

The Eyes Have It : Characteristics of Deep Reading Activities on Desktop and Mobile Devices

ANONYMOUS AUTHOR(S)

Deep reading fosters text comprehension, memory, and critical thinking. But with mobile interfaces and online incentives around attention capture, concerns have grown about deep reading activities being replaced by skimming and sifting through information. Traditionally, reading quality is assessed using comprehension tests, which require readers to explicitly answer a set of carefully composed questions. To quantify and understand reading behaviour in natural settings and at scale, however, implicit measures are needed of deep versus skim reading across desktop and mobile devices. In this paper, we present an approach to systematically induce and detect deep and skim reading patterns using eye movement and interaction data. Based on a user study with 36 participants, we created models that detect deep reading on both devices with up to 0.82 AUC. We present the characteristics of deep reading and discuss how our models can be used to monitor long-term changes in reading behaviours.

CCS Concepts: • **Human-centered computing** → **Ubiquitous and mobile computing systems and tools; Human computer interaction (HCI).**

Additional Key Words and Phrases: Reading mode classification, Digital Devices, Eye tracking, Gaze features

ACM Reference Format:

Anonymous Author(s). 2022. The Eyes Have It : Characteristics of Deep Reading Activities on Desktop and Mobile Devices. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 0, 0, Article 0 (March 2022), 21 pages. <https://doi.org/10.1145/xxxxxx>

1 INTRODUCTION

Over the past five thousand years humans have developed the ability to read and to read deeply, spending years developing an array of sophisticated cognitive processes that include inferential and deductive reasoning, analogical skills, critical analysis, reflection, and insight [49]. In contrast with *skim reading* where readers only process factual and explicit information stated in texts to regurgitate, *deep reading* requires readers to interpret implicit information, make inferences, analyze, and reason, and reflect on it based on prior beliefs [39, 50]. These deep reading skills have followed the development of widespread literacy and the ability to produce written documents in astounding quantities, especially since the Industrial Revolution. Today, more than 86% of the world's population is considered literate¹ and this literacy has changed how we can rapidly communicate complex ideas across communities, countries and nations. While reading has never been a static discipline and reading habits have evolved over time (for example from reading aloud to silent reading), with the advent of the digital revolution, where, how, and what we read is changing rapidly. Chronic information overload, the ability to search and "surf" from one site to another as well as attention-seeking interface design all work against the practice of deeply reading a single, potentially complex piece of writing and drawing inferences and insights from it. While the ability to distill insights from a document takes mere milliseconds for a trained mind, the skill of deep reading takes years to develop. There is currently increasing concern that as a new generation grows

¹<https://ourworldindata.org/literacy>

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.

2474-9567/2022/3-ART0 \$15.00

<https://doi.org/10.1145/xxxxxx>

48 up reading more and more on mobile devices, deep reading skills may become underdeveloped [50]. While this
 49 concern seems well supported, there currently does not exist a reliable implicit method for measuring deep versus
 50 skim reading habits across both desktop and mobile devices. Without being able to measure the potential decline
 51 of deep reading, we will remain uninformed about the true magnitude of the shift and the factors that might be
 52 accelerating it and we will not be able to create an informed strategy to preserve deep reading skills.

53 To address the problem of implicitly measuring deep reading across both desktop and mobile devices we
 54 first developed and tested a method for inducing both deep reading and skim reading on both these platforms,
 55 then performed a series of user tests where we measured eye movement patterns in both conditions and finally
 56 we developed an algorithm that could predict deep or skim reading on either platform using these patterns.
 57 Traditionally, deep versus skim reading is assessed explicitly using reading comprehension tests. These tests
 58 consist of a carefully constructed set of questions that evaluate the reader's literal, inferential, and evaluative
 59 understanding of text [6]. Unfortunately, this evaluation method does not scale as it is infeasible both to create
 60 such questions for every book or article or to expect readers to fill in a set of questions each time they consume a
 61 piece of writing. We instead need an implicit method of measuring reading modes.

62 Since cognition can not be measured directly, technologies need to use surrogate measures to sense and infer
 63 cognitive activities. The eyes have often been described as providing a "window into our mind" [47] and so eye
 64 movements have been used to detect reading [23] as well as differentiate between different reading goals [45].
 65 And while eye fixation durations and saccade lengths have been shown to correlate with readers' comprehension
 66 levels [45, 46], robust models capable of classifying deep and skim reading are yet to be made work on both desktop
 67 and mobile devices, which is where reading increasingly takes place [40]. With the prevalence of front-facing
 68 device cameras, continuous eye tracking becomes increasingly feasible.

69 We evaluated our approach in a lab-based experiment with 36 participants during which we recorded gaze and
 70 interaction data while reading on desktop and a mobile device. We extracted low and mid-level eye-movement
 71 features to train a machine-learning model to differentiate between deep and skim reading. Our results show a
 72 strong correlation between eye-movement patterns and deep reading patterns and the resulting models allow
 73 us to detect deep reading with 0.82 AUC on desktop and 0.73 AUC on a mobile device. We present a detailed
 74 analysis of the differences in the reading patterns on desktop versus mobile, which paves the way for continuous
 75 and unobtrusive tracking and quantification of reading activities.

76 2 RELATED WORK

77 2.1 Deep and Skim Reading Behavior

78 Reading comprehension can be loosely defined as the amount of information or meaning being extracted from
 79 the texts [20, 43]. Moreover, comprehension can be further categorized into different levels based on the amount
 80 of cognitive demands and interactions [31, 42]. Commonly, the underlying theory categorized comprehension
 81 into three levels – literal, inferential, and evaluative [5, 6, 11, 20].

- 82 • **Literal comprehension** requires readers to retrieve explicit information from the texts or recall what is
 stated in the text.
- 83 • **Inferential comprehension** requires readers to interpret the authors' meaning by connecting information
 that is implicit in the text.
- 84 • **Evaluative comprehension** requires readers to analyse information based on previous knowledge or
 experiences and thus relate what is being read to what is known.

85 Based on the level of comprehension, readers' reading behavior can be categorized into *deep* and *skim* reading,
 86 where superficial comprehension level (literal) is achieved in skim reading and deeper comprehension level
 87 (inferential and evaluative) is obtained through deep reading [39, 50]. Identifying deep and skim reading is crucial
 88 in many research and applications in psychology, education, and user interface design [19, 39, 50]. However,
 89

95 differentiating such reading behaviors would require thorough comprehension tests (e.g., via comprehension
 96 questions), hence making it infeasible to be implemented and adopted easily. Moreover, individuals' reading
 97 behaviors are very dynamic [2, 21], which makes the analysis more challenging.
 98

99 2.2 Implicit Assessment of Reading Comprehension using Eye Movements

100 The relationship between eye movements and human cognition has been extensively studied in the past. Many
 101 cognitive processes are infeasible to observe or measure directly but are highly correlated with eye movements
 102 [28, 35]. As a result, much research has used eye movements data to infer or analyse people's cognitive behavior,
 103 such as people's attention [1, 9], reading behavior [8, 26], and reading comprehension [3, 14–16, 30, 46].
 104

105 Eye-tracking research has shown great potential in classifying reading behaviors and predicting readers'
 106 comprehension from their gaze data. For instance, eye movements are correlated with comprehension [23, 24, 34,
 107 41], and more specifically, fixation duration and saccade length are known to be correlated with reading behavior
 108 [45] and comprehension level [45, 46]. The line of work that uses eye movements to classify comprehension was
 109 started by Underwood et al. [46]. They studied the relationship between eye fixations and reading comprehension,
 110 which identified that fixation duration could be an indicator of comprehension. However, the generalizability of
 111 their results was limited as they trained and tested on the same dataset. Copeland et al. [14–16], instead of focusing
 112 on comprehensions, studied predicting the accuracy of the answers. Makowski et al. [30] continued tackling the
 113 comprehension prediction problem. While succeeding in reader identities identification, their approach failed
 114 at predicting comprehension levels. Moreover, their texts were in German with approximately 158 words each,
 115 and there were only three comprehension questions. Hence, their study design may not be the most suitable to
 116 systematically induce deep and skim readings. Recently, Ahn et al. [3] studied predicting people's comprehension
 117 level using gaze patterns as well as other reading parameters such as reading difficulty. The comprehension
 118 was labeled in binary as high and low, based on the percentage of correctly answered questions. For overall
 119 comprehensions, they obtained an accuracy of 64% (11% better than the null accuracy) on testing against unseen
 120 passages (for the same participant), but failed to beat the null model when testing against unseen participants.
 121 Although account for comprehension levels, their study did not consider participants' reading modes and verify
 122 them via question analysis, also their results were limited in terms of generalizability.
 123

124 2.3 Implicit Assessment of Reading Behavior using Eye Movements

125 Our goal is to classify deep reading, where material is being thoughtfully processed, versus skim reading where
 126 the reader either not comprehending the text or skipping through it or scanning for a singular fact. In truth,
 127 modes of reading are not often entirely separate, sometimes skim reading will occur during deep reading and
 128 conversely occasionally a skimming reader will take interest in an article and begin reading more deeply.

129 Very little research has been done in utilizing the correlated eye-movement features to help classify deep
 130 and skim readings. Biedert et al. [8] showed that eye tracking could be effective for identifying skim reading
 131 of paper news articles. In this study, two groups of readers were asked to produce a small set of keywords to
 132 describe the news articles. One group was given a very short time to read the article (skim reading induction)
 133 whereas the other cohort was given unlimited time to read the article (no skim pressure - deep reading assumed).
 134 Using eye saccade metrics these researchers were able to differentiate between the two conditions with 86%
 135 accuracy. In another study, Kelton et al. [26] refined this study and looked at identifying skim versus deep
 136 reading on both a global (entire text) and local (specific region) scale. They followed a similar study design to
 137 Biedert et al. [8] and achieved an accuracy of 82.5% differentiating global skim vs. deep conditions and 72% – 95%
 138 accuracy for differentiating these modes in different local areas of the document. While these studies showed
 139 that time-pressured readers did not stop (measured via saccade) to read the majority of the document, this only
 140 proved a differentiation between skim versus assumed deep reading based on eye movement patterns without
 141

¹⁴² taking into consideration readers' comprehension levels. Moreover, to the best of our knowledge, there is a gap
¹⁴³ in thoroughly evaluating deep reading patterns on mobile devices.

¹⁴⁴

¹⁴⁵ 3 DATA COLLECTION

¹⁴⁶

For our data collection, we factored in three main limitations found in previous studies and introduced a new study design to systematically induce deep reading and skim reading behaviour on both desktop and mobile devices. Our study design differed from previous studies as follows: Firstly, we used reading materials that were longer (≈ 1500 words) compared to previous studies[8, 26] where mostly news articles (≈ 150 words) were used, since it is easier to induce and evaluate deep reading in them. Further, we gave more time to the reading due to our text length. Secondly, we used a set of carefully designed comprehension questions to systematically evaluate different types of comprehension and reading mode. Lastly, to ensure skim reading behavior, we replaced "keywords summarizing" techniques commonly used in previous studies [8, 26] with "finding answers to preview-questions". This was done to counter the limitations of previous works [8, 26] where researchers induced skim reading behaviour by asking participants to select three keywords that best describes the article. However, we believe that keywords of the articles can be easily obtained via inferring from the title and a few paragraphs, which does not explicitly validate skim reading behaviour.

¹⁵⁸

¹⁵⁹

¹⁶⁰ 3.1 Tasks

¹⁶¹

We prepared four reading tasks in total – one deep reading task on mobile, one deep reading task on desktop, one skim reading task on mobile, and one skim reading task on desktop. The deep reading and skim reading tasks were designed as follows:

¹⁶⁴

Deep Reading: The deep reading task's purpose was to make participants thoroughly process and comprehend the reading material. Participants were, therefore, asked to take as much time as they needed to read an article. They were explicitly told that they would have to answer 20 in-depth multiple-choice questions about the text's content. Additionally, we alerted participants to the fact that—once the questions were shown—they could not revisit the article. Participants were encouraged to try and get as many questions correct as they could.

¹⁶⁹

Skim Reading: The skim reading task was designed to induce participants' skim reading behavior, i.e., obtaining literal information in the shortest amount of time. In this task, the participants were provided with three literal questions about the text beforehand and the goal was to find the answers to these three questions using as little time as they could. Participants should be able to answer the three literal questions by retrieving explicit information equally distributed across the text. In the preview stage, only the question text was displayed but not the text. Once they started reading, participants were unable to jump back to the questions. They were meant to act as primers only. Additionally, we imposed a two-minute limit on the reading time. After reading, participants were asked to answer the three previously provided questions plus another 17 multiple-choice to test what else they took away from the test. Participants were informed that there were 20 questions in total but only the three previewed ones matter, and they were unable to refer back the article while answering the questions.

¹⁷⁹

¹⁸⁰

¹⁸¹ 3.2 Reading Materials

¹⁸²

We selected four articles from the easyCBM repository [4], which we randomly assigned across our two-by-two study design: deep versus skim and desktop versus mobile reading. To reduce the impact of prior knowledge and reading difficulty, all articles were stories with approximately 1500 words in length and of 8th-grade reading level. Each article comes with 20 multiple-choice comprehension questions: seven questions tested literal, seven questions inferential, and six questions tested evaluative comprehension. Table 1 and Table 2 lists details for the articles used.

¹⁸⁷

¹⁸⁸

189 Table 1. Readability Statistics of Selected Article. FK stands for “Flesch Kincaid”, GF stands for “Gunning Fog”, SMOG stands
 190 for “Simple Measure of Gobbledygook”, CL stands for “Coleman Liau”, and AR stands for “Automated Readability”.

191

192 Article	FK Reading Ease	FK Grade Level	GF Score	SMOG Index	CL Index	AR Index
194 1	76.4	6.3	8.6	6.7	9.5	6.2
195 2	76.8	7.5	10.2	6.3	9	8.5
196 3	68.5	7.4	10	7.5	11	7.4
197 4	77.8	5.6	7.9	5.9	9.5	5.2
198						

199

200 Table 2. Text Statistics of Selected Article. Sent. stands for “Sentences”

201

202 Article	Sent.	Words	Complex Words	Complex Words %	Words per Sent.	Syllables per Word
204 1	131	1943	169	8.70%	14.83	1.36
205 2	79	1584	88	5.56%	20.05	1.30
206 3	99	1461	160	10.95%	14.76	1.46
207 4	138	1749	132	7.55%	12.67	1.37
208						

209



(a) Desktop Display Example



(b) Mobile Display Example

220

221 Fig. 1. Examples of Desktop and Mobile Displays. Text displays were controlled to be identical across devices.

222

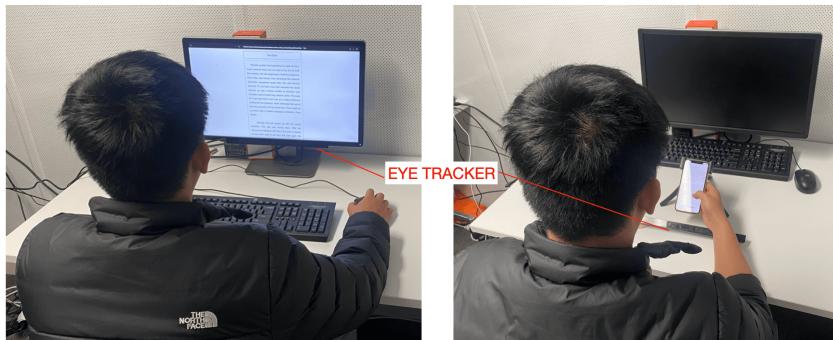
223 The reading-interface parameters were controlled to ensure consistent reading experiences between mobile
 224 and desktop. In particular, we showed the same amount of characters on each line (55 characters per line) for
 225 both devices. To achieve this, we used the same font family (an even spaced sans serif font, Arial), line space
 226 (double spacing), and justification (left justified). Since mobile and desktop have different screen size, we selected
 227 different font size to achieve our primary objective on keeping the line character count identical across devices.
 228 This was compensate by the closer distance between reader and mobile than between reader and desktop. The
 229 line width was set to be 600 px on desktop and full screen on mobile. Our choice of parameters was in alignment
 230 with previous readability studies [17, 32, 33, 37, 38] so that it enhance both readabilities (see Figure 1) and
 231

232

236 eye-tracking accuracy. In both devices, texts navigation worked by vertically scrolling as this interaction is the
 237 de-facto standard in electronic interfaces now.

238 239 3.3 Experimental Setup

240 We recorded participants' eye movement data using the Tobii Pro X3-130 eye tracker² and the Tobii Pro Studio
 241 software. The eye tracker had a sampling rate of 120 Hz. The output data contains raw gaze coordinates (x, y),
 242 left and right pupil diameter in millimeter, as well as the fixation and saccade information generated by the
 243 Pro Studio. The eye-tracker was mounted at the bottom of a 24-inch monitor for desktop reading, and at the
 244 bottom of an iPhone 11 device for mobile reading (see Figure 2). Both devices kept a log of participants' scrolling
 245 behavior, i.e., the x and y offset of the current viewing page.



250 251 252 253 254 255 256 257 Fig. 2. Experiment Setup. Left: desktop reading; Right: mobile reading

261 3.4 Participants

262 We recruited 36 participants, all of them were students and staffs from our university. The data from seven
 263 participants was discarded due to non-native English speaker or insufficient reading, i.e., not being able to go
 264 through at least half of the article because they hit the two-minute time limit during the skim reading task.
 265 The remaining 29 participants (16 women, 13 men) were all native English speakers. Most participants had an
 266 education level of postgraduate degree (N=14), followed by bachelor degree (N=13), while the remaining were
 267 graduate certificate or diploma (N=1) and certificate III or IV (N=1).

268 3.5 Procedure

269 Upon arrival, we informed participants about the purpose and the procedure of this study and asked them to
 270 sign a consent form. Next, we seated participants in a comfortable position and asked them to adjust the seat
 271 such that their heads were centered and approximately 60 – 65 cm away from the desktop screen (and 50 cm
 272 from the mobile phone screen). Before each reading task, a standard, built-in, 9-point eye-tracker calibration
 273 procedure was done. To start the study, participants were asked to fill in a pre-study survey regarding their
 274 demographics and past reading experiences. Then, they finished the desktop reading task and the mobile reading
 275 task in counterbalanced order. On each device, participants were asked to read two articles in counterbalanced
 276 order, one in deep reading mode and the other in skim reading mode. The four articles were randomly assigned
 277 to each reading session. During each reading session, participants were instructed to use the mouse (desktop) or
 278 swipe on the screen (mobile) to scroll down. Further, while reading, participants were free to move their heads

280
 281 ²<https://www.tobiipro.com/product-listing/tobii-pro-x3-120/>

283 slightly as suggested by the robustness of the eye tracker. Lastly, participants were asked to finish a post-study
 284 survey about their experiences in this study. The study took around 45 – 60 minutes to finish and upon completion,
 285 participants were compensated with a 10 dollar coupon.

286

287 4 READING MODE DETECTION

288 In this section, we describe our step-by-step approach for detecting deep versus skim reading. To begin with,
 289 we pre-processed the gaze-data by adjusting the raw gaze data with the scrolling behavior and converting it to
 290 fixations and saccades. Next, we extracted 50 low-level and 25 mid-level [44] gaze features from each sliding
 291 window of the processed data. We then built and test various machine learning classifiers and analyzed its
 292 performance in more detail.
 293

294

295 4.1 Data Preprocessing

296 Although the Tobii Pro Studio software can pre-process the raw gaze-date to extract fixations, it does not consider
 297 the effect of scrolling events. For example, participants may stare at one point while scrolling. Therefore, we first
 298 incorporated the scrolling event into the raw gaze data, such that each gaze point is adjusted by the offset of the
 299 current page. Then, we generated fixations from the adjusted gaze data using the R software package *Saccades*
 300 [48] that uses the velocity-based algorithm for saccade detection proposed by Engbert and Kliegl [18]. Next, we
 301 identified saccades as transitions between consecutive fixations.

302 For most participants, deep reading took longer to complete than skim reading. Hence, classification of the
 303 whole session would be trivial as the duration and fixation number themselves were already strong indicators.
 304 In order to evaluate the robustness of our classifiers, we used the sliding-window technique to partition our
 305 time-series gaze data into consecutive windows. The window size was fixed and all data points within this
 306 window were aggregated to generate features. We defined the start time difference between two consecutive
 307 windows as the step size. Note that if the step size equals window size, then there is no overlap between windows.
 308 Further, window size and step size may have a huge impact on classifications. For example, a smaller window size
 309 could give more fine-grained results but may not able to capture sufficient movement patterns. Previous studies
 310 [10, 27, 29] suggested different window sizes from 20 seconds to 60 seconds. Hence, in this study, we varied the
 311 window size and step size to measure their impact on performance.

312

313 4.2 Feature Extraction

314 To build the classifier based on preprocessed fixation and saccade data, we required a feature set that best captures
 315 the reading behavior. Typically, gaze features can be categorized into three levels – low-level, mid-level, and
 316 high-level [44]. Low-level gaze features are features that can be derived directly from the raw data, such as
 317 fixation duration, saccade length, and pupil diameter. Many similar studies have used mainly low-level features as
 318 their feature set [22, 25]. Mid-level gaze features were proposed by Srivastava et al. [44] to better capture various
 319 eye movement patterns. Mid-level gaze features consider several consecutive fixations and saccades together,
 320 and categorize their shape into patterns such as *compare*, *scan*, and *line reading*. High-level gaze features, on the
 321 other hand, are stimuli-specific features whose analysis normally requires the identification of an area of interest
 322 (AOI). For example, texts that cover a question can be treated as one AOI and the time spent in this AOI can be
 323 considered as one feature. In this study, we included 50 low-level features and 25 mid-level features as shown in
 324 Table 3 and Table 4. We did not consider high-level features as they were hard to extract and to be generalized, as
 325 AOI could vary from text to text.

326 For all fixations, saccades, and pupil diameters over a span of time (window), we calculated metrics, such
 327 as count, mean, standard deviation, variance, minimum value, and maximum value. Also, we derived the rate,
 328 percentage, and slope of all fixations. Fixation rate was defined as the ratio between the total fixation duration
 329

Table 3. Selected Low-level Features

Category	Features	Measurement
Fixation	number, mean, std, var, min, max	value
Duration	rate, percentage, slope	value
	short, medium, long	count
	dispersion area (75%)	value
Saccade	mean, std, var	value
Length	follow, neighbor, opposite, opposite-neighbor	count
	right, left, up, down, up-right, down-right, up-left, down-left	short, medium, long count
Pupil Diameter	mean, std, var	left and right value

Table 4. Selected Mid-level Features. All features are measured by their count.

Category	Sub-category	Features
Shape Based	String	up, right, down, left
	Line	small, medium, long
	Comparison	left-right, right-left, up-down, down-up
	Scan	right-left, left-right, up-down, down-up, right-up, up-right, left-up, up-left, right-down, down-right, left-down, down-left
Distance Based	-	regression, else-where

and the data duration (i.e., window size). Fixation percentage was the ratio between the count of fixations and the data duration. And fixation slope was the slope of the regression line fitted on all fixations. For fixation duration size, we set the thresholds as 200 ms and 400 ms, i.e., fixations within 200 ms were considered as short fixations, fixations between 200 ms and 400 ms were medium fixations, and long fixations had a duration of at least 400 ms. For calculating fixation dispersion, we chose our dispersion area to be 75%, i.e., the area containing 75% fixations. The threshold of saccade size varies on devices: saccades were considered short if their length is less than 200 px on desktop or less than 75 px on mobile, long if length greater than 400 px (desktop) or 175 px (mobile), otherwise medium. The choice of such thresholds was based on the line width, short saccades went through less than 1/4 of the line, while long saccades took approximately half of the line. Moreover, saccades were categorized into eight directions. For instance, a saccade was in up-direction if it was pointed towards up-direction and the angle between it and the vertical line was at most 22.5 degrees. Two consecutive saccades were counted as one follow if they had the same direction. Similar reasoning could be applied to neighbor, opposite, and opposite-neighbor features.

Our mid-level features aligned with Srivastava et al. [44], except we chose a small line to be three right saccades (of any size) followed by one left long saccade, i.e., ArArArLl, where Ar stands for a right saccade of any length, and Ll stands for a left long saccade. Similarly, we defined a medium line to be four right followed by one long left, i.e., ArArArArLl, and a long line to be five right followed by one long left, i.e., ArArArArArLl.

377 **4.3 Classification**

378 To build our classifiers, we first took the features extracted in the previous section from a given window and
 379 use them to predict the reading condition, i.e., deep versus skim. We experimented with Kernel Support Vector
 380 Machines (SVM), Logistic Regression, Random Forest, and XGBoost as our classifiers. We also compared to a
 381 majority vote baseline classifier as the data was imbalanced (75% deep versus 25% skim). We treated the window
 382 size and step size for the feature extraction as hyper-parameters to be optimized as part of the model training.
 383

384 In addition to reading mode classification, we used the same classifier set-up to predict reading devices, i.e.,
 385 which device the participant was reading at. This task would be trivial with saccade-related features, due to the
 386 huge size difference between desktop and mobile. Hence, we only used fixation and pupil features alone. We also
 387 explored models for cross-device reading mode classification, where we trained the models using data from one
 388 device, and tested the models using data from the other device.

389 Hyper parameter tuning is often done a validation set prior to a model being evaluated on a test set. However,
 390 this would significantly reduce the available training data in our setting so we instead used a leave-one participant-
 391 out evaluation using nested cross-validation for the hyperparameter tuning [13]. In nested cross-validation, the
 392 dataset is firstly divided into k parts, such that in each round, one part is picked as the testing set while the
 393 rest are picked as the training set. Furthermore, before evaluating on the test set, a nested k' -cross-validation is
 394 performed within the training set for hyper-parameter tuning. Finally, the best hyper-parameters with the best
 395 AUC (chosen over accuracy due to data imbalance) during cross validation were chosen to evaluate on the test
 396 set. By doing this, the hyper-parameter tuning process does not have access to or is not exposed to the test set,
 397 hence the result is less biased. To evaluate the performance of each model, we used the nested cross-validation
 398 with $k = 29$ and $k' = 10$ and report accuracy, F1 score, and area under the curve (AUC) as our metrics.
 399

400 For reproducability we next describe the parameters used during these experiments. During the hyper-
 401 parameter tuning, we varied the window size from 5 to 120 seconds with a step of 5 seconds. We also varied the
 402 step size from 10, 25, 50, to 100 percentage of the current window size. We set the model hyper-parameters based
 403 on default values and previous similar studies so they were not tuned within the cross-validation. For Kernel
 404 SVM, we chose RBF kernel with $C = 10$ and gamma equals 0.01. For Logistic Regression, we set C to be 1 and
 405 penalty to be l_2 . For Random Forest, we set the number of estimators to be 1000, the maximum feature parameter
 406 to be the square root of feature count, and the minimum number of samples required for internal node splitting
 407 as 2, for leaf splitting as 1. For XGBoost, we set the number of estimators to be 100, the maximum depth to be 3,
 408 gamma to be 0, subsample and subsample ratio of columns to be 1, regularization alpha to be 0, and learning rate
 409 to be 0.1.
 410

411 **5 RESULTS**

412 In this section, we first check the validity of our study design, i.e., whether our task design successfully induced
 413 deep and skim reading. Then, we present our results for the reading mode classification, where we demonstrate
 414 that our classification methods can effectively differentiate between deep reading and skim reading. Furthermore,
 415 we analyze the impact of window size, step size, and mid-level gaze features on our classifiers. We also gained
 416 more insights by looking into important features for classification and misclassified instances. Lastly, we explored
 417 the role of gaze features in some other classification tasks, such as reading device classification, and cross-device
 418 mode classification.

419 **5.1 Comprehension Score and Reading Speed**

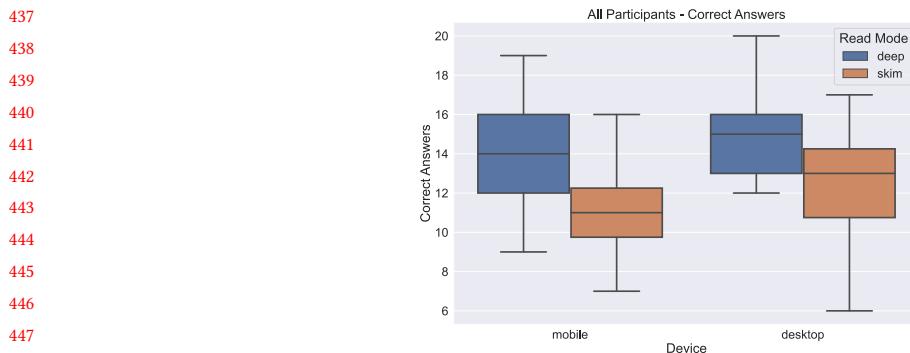
420 Since the deep and skim reading behaviors were induced by our study design, we started by validating it using
 421 participants' comprehension results and reading behaviors. As shown below, during deep reading sessions, most
 422 participants read at a speed that prioritized comprehension and were thus able to achieve a high comprehension
 423

424 scores for all three comprehension levels (i.e., literal, inferential, and evaluative). While in skim reading, participants
 425 read faster (i.e., performing skimming or scanning) but still obtained high-level literal comprehension.
 426 However, their inferential comprehension was significantly poorer in skim reading. Therefore, we concluded that
 427 our study design had successfully induced the expected behavior.

428 Figure 3 shows the total number of correct answers across all participants. In both devices, participants had
 429 significantly better comprehension results in deep reading (desktop $M = 14.89, SD = 1.72$; mobile $M = 14.08, SD =$
 430 2.49) than skim reading (desktop $M = 12.47, SD = 2.88$; mobile $M = 11.11, SD = 2.23$). Paired t-test revealed
 431 statistically significant differences between the conditions for desktop ($t(35) = 4.20, p < 0.001$) as well as mobile
 432 ($t(35) = 5.18, p < 0.001$). Also, when reading shallowly, participants had significant better comprehension on
 433 desktop ($M = 12.47, SD = 2.88$) than mobile ($M = 11.11, SD = 2.23$) (paired t-test $t(35) = 2.37, p = 0.02$). However,
 434 there is no such difference could be observed when reading deeply (paired t-test $t(35) = 1.80, p = 0.08$).

435

436



447 Fig. 3. Comprehension Results of All Participants

448

449

450

451

452 Furthermore, we looked into the comprehension results for each question type, i.e., literal, inferential, and
 453 evaluative. The literal questions were further divided into two sets, the set of questions that were previewed to
 454 participants (Previewed Lit Q), and those that were not. Figure 4 shows the comprehension results for each question
 455 type. When comparing deep versus skim reading on each device, participants showed significantly better results
 456 for non-previewed literal questions (paired t-test: desktop $t(35) = 3.01, p < 0.01$; mobile $t(35) = 4.00, p < 0.01$)
 457 and inferential questions (paired t-test: desktop $t(35) = 6.03, p < 0.001$; mobile $t(35) = 3.20, p < 0.005$), but
 458 not for previewed literal questions (paired t-test $t(35) = 0.81, p > 0.1$) and evaluative questions (paired t-test
 459 $t(35) = 0.53, p > 0.1$) on desktop. When comparing desktop and mobile reading, participants obtained better
 460 comprehension for evaluative questions in skim reading (paired t-test $t(35) = 2.39, p < 0.05$) and inferential
 461 questions in deep reading (paired t-test $t(35) = 2.20, p < 0.05$), but not for others.

462

463

464

465

466

467

468

469

470

469 In addition, we investigated the reading speed of each participant, we found participants read at a much
 470 faster speed in skim reading (paired t-test $t(35) = 9.29, p < 0.001$). As shown in Figure 5, most participants
 471 read at the speed normally considered as reading for comprehension (i.e., 300 words per minute [12]) (desktop
 472 $M = 296.58, SD = 118.38$; mobile $M = 306.37, SD = 131.74$) during the deep reading session. However, in skim
 473 reading, most participants read at the speed normally considered as skimming (i.e., 450 words per minute [12])
 474 and scanning (600 words per minute [12]) (desktop $M = 602.93, SD = 221.38$; mobile $M = 679.35, SD = 326.73$).

475 Two typical examples of deep and skim readings were shown in Figure 6. We plotted the gaze map (i.e., fixations
 476 and saccades) and heatmap overlaid over the reading articles for each condition. As shown in Figure 6(a), when

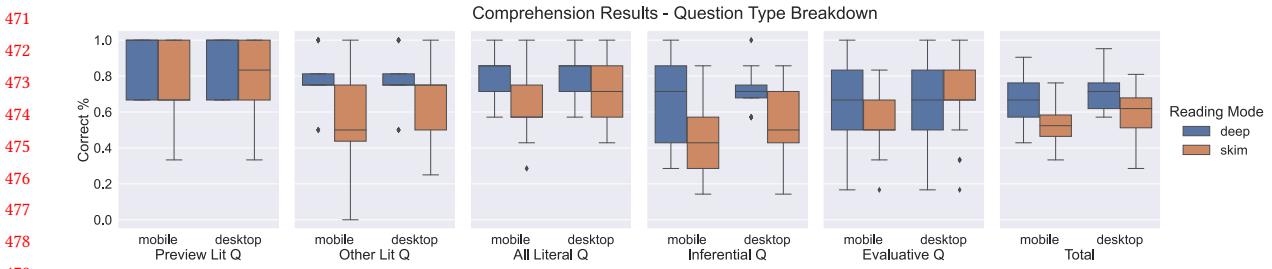


Fig. 4. Comprehension Results with Question Type Breakdown

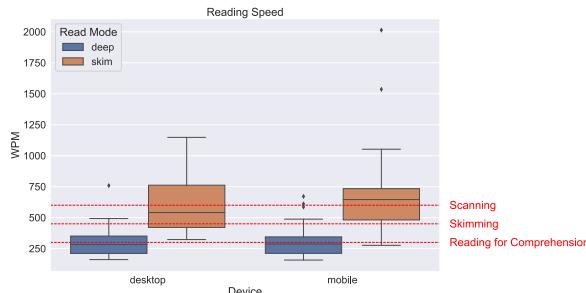


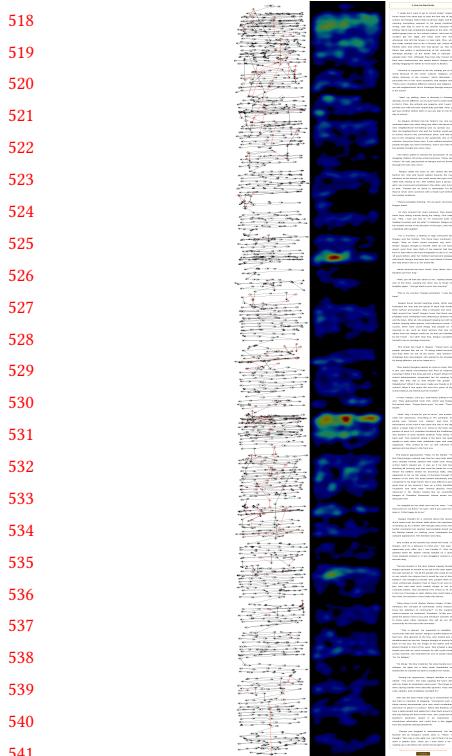
Fig. 5. Reading Speed of All Participants.

reading deeply, there were more fixations and more short saccades, hence the gaze plot looks relatively dense. On the other hand, when reading shallowly (Figure 6(b)), there were fewer fixations but more long saccades, also the regressions appeared more frequently. It is worth noticing that as reflected by the skim-reading heatmap, participants paid much attention on the AOIs of previewed questions. While in deep readings, participants did have particularly focused regions but the overall distribution was more spread out. Another interesting finding is that for each paragraph in deep reading, many participants put more effort into reading at the start of the paragraph than at the end of the paragraph.

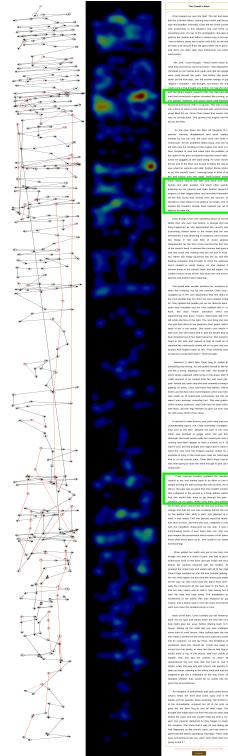
5.2 Reading Mode Classification

5.2.1 Classification Results. The performances of all models are summarized in Table 5, where we used bold text to mark the best classifier regarding one evaluation metric. Note that our baseline, a majority guesser, achieved a relatively high accuracy score due to the imbalanced data from larger number of windows from the longer deep reading sessions. This imbalance makes accuracy less reflective of performance than the AUC and F1 metrics.

All models outperformed our baseline approach, i.e., the majority guesser. Among all the models, XGBoost achieved the best performance in desktop with an accuracy of 88.36%, a F1-score of 0.82, and an AUC of 0.82. In mobile, Logistic Regression had the highest accuracy (86.18%) and F1-score (0.77). Meanwhile, the XGBoost had similar accuracy and F1-score, and it achieved the highest AUC score (0.73). The classifiers typically chose window sizes from 60 to 120 seconds, as well as a step size of 100% during the optimization. During training we also tried variants that balanced the training data via down-sampling the majority class or up-sampling the minority class to account for the imbalance. However, both methods failed to improve our model in terms of our evaluation metrics.



(a) Deep Reading Example (p11, desktop, deep)



(b) Skim Reading Example (p11, desktop, skim)

Fig. 6. Examples of Deep and Skim Readings. From left to right: gaze movement plot, gaze heatmap, and reading article. In skim reading, the AOIs of previewed questions were marked with green boxes. Regressions were marked as red lines.

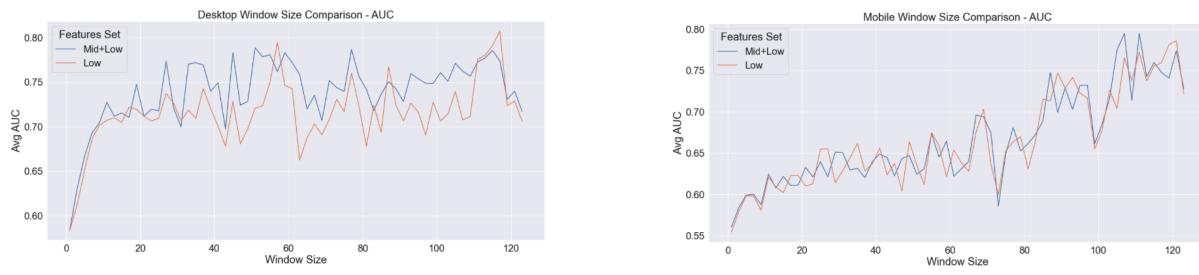
5.2.2 Window Size and Step Size. To better understand the effect of step size and window size on the classifier, we examined them in detail for the best-performed model, XGBoost. We first varied the window size from 1 to 120 seconds fixing the step size to be equal to the windows size. Next we, took the best window size, and then varied the step size from 1 second to the best window size. For simplicity, we plotted only the accuracy and AUC scores, the F1 score graphs were very similar to the AUC graphs.

Figure 7 summarized the classification results of XGBoost model when varying the window size for window extraction. For desktop, the model achieved a relatively good performance with any window size greater than 20 seconds. When we decreased the window size from 10 to 1, performance dropped significantly, indicating a small window size failed to capture any meaningful features. Meanwhile, window sizes had a much bigger influence on mobile. Performance increased as we increased the window size from 1 to 120. Hence a bigger window size might be required for effective mobile classification. Moreover, desktop outperformed mobile on window sizes below 80 seconds, and they performed roughly the same with larger window sizes.

According to the result of our window size evaluation, we fixed the window size to be 120 seconds and varied the step size to measure its influence. Our results for step size comparison was summarized in Figure 8, where we still used XGBoost as our classifier. Similar trends were observed in both desktop and mobile. When increasing the

Table 5. Reading Mode Classification Results

Device	Model	Accuracy	Macro F1 Score	AUC
Desktop	Baseline (Majority Guessing)	0.7733	0.4361	0.5
	Linear SVM	0.8366	0.7501	0.7598
	RBF SVM	0.8824	0.8098	0.7811
	Logistic Regression	0.8797	0.8190	0.8042
	Random Forest	0.8623	0.7649	0.7392
Mobile	XGBoost	0.8836	0.8239	0.8201
	Baseline (Majority Guessing)	0.7666	0.4339	0.5
	Linear SVM	0.8320	0.7414	0.7222
	RBF SVM	0.8279	0.6922	0.6617
	Logistic Regression	0.8618	0.7688	0.7307
	Random Forest	0.8455	0.7416	0.7082
	XGBoost	0.8551	0.7558	0.7311



(a) Desktop Window Size Evaluation – AUC.

(b) Mobile Window Size Evaluation – AUC.

Fig. 7. Effect of Window Size on Classifier Performance

step size from 1 second to 120 seconds, performance of both devices increased. Adding more data via overlapping windows does not help result in a better model.

5.2.3 Important Features. We further analyzed the important features outputted by XGBoost. Table 6 listed the top-ten features, which were chosen based on the proportion of samples they could split with highly correlated features being removed. When the window size is small (i.e., 5 seconds), fixation and pupil features were more important, showing that statistics about fixation duration and pupil diameter was crucial in differentiating deep and skim reading. One possible explanation is that the window size was too small to pick up any meaningful movements, hence the model mainly used statistics features. As the window size increased, more and more saccade-related features were identified as important. When window size equals 60 seconds, more than half of the top ten features were saccade features.

Among them, features that are characteristic of “backtracking” were valued more, this includes medium-length saccades towards up, up-right, and up-left. Fixation duration features remained relatively important in large window sizes, but pupil diameter features became less important. Moreover, the classification process relied

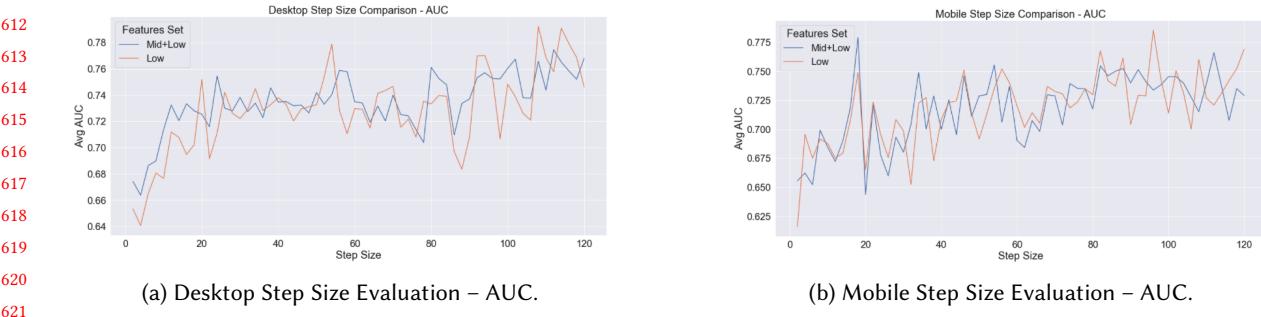


Fig. 8. Effect of Step Size on Classifier Performance

Table 6. Top 10 Important Features in Mode Classification at Different Window Sizes. In pupil diameters, we used (L) for left eye and (R) for right eye. And we omitted “saccades” in saccade-related features, such as Up for Up saccades and Left for left saccades. Also, in fixation and saccade features, we use (S) for short, (M) for Medium, and (L) for Long.

Rank	5 Sec Window		15 Sec Window		60 Sec Window		120 Sec Window	
	Desktop	Mobile	Desktop	Mobile	Desktop	Mobile	Desktop	Mobile
1	Else-where	Fix slope	Else-where	Fix var	Else-where	Up-right (M)	Else-where	Up-right (M)
2	Pup mean (L)	Pup mean (L)	Pup mean (L)	Fix max	Up-left (M)	Up-left (M)	Up-left (M)	Fix slope
3	Pup mean (R)	Pup mean (R)	Sacc mean	Fix slope	Up-right (M)	Up (M)	Up (M)	Up (M)
4	Sacc mean	Sacc mean	Pup mean (R)	Pup mean (L)	Left (M)	Regression	Left (M)	Up-left (M)
5	Pup var (L)	Fix var	Left (M)	Pup mean (R)	Up (M)	Long Fix	Neighbor	Left (M)
6	Pup std (L)	Fix max	Fix var	Up-right (M)	Oppo-neigh	Fix min	Up-right (M)	Long Fix
7	Fix slope	Fix std	Fix max	Sacc mean	Fix max	Fix var	Oppo-neigh	Regression
8	Fix var	Fix min	Fix min	Fix std	Fix var	Fix slope	Fix var	Right (M)
9	Fix max	Fix mean	Fix std	Fix min	Med fix	Fix max	Fix std	Fix min
10	Fix rate	Pup std (R)	Pup var (L)	Else-where	Neighbor	Left (M)	Fix max	Fix var

heavily on distance-based mid-level gaze features such as else-where patterns and regression patterns, but not shape-related features. In addition, there are more important fixation duration features in mobile reading than desktop readings, meaning they are more important in mobile classification. This might be because the screen size of mobile is much smaller than desktop, hence it is harder to identify and extract useful saccade patterns.

5.2.4 Example Instances. Besides the important features, we also looked into each instance (i.e., each window for each participant) and the classification results of XGBoost with a window size and a step size of 60 seconds. Figure 9 shows the typical correctly classified and misclassified instances. Figure 9(a) and Figure 9(b) gave the correctly classified deep and skim reading window. Figure 9(c) showed one of our two misclassified deep windows, which occurred when participant 14 skimmed the article after reading it thoroughly. Figure 9(d) demonstrated one of the misclassified skim windows, and as shown in the figure, participant 25 were actually reading at a very slow speed which appears to be more like deep reading itself.

5.3 Other Classification Tasks

5.3.1 Reading Device Classification. In device classification, we checked if we can predict the reading devices using various feature set. The window size and step size were fixed to 60 seconds, and we varied the feature set

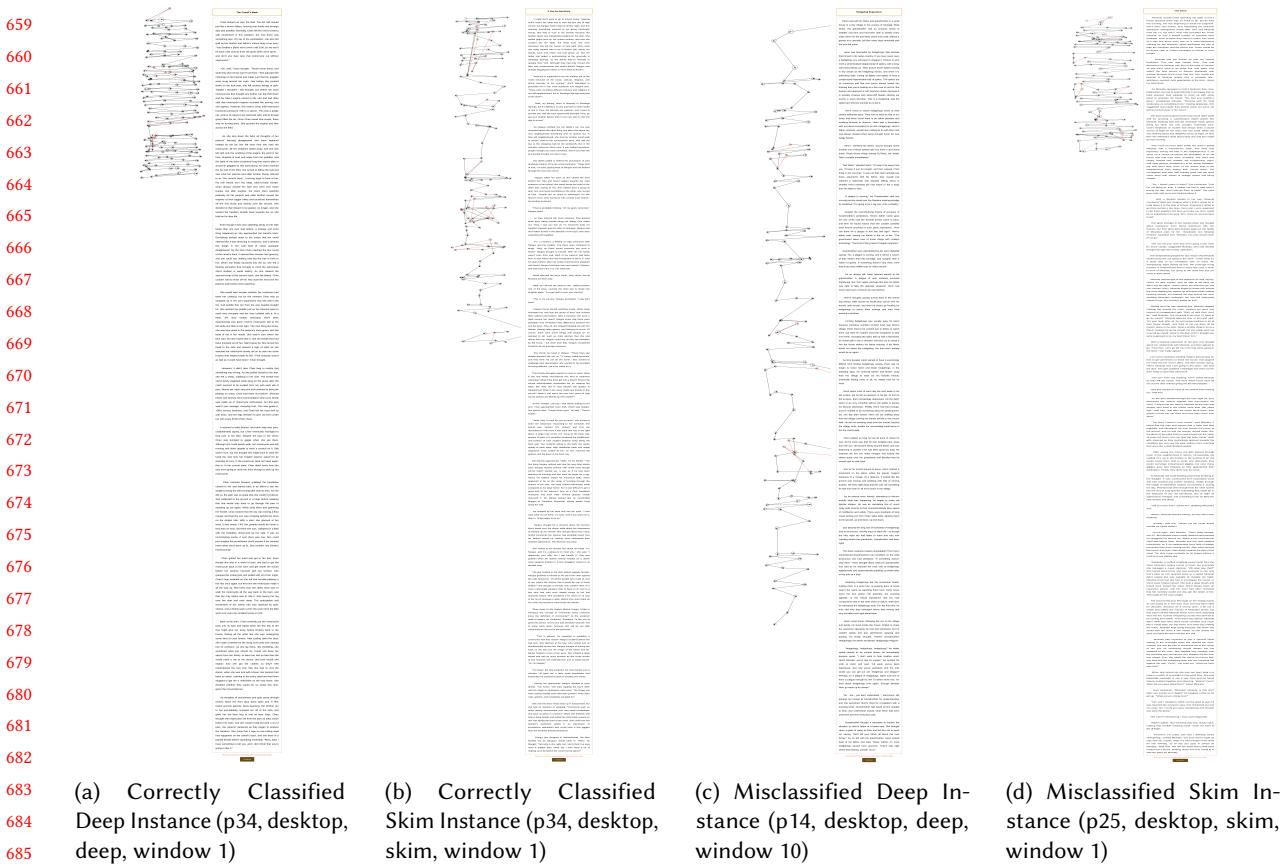


Fig. 9. Classification Results Example Instances

as fixation-related features only, fixation and pupil features, all low-level features, and low-level plus mid-level features. As shown in Table 7, this classification task was trivial when saccade-related features were allowed. The best classifier, Random Forest, achieved an accuracy of 100% with low-level features and low-level plus mid-level features. When using fixation features alone, we could still obtain 69.68% accuracy with Random Forest. Moreover, the classification accuracy could be boosted to 84.16% if pupil features were allowed and XGBoost was used for classification. The usefulness of pupil features could be reflected in Table 8, pupil features were the top two, 7th, and 8th important features.

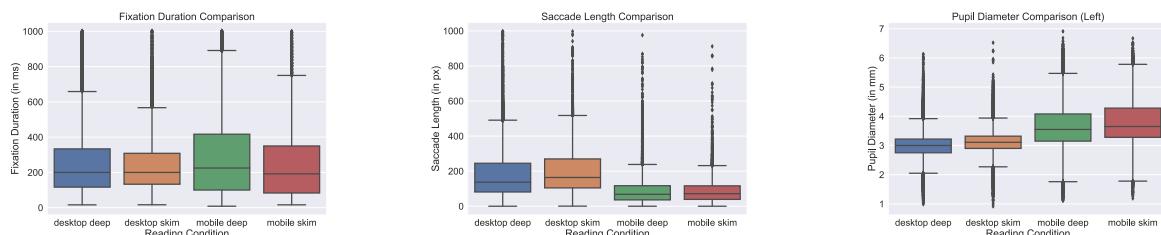
5.3.2 Cross-device Mode Classification. In inter-device mode classification, we failed to obtain a model that significantly outperformed the baseline approach. However, this is expected and can be explained by the inherent difference between devices. Figure 10 demonstrated the statistics summary of fixation duration, pupil diameter, and saccade length in all four conditions. For the same reading device, all features exhibited some differences between deep and skim reading. For the same reading mode, there were a huge difference between desktop and mobile in all features. These difference explained why our device and mode classification succeed, and inter-device mode classification failed.

706
707
708 Table 7. Device Classification Accuracy.
709

Model	Fixation Only	Fixation + Pupil	Low Level	Low + Mid Level
Baseline (Majority Guessing)	0.5111	0.5111	0.5111	0.5111
RBF SVM	0.6610	0.6525	0.9949	0.9966
Logistic Regression	0.6678	0.6746	0.9915	0.9915
Random Forest	0.6968	0.8262	1.0000	1.0000
XGBoost	0.6934	0.8416	0.9949	0.9966

716
717 Table 8. Top 10 Important Features in Device Classification.
718

Rank	Fix Only	Fix + Pupil	Low level	Low + Mid level
1	Fix slope	Pup mean (L)	Right (M)	Right (M)
2	Med Fix	Pup mean (R)	Sacc var	Else-where
3	Fix std	Fix max	Sacc std	Sacc var
4	Fix max	Fix var	Left (M)	Left (M)
5	Fix var	Med fix	Sacc mean	Sacc std
6	Fix rate	Fix slope	Left (L)	Sacc mean
7	Fix mean	Pup std (R)	Right (L)	Left (L)
8	Fix min	Pup var (R)	Down-right (M)	Right (L)
9	Short fix	Fix min	Left (S)	Down-right (M)
10	Fix per	Fix mean	Opposite	Up-left (M)

734
735 (a) Fixation Duration Comparison.736
737 (b) Saccade Length Comparison.738
739 (c) Pupil Diameter Comparison.

740 Fig. 10. Fixation, Saccades, and Pupil Data Summary

741 6 DISCUSSION

742 To advance reading research we need to be able to track reading behaviours in a scalable and non-intrusive
 743 manner. To develop implicit metrics about the quality of reading, we created a method to reliably induce skim
 744 and deep reading using three levels of comprehension. Since reading increasingly takes place electronically [40],
 745 we developed classifiers for both desktop and mobile devices.

753 6.1 Inducing and Validating Deep vs. Skim Reading

754 Our evaluation showed that inferential comprehension is most prominently linked to deep reading activities
755 regardless of reading device as it requires readers to empathise with the context of a story in order to derive
756 insights not explicitly stated in the text. Inferential comprehension is essential for passage understanding as it
757 requires readers to connect events in a narrative or understand a character's motive for a particular action [5].
758 The tasks we developed to induce deep and skim reading were most effective in their discriminative power with
759 regard to this crucial level of comprehension.

760 Readers' ability to answer non-reviewed literal and inferential questions was further significantly impeded
761 by skim reading on both desktop and mobile. For previewed and evaluative questions we did, however, not see a
762 difference between reading modes, which may imply that skim readers have not just blindly skipped text but
763 applied skim reading as a selective form of text processing.

764 When reading on mobile devices, our results show better comprehension for evaluative questions in skim
765 reading mode. Especially high-level questions like "What kind of person is Kim?" could be answered well during
766 skim reading. This may indicate that readers have successfully adapted their reading ability to quickly determine
767 how new information fit into their existing knowledge network and make use of prior experiences to quickly
768 judge a passage or character.

769 Our method has been shown to effectively induce skim and deep reading and thus can be reused by researchers
770 to further the research on these two reading modes. The reading materials can be found in the Appendix along
771 with the comprehension questions and their groupings into three levels. We further open-sourced our reading
772 interface, which is listed in our Git repository ³.

773 6.2 Linking Deep Reading to Eye Gaze Pattern

774 By combining participants' comprehension results and their reading speed with eye movement data we were
775 able to determine eye gaze patterns that are characteristic of deep reading activities. In general, participants
776 read slowly at a speed for comprehension in deep reading, and relatively fast when doing skim reading. By
777 looking at the gaze plot and heatmap of each participant, we found that participants tend to read line by line
778 carefully in deep reading, i.e., there were much more short saccades and less long saccades, and the overall reading
779 patterns looked more like scanning slowly from left to right, line by line. While in skim reading, we found that
780 participants jump quickly between paragraphs and only spent relatively more time on texts related to previewed
781 questions. Their reading patterns were more like scanning in a zig-zag manner with longer saccades and less
782 fixations. In addition, these features played an important role in the later classification. We found that, classifiers
783 saw a task as deep reading if it had less upwards and left saccades (i.e., backtracking-like saccades), more and
784 longer fixations, and smaller pupil diameter. The bigger the window size, the more important saccade-related
785 features were. Desktop classifiers utilized saccade features more, as it were easier to pick up on larger screen size.
786 Moreover, distance-based mid-level features were also considered by the classifiers as they identified regression
787 patterns.

788 6.3 Reading Mode Classification

789 Our results demonstrate the existence of a strong relationship between eye movements and reading mode. This
790 allowed us to build classifiers to distinguish deep from skim reading. One of the central contributions of our work
791 is to provide reliable classifiers for reading on desktop as well as on mobile devices.

792 Our classifiers achieved high performance in recognizing deep versus skim reading in general settings. In
793 desktop classification, among all four classifiers, the XGBoost achieved an accuracy of 88.36%, a F1 score of
794 0.82 and an AUC of 0.82, which outperformed our baseline measurement, majority guesser (accuracy = 77.33%,

³link anonymized

800 F1 score = 0.44 and AUC = 0.50), by the most. Note our baseline approach achieved a high accuracy due to
 801 imbalanced dataset. In mobile, Logistic Regression performed the best with accuracy = 86.18%, F1 score = 0.77
 802 and AUC = 0.73. XGBoost achieved a similar performance with accuracy = 85.51%, F1 score = 0.76 and AUC =
 803 0.73. The majority guesser had accuracy = 76.66%, F1 score = 0.43, and AUC = 0.50.

804 Misclassifications may be due to the large variance among individuals' reading behaviors, which has been
 805 widely reported about [36], and the natural mix of deep and skim reading in daily reading activities. This was
 806 verified by looking at each individual misclassified instance. During deep reading session, some participants (e.g.,
 807 participant 14 and 32) still performed a "revision" after the reading the entire article once, and the "revision" was
 808 more like skimming and scanning than deeply reading. Similarly, when reading shallowly, some participants (e.g.,
 809 participant 22, 25, 36, etc.) were not able to perform skim reading due to their deep engagement with the presented
 810 stories, hence they still read slowly and deeply. Moreover, the performance difference between desktop and
 811 mobile may be exacerbated by the smaller screen size and closer distance to the eye tracker. Hence the abilities
 812 of both our eye tracking hardware (i.e., eye tracker's accuracy in terms of correctly identifying fixations and
 813 saccades) and software (i.e., correctly picking up useful movements by feature extraction) were limited. To further
 814 improve performance on mobile, more fine-grained eye-tracking technique and mobile-tailored feature set were
 815 required. Moreover, we experimented with various window size and step size and found that a size of 120 seconds
 816 is the most suitable one for both. This aligns with previous studies [10, 27, 29] indicating that a large window
 817 size is required for such high-level activity recognition. A shorter window size may not be sufficient to capture
 818 important movements. The claim on shorter window was confirmed by important features outputted by XGBoost
 819 – when we increased the window size from 5 to 60 seconds, fixation and pupil features were out-weighted by
 820 saccades and mid-level features, indicating the presence of movement-related features and their importance to
 821 our model. Our results on step size indicated that adding more data via sliding window overlapping did not work
 822 well. To obtain a more robust and accurate, more data needs to be collected and analyzed. In addition, we found
 823 that statistics about pupil diameter size played an important role in our model, which was not often considered
 824 in previous works. Also, we found that less saccade features were considered as crucial in mobile. This aligns
 825 with our suspicion that our hardware and software ability was limited on the mobile device.

826 Additionally, our results on device classification showed that we could tell the reading device effectively
 827 (accuracy 69.68%) by using fixation features along. It was mainly because fixation duration was shorter in desktop
 828 reading. One possible explanation for this is the text size. In order to keep the same amount of characters on each
 829 line, the text size on mobile was smaller than desktop. Previous studies revealed that the smaller the font size, the
 830 longer the fixation duration [7, 37]. Hence it was expected that fixation duration could help identify deep and
 831 skim reading. With the help of pupil features, the performance of our device classifier could be boosted to 84.16%.
 832 The pupil features were so useful because the pupil diameter was larger and had larger variance in mobile. This
 833 might be because of the text size as well as screen settings such as brightness. Further study was required to
 834 examine it. The task became trivial when considering saccade features, due to the huge difference in screen size
 835 between desktop and mobile.

836

837

838 6.4 Limitations

839 First of all, as an in-lab study participants read under controlled environmental and device conditions. Additionally,
 840 we imposed several soft constraints on participants for accurate data recording (e.g., do not move head too much).
 841 Hence, a lot of varying environmental influences (e.g., distractions, differences in lighting, etc.) encountered
 842 during everyday reading activities are not reflected in our data. The stationary set-up of using a Tobii eye tracker
 843 under controlled condition was a necessary first step to validate the successful induction of deep and skim reading
 844 as well as link deep reading activities to eye gaze patterns. The resulting classifiers build the basis for building
 845

846

847 more robust and mobile tracking solutions, which we also plan to expand to other sensor types (e.g., eye tracking
 848 through the front-facing rgb camera).

849 Moreover, in our study design we limited the skim reading time to two minutes, which may have caused
 850 changes in reading behavior due to the explicit time pressure imposed. The high scores in literal comprehension,
 851 however, indicate that participants were able to successfully undertake the reading tasks.

852 Regarding our feature selection, we suspect that a better set could be obtained if we tailored our features
 853 specifically for reading task, such as developing more appropriate mid-level and high-level features, which we
 854 left for future work. Moreover, since a large window size was required for accurate classification, our current
 855 approach is more appropriate for post-hoc than real-time classification, which we deem feasible considering our
 856 goal of tracking reading behaviour changes long-term.

857

858 6.5 Future Directions

859 In future work, we will explore more accurate and less constrained methods for classifying reading behaviour
 860 in-the-wild. By improving our classifiers through richer mid-level features along with collecting larger datasets,
 861 we intend to build more robust classifiers that run with alternative sensors, such as an rgb camera or Apple's
 862 true-depth camera. This would allow us to take the tracking of deep reading episodes into the field by deploying
 863 our models on consumer devices.

864 By providing a method to induce deep and skim reading activities along with the link between eye gaze patterns
 865 and deep reading, our work paves the way towards implicitly tracking reading behaviour over time. Follow-up
 866 studies that we can now undertake include investigating the effect of reading content, layout, or style on reading
 867 behaviour in natural settings. For example, are users conditioned to various reading behaviours based on the
 868 website they visit (i.e. social media vs news feeds vs blogs)? Similarly, with robust classification we can study
 869 different settings where reading behaviours may be influenced by contextual factors, such as while at work,
 870 home, or commuting.

871

872 7 CONCLUSIONS

873 Understanding reading behaviour on desktop and mobile devices through eye-movement is important to foster
 874 and understand deep reading in natural settings. Our study is the first to perform a detailed investigation into
 875 reading modes for both desktop and mobile devices. This was enabled through a study design we introduced
 876 that could successfully induce people's deep and skim reading behavior as evidenced through comprehension
 877 results and reading speed. Moreover, we showed that these two reading modes could be effectively identified
 878 with classifiers trained on features derived from eye-movement data. Finally, we discussed key characteristics of
 879 deep and skim readings through the analysis of the most important features for our classifiers. Our work paves
 880 the way to study long-term changes in reading behaviour through implicit metrics in natural settings.

881

882

883

884

885

886

887

888

889

890

891

892

893

894 REFERENCES

- 895 [1] Yomna Abdelrahman, Anam Ahmad Khan, Joshua Newn, Eduardo Velloso, Sherine Ashraf Safwat, James Bailey, Andreas Bulling, Frank
 896 Vetere, and Albrecht Schmidt. 2019. Classifying attention types with thermal imaging and eye tracking. *Proceedings of the ACM on*
 897 *Interactive, Mobile, Wearable and Ubiquitous Technologies* 3, 3 (2019), 1–27.
- 898 [2] Peter Aflerbach. 2015. *Handbook of individual differences in reading: Reader, text, and context*. Routledge.
- 899 [3] Seoyoung Ahn, Conor Kelton, Aruna Balasubramanian, and Greg Zelinsky. 2020. Towards predicting reading comprehension from gaze
 900 behavior. In *ACM Symposium on Eye Tracking Research and Applications*. 1–5.
- 901 [4] J Alonzo, G Tindal, K Ulmer, and A Glasgow. 2006. easyCBM online progress monitoring assessment system. *Eugene, OR: Center for*
Educational Assessment Accountability (2006).
- 902 [5] Mary DeKonty Applegate, Kathleen Benson Quinn, and Anthony J Applegate. 2002. Levels of thinking required by comprehension
 903 questions in informal reading inventories. *The Reading Teacher* 56, 2 (2002), 174–180.
- 904 [6] Deni Basaraba, Paul Yovanoff, Julie Alonzo, and Gerald Tindal. 2013. Examining the structure of reading comprehension: do literal,
 905 inferential, and evaluative comprehension truly exist? *Reading and Writing* 26, 3 (2013), 349–379.
- 906 [7] D Beymer, D Russell, and P Orton. 2008. An eye tracking study of how font size and type influence online reading. *People and computers*
 907 XXII: culture, creativity, interaction: proceedings of HCI 2008. In *The 22nd British HCI Group annual conference*, Vol. 2.
- 908 [8] Ralf Biedert, Jörn Hees, Andreas Dengel, and Georg Buscher. 2012. A robust realtime reading-skimming classifier. In *Proceedings of the*
Symposium on Eye Tracking Research and Applications. 123–130.
- 909 [9] Mark R Blair, Marcus R Watson, R Calen Walshe, and Fillip Maj. 2009. Extremely selective attention: eye-tracking studies of the dynamic
 910 allocation of attention to stimulus features in categorization. *Journal of Experimental Psychology: Learning, Memory, and Cognition* 35, 5
 (2009), 1196.
- 911 [10] Andreas Bulling, Jamie A Ward, Hans Gellersen, and Gerhard Tröster. 2010. Eye movement analysis for activity recognition using
 912 electrooculography. *IEEE transactions on pattern analysis and machine intelligence* 33, 4 (2010), 741–753.
- 913 [11] Douglas Carnine, Jerry Silbert, Edward J Kameenui, and Sara G Tarver. 1997. Direct instruction reading. (1997).
- 914 [12] Ronald P Carver. 1990. *Reading rate: A review of research and theory*. Academic Press.
- 915 [13] Gavin C Cawley and Nicola LC Talbot. 2010. On over-fitting in model selection and subsequent selection bias in performance evaluation.
The Journal of Machine Learning Research 11 (2010), 2079–2107.
- 916 [14] Leana Copeland and Tom Gedeon. 2013. Measuring reading comprehension using eye movements. In *2013 IEEE 4th International*
917 Conference on Cognitive Infocommunications (CogInfoCom). IEEE, 791–796.
- 918 [15] Leana Copeland, Tom Gedeon, and B Sumudu U Mendis. 2014. Predicting reading comprehension scores from eye movements using
 919 artificial neural networks and fuzzy output error. *Artif. Intell. Res.* 3, 3 (2014), 35–48.
- 920 [16] Leana Copeland, Tom Gedeon, and Sumudu Mendis. 2014. Fuzzy output error as the performance function for training artificial neural
 921 networks to predict reading comprehension from eye gaze. In *International Conference on Neural Information Processing*. Springer,
 586–593.
- 922 [17] Mary C Dyson. 2004. How physical text layout affects reading from screen. *Behaviour & information technology* 23, 6 (2004), 377–393.
- 923 [18] Ralf Engbert and Reinhold Kliegl. 2003. Microsaccades uncover the orientation of covert attention. *Vision research* 43, 9 (2003),
 1035–1045.
- 924 [19] Maureen P Hall, Aminda O'Hare, Nicholas Santavicca, and Libby Falk Jones. 2015. The power of deep reading and mindful literacy: An
 925 innovative approach in contemporary education. *Innovación educativa (México, DF)* 15, 67 (2015), 49–60.
- 926 [20] Harold L Herber. 1978. *Teaching reading in content areas*. Prentice Hall.
- 927 [21] Edmund Burke Huey. 1908. The psychology and pedagogy of reading: With a review of the history of reading and writing and of
 928 methods, texts, and hygiene in reading. (1908).
- 929 [22] Shoya Ishimaru, Kensuke Hoshika, Kai Kunze, Koichi Kise, and Andreas Dengel. 2017. Towards reading trackers in the wild: detecting
 930 reading activities by EOG glasses and deep neural networks. In *Proceedings of the 2017 ACM International Joint Conference on Pervasive*
and Ubiquitous Computing and Proceedings of the 2017 ACM International Symposium on Wearable Computers. 704–711.
- 931 [23] Halszka Jarodzka and Saskia Brand-Gruwel. 2017. Tracking the reading eye: Towards a model of real-world reading.
- 932 [24] Marcel A Just and Patricia A Carpenter. 1980. A theory of reading: From eye fixations to comprehension. *Psychological review* 87, 4
 (1980), 329.
- 933 [25] Shian-Ru Ke, Hoang Le Uyen Thuc, Yong-Jin Lee, Jenq-Neng Hwang, Jang-Hee Yoo, and Kyoung-Ho Choi. 2013. A review on video-based
 934 human activity recognition. *Computers* 2, 2 (2013), 88–131.
- 935 [26] Conor Kelton, Zijun Wei, Seoyoung Ahn, Aruna Balasubramanian, Samir R Das, Dimitris Samaras, and Gregory Zelinsky. 2019. Reading
 936 detection in real-time. In *Proceedings of the 11th ACM Symposium on Eye Tracking Research & Applications*. 1–5.
- 937 [27] Peter Kiefer, Ioannis Giannopoulos, and Martin Raubal. 2013. Using eye movements to recognize activities on cartographic maps. In
 938 *Proceedings of the 21st ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*. 488–491.

- 941 [28] Reinhold Kliegl, Antje Nuthmann, and Ralf Engbert. 2006. Tracking the mind during reading: The influence of past, present, and future
 942 words on fixation durations. *Journal of experimental psychology: General* 135, 1 (2006), 12.
- 943 [29] Kai Kunze, Yuzuko Utsumi, Yuki Shiga, Koichi Kise, and Andreas Bulling. 2013. I know what you are reading: recognition of document
 944 types using mobile eye tracking. In *Proceedings of the 2013 international symposium on wearable computers*. 113–116.
- 945 [30] Silvia Makowski, Lena A Jäger, Ahmed Abdelwahab, Niels Landwehr, and Tobias Scheffer. 2018. A discriminative model for identifying
 946 readers and assessing text comprehension from eye movements. In *Joint European Conference on Machine Learning and Knowledge
 Discovery in Databases*. Springer, 209–225.
- 947 [31] Sandra McCormick. 1992. Disabled readers' erroneous responses to inferential comprehension questions: Description and analysis.
 948 *Reading Research Quarterly* (1992), 55–77.
- 949 [32] Aliaksei Miniukovich, Antonella De Angeli, Simone Sulpizio, and Paola Venuti. 2017. Design guidelines for web readability. In *Proceedings
 950 of the 2017 Conference on Designing Interactive Systems*. 285–296.
- 951 [33] Aliaksei Miniukovich, Michele Scaltritti, Simone Sulpizio, and Antonella De Angeli. 2019. Guideline-based evaluation of web readability.
 952 In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–12.
- 953 [34] Gary E Raney, Spencer J Campbell, and Joanna C Bovee. 2014. Using eye movements to evaluate the cognitive processes involved in
 954 text comprehension. *Journal of visualized experiments: JoVE* 83 (2014).
- 955 [35] Keith Rayner. 1998. Eye movements in reading and information processing: 20 years of research. *Psychological bulletin* 124, 3 (1998),
 956 372.
- 957 [36] Keith Rayner, Alexander Pollatsek, Jane Ashby, and Charles Clifton Jr. 2012. Psychology of reading. (2012).
- 958 [37] Luz Rello and Mari-Carmen Marcos. 2012. An eye tracking study on text customization for user performance and preference. In *2012
 959 Eighth Latin American Web Congress*. IEEE, 64–70.
- 960 [38] Luz Rello, Martin Pielot, and Mari-Carmen Marcos. 2016. Make it big! The effect of font size and line spacing on online readability. In
 961 *Proceedings of the 2016 CHI conference on Human Factors in Computing Systems*. 3637–3648.
- 962 [39] Judith C Roberts and Keith A Roberts. 2008. Deep reading, cost/benefit, and the construction of meaning: Enhancing reading compre-
 963 hension and deep learning in sociology courses. *Teaching Sociology* 36, 2 (2008), 125–140.
- 964 [40] Ellen Rose. 2011. The phenomenology of on-screen reading: University students' lived experience of digitised text. *British Journal of
 965 Educational Technology* 42, 3 (2011), 515–526.
- 966 [41] Ladislao Salmerón, Johannes Naumann, Victoria García, and Inmaculada Fajardo. 2017. Scanning and deep processing of information in
 967 hypertext: an eye tracking and cued retrospective think-aloud study. *Journal of Computer Assisted Learning* 33, 3 (2017), 222–233.
- 968 [42] VE Snider. 1988. The role of prior knowledge in reading comprehension: A test with LD adolescents. *Direct Instruction News* 611 (1988).
- 969 [43] Catherine Snow. 2002. *Reading for understanding: Toward an R&D program in reading comprehension*. Rand Corporation.
- 970 [44] Namrata Srivastava, Joshua Newn, and Eduardo Veloso. 2018. Combining low and mid-level gaze features for desktop activity
 971 recognition. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 2, 4 (2018), 1–27.
- 972 [45] Alexander Strukelj and Diederick C Niehorster. 2018. One page of text: Eye movements during regular and thorough reading, skimming,
 973 and spell checking. *Journal of Eye Movement Research* 11, 1 (2018).
- 974 [46] Geoffrey Underwood, Alison Hubbard, and Howard Wilkinson. 1990. Eye fixations predict reading comprehension: The relationships
 975 between reading skill, reading speed, and visual inspection. *Language and speech* 33, 1 (1990), 69–81.
- 976 [47] Boris M. Velichkovsky and John Paulin Hansen. 1996. *New Technological Windows into Mind: There is More in Eyes and Brains for Human-
 977 Computer Interaction*. Association for Computing Machinery, New York, NY, USA, 496–503. <https://doi.org/10.1145/238386.238619>
- 978 [48] Titus von der Malsburg. 2015. Saccades: An R package for detecting fixations in raw eye tracking data.
- 979 [49] Maryanne Wolf. 2018. *Reader, come home: The reading brain in a digital world*. Harper New York, NY.
- 980 [50] Maryanne Wolf, Mirit Barzillai, and John Dunne. 2009. The importance of deep reading. *Challenging the whole child: reflections on best
 981 practices in learning, teaching, and leadership* 130 (2009), 21.
- 982
- 983
- 984
- 985
- 986
- 987