visual attention to the
 633 target word, “liberation” where the fixation labels are 14 and 13 for P4 and P5, respectively. They then shift their
 634 gazes. Participant P2 gazes back at the previously read word, “pleasant”, while Participant P5 gazes forward to

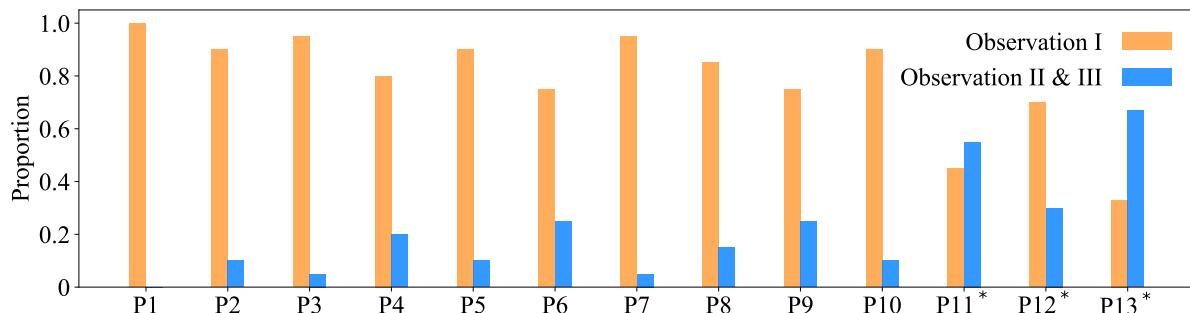


Fig. 13. Profile of word-level comprehension failures for thirteen readers. * indicates native participants.

the word, “promised”. Both of these words are semantically relevant to the difficult word, “liberation”, as shown in Figure 14 (top row).

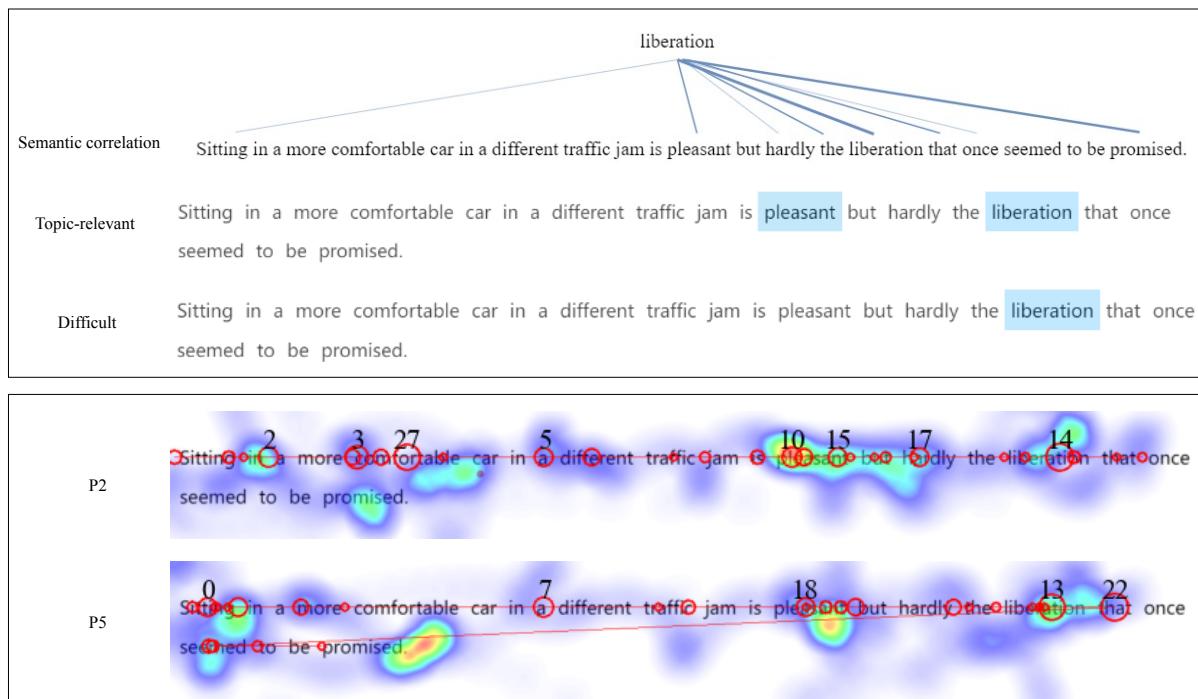


Fig. 14. An example in which two non-native users both have difficulty understanding the word, “liberation”.

Figure 15 provides an example for two native readers, P12* and P13*. They face the same reading difficulty in comprehending the challenging word “realism”. Clearly, P12* and P13* first direct their gaze to “realism” with fixation labels 14 and 11, respectively, and then shift their gazes. P12* gazes forward to an antonym word “antirealism”. Differently, P13* gazes back at the already-read word “internal”, a modifier of the target word.

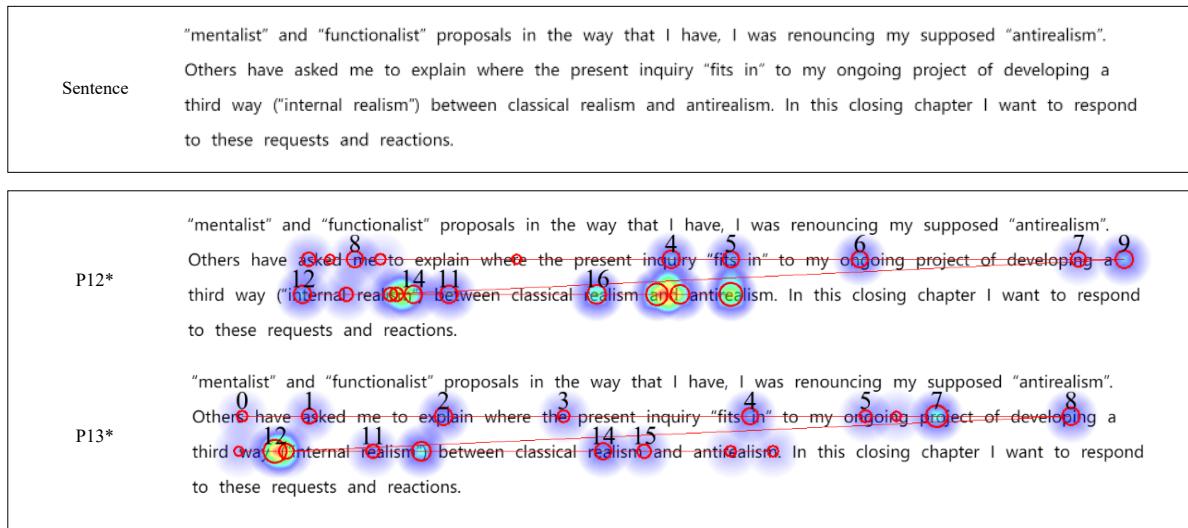


Fig. 15. An example in which two native users both have difficulty understanding the word, "realism".

641 6.2.2 *Observations at Sentence Level.* This section focuses on two modes of comprehending sentences: interpretive
642 (semantic) and structural (syntactic).

643 *Observation IV: People incrementally comprehend the semantics of a sentence as they read each word, while with*
644 *different gaze time series.* Figure 16 shows the inter-reader differences in gaze time series for the same sentence.
645 P1 focuses on the first parts of sentences (with more fixation, labels 0–12) while P4 focuses on other parts of
646 sentences (fixation labels 9–11).

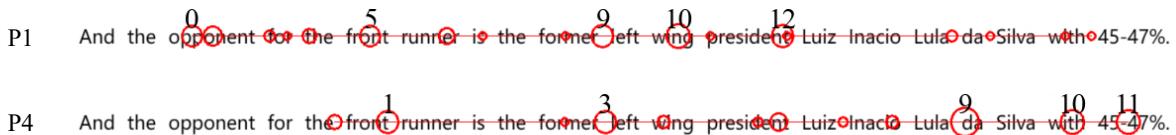


Fig. 16. An example of two **non-native** users in the same sentence but with different visual attention.

647 Figure 17 shows the inter-readers differences for two native readers (P11* and P13*), in the gaze timeseries.
648 They both read the sentence sequentially but with different visual focuses. P13* focuses on the first parts of the
649 sentence, e.g., with more distinct locations of focus, while P11* focuses on other parts of sentences.

650 *Observation V: Readers enter the "rereading" or "reanalysis" state at different times when having difficulty with*
651 *the same sentence, as illustrated in Figure 18.* P8 backtracks 3–4 words (with a fixation label starting from 28)
652 when reading the middle of the sentence, and then continues reading the sentence; while P3 rereads the sentence
653 from the beginning when reading the middle of the sentence (fixation label 12).

654 Figure 19 depicts such an example for native users. The two users both face challenges in comprehending
655 the sentence "As countless boards and and overall performance." We observe that P11* rereads the sentence
656 (fixation label 23) right after completing the first-pass reading (fixation label 22); while P13* rereads the sentence
657 from the beginning of the sentence (fixation label 6) when finishing the middle of the sentence (fixation label 5).

P11* 0 have been arrested in relation to the unrest which reportedly started after an India-Pakistan cricket match.

P13* 0 have been arrested in relation to the unrest which reportedly started after an India-Pakistan cricket match.

Fig. 17. An example of two native users in the same sentence but with different visual attention.

P8 The authors, who focus on education in England, offer a number of sensible recommendations, some of which are an attempt to alleviate the uninspiring and fact-based approach to education that has crept into policy in recent years. When children are regarded as vessels to be filled with facts, creativity does not prosper; nor does

P3 The authors, who focus on education in England, offer a number of sensible recommendations, some of which are an attempt to alleviate the uninspiring and fact-based approach to education that has crept into policy in recent years. When children are regarded as vessels to be filled with facts, creativity does not prosper; nor does

Fig. 18. An example of two non-native users having different “reread” behaviors when in the same reading state: encountering comprehension difficulties on the sentence.

P11* But these provisions create difficulties for businesses when applied to highly paid managers and executives. As countless boards and business owners will attest, constraining firms from firing poorly performing, high-earning managers is a handbrake on boosting productivity and overall performance. The difference between C-grade and

P13* But these provisions create difficulties for businesses when applied to highly paid managers and executives. As countless boards and business owners will attest, constraining firms from firing poorly performing, high-earning managers is a handbrake on boosting productivity and overall performance. The difference between C-grade and

Fig. 19. An example of two native users having different “rereading” behaviors when in the same reading state: encountering comprehension difficulties on the sentence.

658 *Observation VI: Different people “reread” the same sentence with different reading states.* Figure 20 shows two
 659 non-native users reading the same sentence twice. P1 gets distracted (i.e., enters the mind wandering state) during
 660 the first reading of the sentence (typical fixation labels 2, 6, and 14); therefore, P1 spends more time and has more
 661 fixations on the sentence in the second reading (fixation labels 21, 24, 26, 28) than in the first pass. In contrast, P4
 662 spends more time when reading the sentence the first time (fixation labels 3, 17, and 18), but he quickly skims it
 663 the second time (fixation labels 26 and 37).

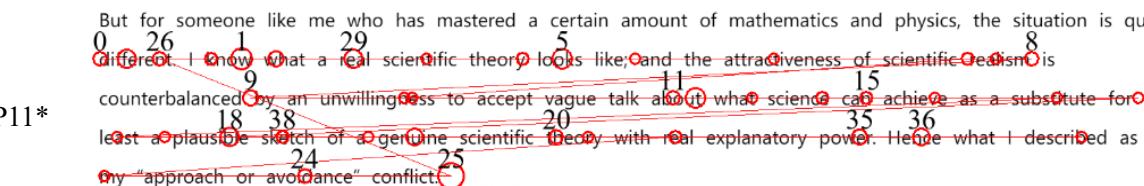
664 Figure 21 shows two native users, P11* and P12*, reading the same sentence twice. They label their reading
 665 states as sentence-level processing difficulty and mind wandering, respectively. P11* spends more time and more
 666 fixations when reading the sentence in the first pass (typical fixation labels 1, 5, and 18) than in the second
 667 pass (typical fixation labels 26, 29, and 36). Also, P11* rereads the sentence after completing the next sentence

P1  policy sets the stage for the Department of Commerce to take over the management of traffic in space.

P4  This policy sets the stage for the Department of Commerce to take over the management of traffic in space.

Fig. 20. An example of two **non-native** users “rereading” the same sentence with different reading states. P1 gets distracted during the first reading and he then rereads the sentence.

668 (fixation label 25). Differently, P12* gets distracted during the first pass of the sentence (typical fixation labels
 669 5–15); therefore, P12* rereads the sentence with more time and has more fixations on the sentence in the second
 670 pass (fixation labels 18, 21, 26 and 27).

P11*  But for someone like me who has mastered a certain amount of mathematics and physics, the situation is quite
 different. I know what a real scientific theory looks like; and the attractiveness of scientific realism is
 counterbalanced by an unwillingness to accept vague talk about what science can achieve as a substitute for at
 least a plausible sketch of a genuine scientific theory with real explanatory power. Hence what I described as
 my “approach or avoidance” conflict.

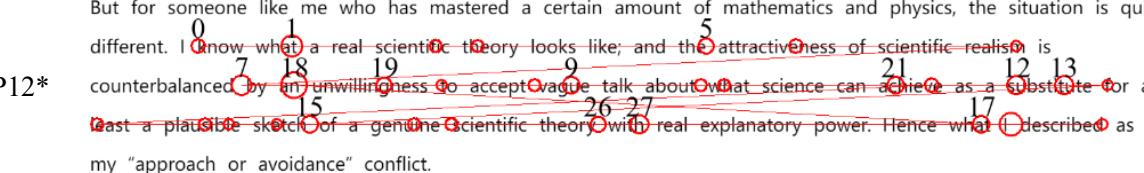
P12*  But for someone like me who has mastered a certain amount of mathematics and physics, the situation is quite
 different. I know what a real scientific theory looks like; and the attractiveness of scientific realism is
 counterbalanced by an unwillingness to accept vague talk about what science can achieve as a substitute for at
 least a plausible sketch of a genuine scientific theory with real explanatory power. Hence what I described as
 my “approach or avoidance” conflict.

Fig. 21. An example of two native users “reread” the same sentence with different reading states. P11* encounters processing difficulty with the sentence, and then P11* rereads the sentence; P12* gets distracted during the first reading, and then P12* rereads the sentence.

671 6.3 Evaluation of EYEReader in Practice

672 Section 5 shows that CASES-Net accurately detects reading states. This section evaluates the ability of EYEReader
 673 to promote reading comprehension by detecting reading states implying processing difficulties and making
 674 real-time interventions.

675 To make quantitative assessment of EYEReader, we define reading comprehension improvement as $(s_{past} -$
 676 $s_{present})/s_{past}$, where $s_{present}$ and s_{past} denote the number of challenging words or sentences at present and in
 677 the past, respectively. The higher the $(s_{past} - s_{present})/s_{past}$, the higher the reading comprehension improvement.
 678 This definition is used to identify challenging words and sentences. After pilot studies, we ask participants to
 679 indicate whether they still face challenges in comprehending these words and sentences. Figure 22 shows the
 680 results. All thirteen participants have positive reading gains, which means EYEReader is effective in helping users
 681 to overcome unfamiliar words and complex sentences.

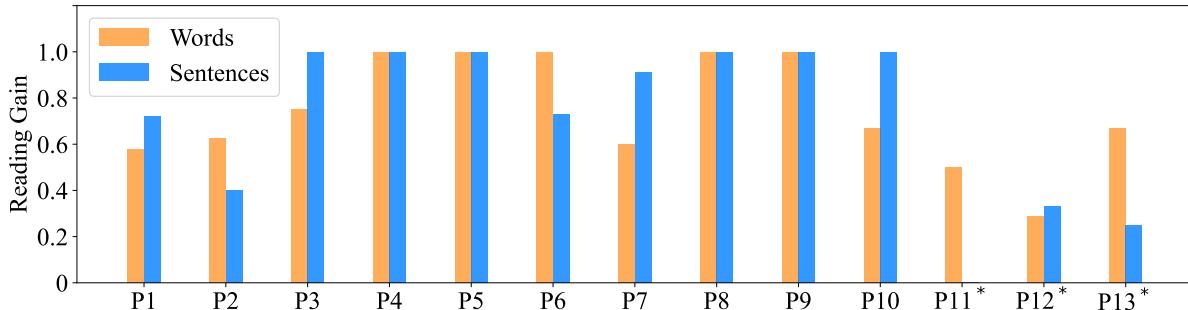


Fig. 22. Profile of reading gain for 13 participants in pilot studies. * indicates native readers.

682 6.4 Feedback from Participants

683 We designed several open-ended questionnaires to qualitatively evaluate EYEReader. Thirteen questionnaires
 684 were sent to participants, twelve of which were returned. Among them, 10/12 of the participants positively
 685 commented on word-level intervention. They believe that fine-grained intervention at the word level can precisely
 686 pinpoint the reading difficulties they are experiencing. There are 9/12 of the participants found sentence-level
 687 intervention helpful. In particular, when facing challenging sentences with complex syntactic structures, it was
 688 difficult to comprehend the sentence even though they were familiar with all the words. In this case, EYEReader
 689 helped them overcome this reading difficulty by highlighting the sentence and explaining it. In addition, 9/12 of
 690 the participants found EYEReader valuable in reminding them when their minds wandered; these participants
 691 stated that they usually do not realize when they are distracted. Timely reminders can make their reading more
 692 focused and efficient.

693 Furthermore, we collect participants' opinions regarding whether the eyewear hardware will negatively affect
 694 their reading process. Through conducting a questionnaire, we ask the participants in the pilot study to give
 695 scores for the comfortable level of hardware on a scale of 1-5; the corresponding description is listed in Table 1. A
 696 total of 13 questionnaires are sent out, and 12 are recalled. Statistical results show that participants generally
 697 think the hardware has a media or negligible impact on their reading process (mean = 3.33, std = 0.75).

698 In addition, once the trained CASES is applied to practice, we design the proper system intervention so as not
 699 to break the chain of thoughts of users. That is, the intervention can help readers avoid interrupting reading due
 700 to the encountered processing difficulties that lead them to seek help from other means [32], such as a dictionary.
 701 Therefore, we design the intervention process with minimal interaction cost and encourage readers to focus
 702 on the current reading. To assess the impact of the intervention on reading, we also conduct a questionnaire
 703 to collect readers' opinions regarding the user-friendliness of intervention interaction. Similarly, we ask the
 704 participants in the pilot studies to give scores on a scale of 1-5; the corresponding description is listed in Table 2.
 705 Statistical results show that most participants generally think the intervention process is user-friendly (mean =
 706 3.92, std = 0.76).

Table 1. Rating scale regarding the comfort of the hardware that does not have a negative impact on the reading process.

Score	Description
1	Severe impact
2	Significant impact
3	Neutral
4	Negligible impact
5	No impact totally

Table 2. Rating scale used to describe whether the intervention is user-friendly that does not affect the reading process negatively.

Score	Description
1	Very unfriendly
2	Unfriendly
3	Neutral
4	Friendly
5	Very friendly

708 6.5 Discussion and Future Work

709 CASES has the goal of accurately estimating and providing semantic explanations of reading states over time,
710 which can facilitate the scientific study of reading by enabling a deeper understanding of the cognitive processes
711 involved in learning to read, disentangling the complex combination of cognitive skills and their impact on
712 reading fluency, and measuring the efficacy of methods for teaching reading and beneficial reading habits.

713 This section first revisits the proposed research questions and hypotheses. Then, it briefly discusses the potential
714 future works that will improve CASES.

715 **6.5.1 Revisiting Research Questions and Hypotheses.** We confirm the hypotheses for the two presented research
716 questions (RQ) based on the results and observations, which we detail below.

717 **RQ1: Do readers in the same reading states show different visual attention distributions on the reading text?**

718 Confirming hypothesis 1: Readers in the same reading state do show varying visual attention histories. As
719 inter-person variation, i.e., individual difference, is ubiquitous, the visual attention histories of readers in the
720 same reading states indeed differ from each other, which can be found from Observation II, Observation III,
721 Observation IV, Observation V, and Observation VI.

722 **RQ2: When readers are in the same reading states, e.g., encountering difficulty progressing, how does reader visual
723 attention interact with semantic cues in the text?**

724 Conforming hypothesis 2: When readers encounter the same processing difficulties, they shift their visual
725 attention to the surrounding text to fetch contextual semantic cues. In other words, when readers' reading
726 progress is blocked, easy text that is semantically related to complex text also receives more visual attention and
727 cognitive effort, which can be found from Observation II and Observation III.

728 6.5.2 Discussion and Future Work.

729 **(1) Science of reading.** This work investigates the human cognitive reading process by exploring the com-
730plementarity of eye movements and text. However, it is also important to integrate illustration information to
731 understand how people read. A recent study has shown that text-diagram instructions can improve reading
732 comprehension [37]. Thus, our future work aims to exploit semantic information, including text and illustrations,
733 and integrate them with eye movements. In addition, we aim to investigate more reading states that might provide
734 a complete picture of the reading cognitive progress. In addition to determining the reading states at the word
735 and sentence levels, it would be valuable to measure how people read at the entire passage level. This could
736 deepen our understanding of how people summarize and reflect on learned knowledge during reading.

737 **(2) Interactive Reading System.** Our system is still an early-stage prototype. A longer user study would
738 enable the collection of more data and user feedback to improve the interactive design and user experience.

739 This could help us to build a mature reading assistance system that contributes to educational applications, HCI
 740 studies, etc.

741 **(3) Reading Contexts.** We would like to emphasise that we presume that the system will be well-migrated
 742 to various reading scenarios, and therefore, we use the eyeglasses form to study reading. We believe wearing
 743 eyeglasses to read is a portable way in numerous reading contexts, including computerized reading and physical
 744 reading (e.g., reading newspapers). However, since our eyeglasses are still in the early prototype stage, in this
 745 work, we did not experimentally cover all the scenarios. The system presented in this work is currently used in
 746 a computerized-reading context, as reading using electronic devices has become common in our modern lives
 747 and has been widely studied by a large body of researchers [13, 30, 54]. We are aware that investigating the
 748 physical-reading context is also important, and we are interested in applying our eyeglasses to investigate the
 749 reading (cognitive) states under this context in our future work.

750 **(4) Brain-Sensing Methods in Reading.** In addition to eye-tracking in reading, brain-sensing via electroen-
 751 cephalograph (EEG) can determine the level of cognitive workload under different rapid serial visual presentation
 752 settings, as demonstrated in [45]. It can be utilized to determine the cognitive workload or attention of texts at
 753 different granularity levels. However, this has to be done with the eye movement data jointly to accurately locate
 754 the positions of text being read and allow fine-grained analysis on processing difficulty of words. We believe that
 755 it is a direction that is worth exploring in the future to further improve the performance of our system.

756 7 CONCLUSIONS

757 This work describes CASES, a cognition-aware smart eyewear system that automatically recognizes (cognitive)
 758 reading state timeseries using eye tracking based visual attention and text semantic context. Ablation studies
 759 demonstrate that CASES significantly improves the accuracy of reading state recognition over the conventional
 760 approach using only eye tracking. Furthermore, in-field studies enabled several observations about how individual
 761 reading state timeseries are related to text semantic context at different granularities. The ability to track semantic
 762 context cues enables better understanding of progressive reading states. We embodied CASES in an interactive
 763 reading assistant system, which provides just-in-time interventions when reading difficulty is encountered.
 764 A two-month deployment demonstrates the benefits of the system in promoting self-awareness of cognitive
 765 processes while reading and helping to improve reading habits. We believe CASES will be of use in the scientific
 766 study of reading, cognition, and human-computer interfaces.

767 REFERENCES

- 768 [1] Ugo Ballenghein, Johanna K Kaakinen, Geoffrey Tissier, and Thierry Baccino. 2020. Cognitive engagement during reading on digital
 769 tablet: Evidence from concurrent recordings of postural and eye movements. *Quarterly Journal of Experimental Psychology* 73, 11 (2020),
 770 1820–1829.
- 771 [2] Gregory S Berns, Kristina Blaine, Michael J Prietula, and Brandon E Pye. 2013. Short-and long-term effects of a novel on connectivity in
 772 the brain. *Brain connectivity* (2013), 590–600.
- 773 [3] Yevgeni Berzak, Chie Nakamura, Amelia Smith, Emily Weng, Boris Katz, Suzanne Flynn, and Roger Levy. 2022. CELER: A 365-participant
 774 corpus of eye movements in L1 and L2 English reading. *Open Mind* (2022), 1–10.
- 775 [4] Stephen Bottos and Balakumar Balasingam. 2019. Tracking the progression of reading through eye-gaze measurements. In *2019 22th
 776 International Conference on Information Fusion (FUSION)*. IEEE, 1–8.
- 777 [5] Carl Burch. 2010. Django, a web framework using python: Tutorial presentation. *Journal of Computing Sciences in Colleges* 25, 5 (2010),
 778 154–155.
- 779 [6] Ricardo Campos, Vítor Mangaravite, Arian Pasquali, Alípio Jorge, Célia Nunes, and Adam Jatowt. 2020. YAKE! Keyword extraction from
 780 single documents using multiple local features. *Information Sciences* 509 (2020), 257–289.
- 781 [7] Jon W Carr, Valentina N Pescuma, Michela Furlan, Maria Ktori, and Davide Crepaldi. 2022. Algorithms for the automated correction of
 782 vertical drift in eye-tracking data. *Behavior Research Methods* 54, 1 (2022), 287–310.
- 783 [8] Benjamin T Carter and Steven G Luke. 2020. Best practices in eye tracking research. *International Journal of Psychophysiology* 155
 784 (2020), 49–62.

- 785 [9] Yuhu Chang, Yingying Zhao, Mingzhi Dong, Yujiang Wang, Yutian Lu, Qin Lv, Robert P. Dick, Tun Lu, Ning Gu, and Li Shang. 2021.
 786 MemX: An Attention-Aware Smart Eyewear System for Personalized Moment Auto-Capture. *Proc. ACM Interact. Mob. Wearable
 787 Ubiquitous Technol.* 5, 2, Article 56 (June 2021), 23 pages. <https://doi.org/10.1145/3463509>
- 788 [10] Shiwei Cheng, Zhiqiang Sun, Lingyun Sun, Kirsten Yee, and Anind K Dey. 2015. Gaze-based annotations for reading comprehension. In
 789 *Proceedings of the 33rd annual ACM conference on human factors in computing systems*. 1569–1572.
- 790 [11] Michelene TH Chi and Ruth Wylie. 2014. The ICAP framework: Linking cognitive engagement to active learning outcomes. *Educational
 791 psychologist* (2014), 219–243.
- 792 [12] KR1442 Chowdhary. 2020. Natural language processing. *Fundamentals of artificial intelligence* (2020), 603–649.
- 793 [13] Virginia Clinton. 2019. Reading from paper compared to screens: A systematic review and meta-analysis. *Journal of research in reading*
 794 42, 2 (2019), 288–325.
- 795 [14] Max Coltheart. 1981. The MRC psycholinguistic database. *The Quarterly Journal of Experimental Psychology Section A* 33, 4 (1981),
 796 497–505.
- 797 [15] Anne Cunningham and Keith Stanovich. 2003. Reading Can Make You Smarter!. *Principal* 83, 2 (2003), 34–39.
- 798 [16] Pablo Delgado and Ladislao Salmerón. 2022. Cognitive Effort in Text Processing and Reading Comprehension in Print and on Tablet: An
 799 Eye-Tracking Study. *Discourse Processes* 59, 4 (2022), 237–274.
- 800 [17] Ekaterina Denkova, Jason S Nomi, Lucina Q Uddin, and Amishi P Jha. 2019. Dynamic brain network configurations during rest and an
 801 attention task with frequent occurrence of mind wandering. *Human brain mapping* 40, 15 (2019), 4564–4576.
- 802 [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for
 803 Language Understanding. *arXiv preprint arXiv:1810.04805* (2018).
- 804 [19] Sidney K D'Mello, Caitlin Mills, Robert Bixler, and Nigel Bosch. 2017. Zone out No More: Mitigating Mind Wandering during
 805 Computerized Reading. *International Educational Data Mining Society* (2017).
- 806 [20] Myrthe Faber, Robert Bixler, and Sidney K D'Mello. 2018. An automated behavioral measure of mind wandering during computerized
 807 reading. *Behavior Research Methods* 50, 1 (2018), 134–150.
- 808 [21] Aisha Farid, Muhammad Ishtiaq, and Muhammad Saboor Hussain. 2020. A Review of Effective Reading Strategies to Teach Text
 809 Comprehension to Adult English Language Learners. *Global Language Review* (2020), 77–88.
- 810 [22] Michael S Franklin, Jonathan Smallwood, and Jonathan W Schooler. 2011. Catching the mind in flight: Using behavioral indices to
 811 detect mindless reading in real time. *Psychonomic bulletin & review* 18, 5 (2011), 992–997.
- 812 [23] Edward Gibson and Evelina Fedorenko. 2013. The need for quantitative methods in syntax and semantics research. *Language and Cognitive
 813 Processes* 28, 1-2 (2013), 88–124. <https://doi.org/10.1080/01690965.2010.515080> arXiv:<https://doi.org/10.1080/01690965.2010.515080>
- 814 [24] Amber Gove and Peter Cvelich. 2011. Early reading: Igniting education for all. *A report by the Early Grade Learning Community of
 815 Practice. Revised Edition*. Research Triangle Park, NC: Research Triangle Institute (2011).
- 816 [25] Malcolm Haynes and Thad Starner. 2018. Effects of lateral eye displacement on comfort while reading from a video display terminal.
 817 *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 1, 4 (2018), 1–17.
- 818 [26] John M Henderson and Fernanda Ferreira. 1993. Eye movement control during reading: fixation measures reflect foveal but not parafoveal
 819 processing difficulty. *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie expérimentale* 47, 2 (1993), 201.
- 820 [27] Riku Higashimura, Andrew Vargo, Motoi Iwata, and Koichi Kise. 2022. Helping Mobile Learners Know Unknown Words through their
 821 Reading Behavior. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*. 1–5.
- 822 [28] Jin Huang, Chun Yu, Yuntao Wang, Yuhang Zhao, Siqi Liu, Chou Mo, Jie Liu, Lie Zhang, and Yuanchun Shi. 2014. FOCUS: enhancing
 823 children's engagement in reading by using contextual BCI training sessions. In *Proceedings of the SIGCHI Conference on Human Factors
 824 in Computing Systems*. 1905–1908.
- 825 [29] Michael Xuelin Huang, Tiffany CK Kwok, Grace Ngai, Hong Va Leong, and Stephen CF Chan. 2014. Building a self-learning eye gaze
 826 model from user interaction data. In *Proceedings of the 22nd ACM international conference on Multimedia*. 1017–1020.
- 827 [30] Stephen Hutt, Kristina Krasich, Caitlin Mills, Nigel Bosch, Shelby White, James R Brockmole, and Sidney K D'Mello. 2019. Automated
 828 gaze-based mind wandering detection during computerized learning in classrooms. *User Modeling and User-Adapted Interaction* 29
 829 (2019), 821–867.
- 830 [31] Jukka Hyönä and Richard K Olson. 1995. Eye fixation patterns among dyslexic and normal readers: effects of word length and word
 831 frequency. *Journal of Experimental Psychology: Learning, Memory, and Cognition* 21, 6 (1995), 1430.
- 832 [32] Aulikki Hyrskykari. 2006. Utilizing eye movements: Overcoming inaccuracy while tracking the focus of attention during reading.
 833 *Computers in human behavior* 22, 4 (2006), 657–671.
- 834 [33] Aulikki Hyrskykari, Päivi Majaranta, Antti Aaltonen, and Kari-Jouko Räihä. 2000. Design issues of iDict: a gaze-assisted translation aid.
 835 In *Proceedings of the 2000 symposium on Eye tracking research & applications*. 9–14.
- 836 [34] Md Rabiul Islam, Shuji Sakamoto, Yoshihiro Yamada, Andrew W Vargo, Motoi Iwata, Masakazu Iwamura, and Koichi Kise. 2021.
 837 Self-supervised learning for reading activity classification. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous
 838 Technologies* 5, 3 (2021), 1–22.

- 839 [35] Philip C. Jackson. 2018. Natural language in the Common Model of Cognition. *Procedia Computer Science* 145 (2018), 699–709.
 840 <https://doi.org/10.1016/j.procs.2018.11.051> Postproceedings of the 9th Annual International Conference on Biologically Inspired
 841 Cognitive Architectures, BICA 2018 (Ninth Annual Meeting of the BICA Society), held August 22-24, 2018 in Prague, Czech Republic.
- 842 [36] Ariel N James, Scott H Fraundorf, Eun-Kyung Lee, and Duane G Watson. 2018. Individual differences in syntactic processing: Is there
 843 evidence for reader-text interactions? *Journal of memory and language* 102 (2018), 155–181.
- 844 [37] Yu-Cin Jian. 2021. The immediate and delayed effects of text-diagram reading instruction on reading comprehension and learning
 845 processes: evidence from eye movements. *Reading and Writing* 34, 3 (2021), 727–752.
- 846 [38] Yu-Cin Jian. 2022. Influence of science text reading difficulty and hands-on manipulation on science learning: An eye-tracking study.
 847 *Journal of Research in Science Teaching* (2022), 358–382.
- 848 [39] Marcel A Just and Patricia A Carpenter. 1980. A theory of reading: from eye fixations to comprehension. *Psychological review* 87, 4
 849 (1980), 329.
- 850 [40] Yuki Kamide and Anuenue Kukona. 2018. The influence of globally ungrammatical local syntactic constraints on real-time sentence
 851 comprehension: Evidence from the visual world paradigm and reading. *Cognitive Science* (2018), 2976–2998.
- 852 [41] Moritz Kassner, William Patera, and Andreas Bulling. 2014. Pupil: an open source platform for pervasive eye tracking and mobile
 853 gaze-based interaction. In *Proceedings of the 2014 ACM international joint conference on pervasive and ubiquitous computing: Adjunct
 854 publication*. 1151–1160.
- 855 [42] Diksha Khurana, Aditya Koli, Kiran Khatter, and Sukhdev Singh. 2022. Natural language processing: State of the art, current trends and
 856 challenges. *Multimedia tools and applications* (2022), 1–32.
- 857 [43] Jumpei Kobayashi and Toshio Kawashima. 2019. Paragraph-based faded text facilitates reading comprehension. In *Proceedings of the
 858 2019 CHI Conference on Human Factors in Computing Systems*. 1–12.
- 859 [44] Arnout W Koornneef and Jos JA Van Berkum. 2006. On the use of verb-based implicit causality in sentence comprehension: Evidence
 860 from self-paced reading and eye tracking. *Journal of Memory and Language* (2006), 445–465.
- 861 [45] Thomas Kosch, Albrecht Schmidt, Simon Thanheiser, and Lewis L Chuang. 2020. One does not simply RSVP: mental workload to select
 862 speed reading parameters using electroencephalography. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing
 863 Systems*. 1–13.
- 864 [46] Kyle Krafcik, Aditya Khosla, Petr Kellnhofer, Harini Kannan, Suchendra Bhandarkar, Wojciech Matusik, and Antonio Torralba. 2016. Eye
 865 tracking for everyone. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2176–2184.
- 866 [47] Richard L Lewis, Shravan Vasishth, and Julie A Van Dyke. 2006. Computational principles of working memory in sentence comprehension.
 867 *Trends in cognitive sciences* 10, 10 (2006), 447–454.
- 868 [48] Tal Linzen. 2018. What can linguistics and deep learning contribute to each other? *arXiv preprint arXiv:1809.04179* (2018).
- 869 [49] Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. 2016. Assessing the ability of LSTMs to learn syntax-sensitive dependencies.
 870 *Transactions of the Association for Computational Linguistics* 4 (2016), 521–535.
- 871 [50] Steven G Luke and Kiel Christianson. 2018. The Provo Corpus: A large eye-tracking corpus with predictability norms. *Behavior research
 872 methods* (2018), 826–833.
- 873 [51] Robert A Mason and Marcel Adam Just. 2007. Lexical ambiguity in sentence comprehension. *Brain research* (2007), 115–127.
- 874 [52] Joseph E Michaelis and Bilge Mutlu. 2017. Someone to read with: Design of and experiences with an in-home learning companion robot
 875 for reading. In *Proceedings of the 2017 CHI conference on human factors in computing systems*. 301–312.
- 876 [53] Brian W Miller. 2015. Using reading times and eye-movements to measure cognitive engagement. *Educational psychologist* 50, 1 (2015),
 877 31–42.
- 878 [54] Caitlin Mills, Julie Gregg, Robert Bixler, and Sidney K D'Mello. 2021. Eye-mind reader: An intelligent reading interface that promotes
 879 long-term comprehension by detecting and responding to mind wandering. *Human–Computer Interaction* 36, 4 (2021), 306–332.
- 880 [55] Simon Moe and Michael Wright. 2013. Can accessible digital formats improve reading skills, habits and educational level for dyslexic
 881 youngsters? In *International Conference on Universal Access in Human-Computer Interaction*. 203–212.
- 882 [56] Robert E Morrison. 1984. Manipulation of stimulus onset delay in reading: evidence for parallel programming of saccades. *Journal of
 883 Experimental psychology: Human Perception and performance* 10, 5 (1984), 667.
- 884 [57] AB MySQL. 2001. MySQL.
- 885 [58] Rustam Nazuryt, Nurullaningsih Priyanto, Sarmandan Anggia Pratiwi, and Amirul Mukminin. 2019. Learning strategies in reading: The
 886 case of Indonesian language education student teachers. *Universal Journal of Educational Research* (2019), 2536–2543.
- 887 [59] Pablo Oyarzo, David D Preiss, and Diego Cosmelli. 2022. Attentional and meta-cognitive processes underlying mind wandering episodes
 888 during continuous naturalistic reading are associated with specific changes in eye behavior. *Psychophysiology* 59, 4 (2022), e13994.
- 889 [60] Alexandra Papoutsaki, Patsorn Sangkloy, James Laskey, Nediyana Daskalova, Jeff Huang, and James Hays. 2016. Webgazer: Scalable
 890 Webcam Eye Tracking Using User Interactions. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence
 891 (New York, New York, USA) (IJCAI'16)*. AAAI Press, 3839–3845.
- 892 [61] Ellie Pavlick. 2022. Semantic structure in deep learning. *Annual Review of Linguistics* (2022), 447–471.

- 893 [62] Charles A Perfetti et al. 1999. Comprehending written language: A blueprint of the reader. *The neurocognition of language* 167 (1999),
 894 208.
- 895 [63] Hilary Putnam. 1988. *Representation and reality*. MIT press.
- 896 [64] Keith Rayner. 1998. Eye movements in reading and information processing: 20 years of research. *Psychological bulletin* 124, 3 (1998), 372.
- 897 [65] Erik D Reichle. 2006. Computational models of eye-movement control during reading: Theories of the "eye-mind" link. (2006).
- 898 [66] Erik D Reichle, Keith Rayner, and Alexander Pollatsek. 2003. The EZ Reader model of eye-movement control in reading: Comparisons to
 899 other models. *Behavioral and brain sciences* 26, 4 (2003), 445–476.
- 900 [67] Erik D Reichle, Andrew E Reineberg, and Jonathan W Schooler. 2010. Eye movements during mindless reading. *Psychological science*
 901 (2010), 1300–1310.
- 902 [68] Margarete Sandelowski. 1995. Qualitative analysis: What it is and how to begin. *Research in nursing & health* 18, 4 (1995), 371–375.
- 903 [69] Martin Schrimpf, Idan Asher Blank, Greta Tuckute, Carina Kauf, Eghbal A Hosseini, Nancy Kanwisher, Joshua B Tenenbaum, and
 904 Evelina Fedorenko. 2021. The neural architecture of language: Integrative modeling converges on predictive processing. *Proceedings of
 905 the National Academy of Sciences* 118, 45 (2021), e2105646118.
- 906 [70] Mona L Scott and MONA L SCOTT. 1998. Dewey decimal classification. *Libraries Unlimited* (1998).
- 907 [71] Prafull Sharma and Yingbo Li. 2019. Self-supervised contextual keyword and keyphrase retrieval with self-labelling. (2019).
- 908 [72] John L Sibert, Mehmet Gokturk, and Robert A Levine. 2000. The reading assistant: eye gaze triggered auditory prompting for reading
 909 remediation. In *Proceedings of the 13th annual ACM symposium on User interface software and technology*. 101–107.
- 910 [73] Subroto Singha. 2021. *Gaze Based Mind Wandering Detection Using Deep Learning*. Ph. D. Dissertation. Texas A&M University-Commerce.
- 911 [74] Matthew S Starr and Keith Rayner. 2001. Eye movements during reading: Some current controversies. *Trends in cognitive sciences* 5, 4
 912 (2001), 156–163.
- 913 [75] Mikhail Startsev, Ioannis Agtzidis, and Michael Dorr. 2019. 1D CNN with BLSTM for automated classification of fixations, saccades, and
 914 smooth pursuits. *Behavior Research Methods* 51, 2 (2019), 556–572.
- 915 [76] Sawitri Suwanaroa et al. 2021. Factors and Problems Affecting Reading Comprehension of Undergraduate Students. *International
 916 Journal of Linguistics, Literature and Translation* 4, 12 (2021), 15–29.
- 917 [77] Amos Van Gelderen, Rob Schoonen, Reinoud D Stoel, Kees De Glopper, and Jan Hulstijn. 2007. Development of adolescent reading
 918 comprehension in language 1 and language 2: A longitudinal analysis of constituent components. *Journal of Educational Psychology*
 919 (2007), 477.
- 920 [78] Vladimir Vapnik. 1999. *The nature of statistical learning theory*. Springer science & business media.
- 921 [79] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017.
 922 Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- 923 [80] Shaun Wallace, Zoya Bylinskii, Jonathan Dobres, Bernard Kerr, Sam Berlow, Rick Treitman, Nirmal Kumawat, Kathleen Arpin, Dave B
 924 Miller, Jeff Huang, et al. 2022. Towards Individuated Reading Experiences: Different Fonts Increase Reading Speed for Different
 925 Individuals. *ACM Transactions on Computer-Human Interaction (TOCHI)* (2022), 1–56.
- 926 [81] Shang Wang and Erin Walker. 2021. Providing Adaptive Feedback in Concept Mapping to Improve Reading Comprehension. In
 927 *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–11.
- 928 [82] Edward W Wlotko and Kara D Federmeier. 2015. Time for prediction? The effect of presentation rate on predictive sentence comprehen-
 929 sion during word-by-word reading. *Cortex* (2015), 20–32.
- 930 [83] Fang-Ying Yang. 2017. Examining the reasoning of conflicting science information from the information processing perspective—an eye
 931 movement analysis. *Journal of Research in Science Teaching* 54, 10 (2017), 1347–1372.
- 932 [84] Li Yang and Tam Shu Sim. 2017. Metacognitive Awareness of Reading Strategies among EFL High School Students in China. *AJELP:
 933 Asian Journal of English Language and Pedagogy* (2017), 34–45.
- 934 [85] Shun-nan Yang. 2006. An oculomotor-based model of eye movements in reading: The competition/interaction model. *Cognitive Systems
 935 Research* 7, 1 (2006), 56–69.
- 936 [86] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive
 937 pretraining for language understanding. *Advances in neural information processing systems* 32 (2019).
- 938 [87] Zhuosheng Zhang, Yuwei Wu, Hai Zhao, Zuchao Li, Shuailiang Zhang, Xi Zhou, and Xiang Zhou. 2020. Semantics-aware BERT for
 939 language understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 9628–9635.
- 940 [88] Yingying Zhao, Yuhu Chang, Yutian Lu, Yujiang Wang, Mingzhi Dong, Qin Lv, Robert P Dick, Fan Yang, Tun Lu, Ning Gu, and Li Shang.
 941 2022. Do Smart Glasses Dream of Sentimental Visions? Deep Emotionship Analysis for Eyewear Devices. *Proceedings of the ACM on
 942 Interactive, Mobile, Wearable and Ubiquitous Technologies* 6, 1 (2022), 1–29.