

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20

The jingle fallacy in comprehension tests for reading

Charlotte E. Lee*, Hayward J. Godwin, and Denis Drieghe

School of Psychology, University of Southampton, United Kingdom.

*Corresponding author

E-mail: c.lee@soton.ac.uk

Open Practice Statement: All data and materials are available online at:

https://osf.io/dk82u/?view_only=88d54bdc1f9b4793ba4b2dfef106085

This experiment was not pre-registered.

21 **Abstract**

22 The *Jingle fallacy* is the false assumption that instruments which share the same name
23 measure the same underlying construct. In this experiment, we focus on the comprehension
24 subtests of the Nelson Denny Reading Test (NDRT) and the Wechsler Individual
25 Achievement Test (WIAT-II). 91 university students read passages for comprehension whilst
26 their eye movements were recorded. Participants took part in two experimental blocks of
27 which the order was counterbalanced, one with higher comprehension demands and one with
28 lower comprehension demands. We assumed that tests measuring comprehension would be
29 able to predict differences observed in eye movement patterns as a function of varying
30 comprehension demands. Overall, readers were able to adapt their reading strategy to read
31 more slowly, making more and longer fixations, coupled with shorter saccades when
32 comprehension demands were higher. Within an experimental block, high scorers on the
33 NDRT were able to consistently increase their pace of reading over time for both higher and
34 lower comprehension demands, whereas low scorers approached a threshold where they
35 could not continue to increase their reading speed or further reduce the number of fixations to
36 read a text, even when comprehension demands were low. Individual differences based on
37 the WIAT-II did not explain similar patterns. The NDRT comprehension test was therefore
38 more predictive of differences in the reading patterns of skilled adult readers in response to
39 comprehension demands than the WIAT-II (which also suffered from low reliability). Our
40 results revealed that these different comprehension measures should not be used
41 interchangeably, and researchers should be cautious when choosing reading comprehension
42 tests for research.

43 Keywords: Jingle fallacy, Comprehension Demands, Individual Differences

44 **Introduction**

45 Reading comprehension is a complex task made up of interactions between the features
46 of a text and the skill and strategies of the reader [1,3]. The Simple View of Reading [4]
47 describes the basic requirements for reading as the ability to decode and identify words in
48 text by converting graphemes into phonemes combined with the ability to understand
49 information presented orally (language comprehension). However, in the more complex
50 Construction-Integration (CI) model of reading comprehension [1], text is represented by a
51 surface structure (semantic representations of words within a text), a textbase (a
52 representation of the explicit meaning of the whole text, coherently integrating each word
53 meaning) [5], and a situation model (where a reader creates a model of the situation,
54 integrating the explicit meaning of the text with their own world knowledge). For shallow
55 comprehension, a textbase is sufficient, however for deeper understanding a situation model
56 is required. Differences in theoretical conceptualisation of comprehension can result in
57 differences in the underlying mechanisms measured by comprehension tests based upon
58 them. Indeed, inconsistencies in research where skills measured by cognitive tasks are used to
59 predict readers' performance on reading comprehension measures have been suggested to
60 reflect differences in underlying cognitive mechanisms [6-9]. The current paper strives to
61 shed some light on the problems that researchers may face when selecting reading
62 comprehension tests, and the direct impact that test selection can have on conclusions based
63 upon them in eye tracking investigations.

64 **Evidence for a jingle fallacy**

65 In some of our previous eye movement investigations of average-to-very-skilled
66 readers [10,11] we found that two often-used reading comprehension subtests from

67 standardised reading ability measures failed to load together in a principal components
68 analysis and were only weakly correlated ($r = 0.21$, [10]; $r = 0.15$, [11]). These subtests were
69 from the Wechsler Individual Achievement Test (WIAT-II UK [12]) and the Nelson Denny
70 Reading Test (NDRT [13]). We concluded that these comprehension tests might be assessing
71 different underlying skills. Since these measures are both named ‘reading comprehension’,
72 this would present a clear example of Thorndike’s *Jingle fallacy*: that is, the misleading
73 assumption that two measures assess a single underlying construct because they share the
74 same name [14]. Although not uncommon in psychological research, where a variety of tests
75 are available to assess common constructs, problems when selecting and reporting
76 appropriate measures can lead to questionable research practices when used for scientific
77 purposes [15]. The aim of the current paper is to extend our previous investigations to
78 directly test the differences between the two tests by using them to predict differences in eye
79 movement patterns reflective of different comprehension demands and to further highlight
80 the pitfalls of comparing research that uses either test for this area of research.

81 **Differences in test format**

82 We start by discussing some qualitative differences in the format of the two
83 comprehension tests that may provide some insight into the underlying constructs that are
84 being tapped into by each one. First, the NDRT exclusively features non-fiction passages
85 whereas the WIAT-II features more varied text formats, with some fiction and non-fiction
86 passages as well as single sentences. Previous studies have noted that differences in the
87 format of reading materials (sentences vs paragraphs [16], fiction vs non-fiction [e.g., 17-19])
88 can impact reading behaviour as reflected in eye movement measures. Reading times are
89 longer and rereading is more common for sentences presented within paragraphs than for
90 sentences presented alone, which suggests that text format influences the reading strategy

91 used to comprehend the text [16]. Best et al. [17] also showed that comprehension accuracy
92 was higher for narrative texts than expository texts (non-fiction/scientific) and performance
93 on each was predicted by different individual skills. Decoding skills were a key element for
94 successful narrative text comprehension, whereas world knowledge was more important for
95 successful expository text comprehension. While this may suggest that narrative texts
96 included in the WIAT-II where comprehension is suggested to be higher, might be ‘easier’
97 for skilled readers, it also suggests that a reliance on non-fiction passages in the NDRT may
98 result in greater overlap with general knowledge. This was also suggested by Ready et al.
99 [20] following work by Coleman et al, [21] who found that college students could answer the
100 questions on NDRT comprehension tests and achieve a greater-than-chance level of accuracy
101 without actually reading the associated passages. However, we note that accuracy was 44 –
102 47 % whereas chance level was 20 % so the test is clearly measuring more than just general
103 knowledge.

104 Both tests also feature explicit differences in reading instruction since the WIAT-II
105 includes a combination of silent and oral reading, whereas the NDRT only features silent
106 reading. Reading aloud involves articulating the text as well as the standard process of
107 reading, and evidence from the eye-voice span (the distance between the location of a
108 fixation and the articulated word) demonstrates that oral reading involves additional working
109 memory processes [22]. Hale et al. [23] investigated differences in reading aloud and silently
110 and found that for children across grades 4-12 reading comprehension was higher when
111 reading aloud than when reading silently. In addition, some prior research suggests that
112 changing oral and silent reading tasks in comprehension tests may lead to different outcomes,
113 though this has been noted specifically in relation to differences between children with
114 reading difficulties and average readers [24]. Much less is known about how comprehension
115 changes when adults read aloud. A survey by Duncan and Freeman [25] reported that, of 529

116 respondents, 67.5 % said that they read aloud to understand difficult text (though they noted
117 that this was usually only brief). Gambrell and Heathington [26] reported that 36 % of poor
118 adult readers said that they could read more quickly when reading aloud compared to just 4 %
119 of good readers. It may be that an oral reading task to assess comprehension is less
120 informative for adult readers due to individual differences, though more research is needed to
121 investigate this.

122 Another notable difference is that testing in the WIAT-II is administered by an
123 experimenter who asks questions aloud to the participant and records their spoken responses
124 on paper, whereas the NDRT is administered independently. This procedural difference could
125 lead to performance anxiety for participants when taking the WIAT-II and may introduce
126 noise into data collected under these conditions. This may be especially important where
127 participants are sometimes asked to read aloud. In contrast it may mean that the NDRT has
128 comparatively less control to determine whether a participant is properly engaging with the
129 task. This aspect highlights another qualitative difference in the administration of the WIAT-
130 II in comparison to the NDRT.

131 A good comprehension test should be able to predict differences in behaviour between
132 tasks that vary in comprehension demands. Eye movement measures reflect complex
133 cognitive processes active during reading [27,28]. The current paper therefore investigated
134 global reading strategies for paragraph reading and aims to examine whether the differences
135 we described between the comprehension subtests of the WIAT-II and the NDRT impact
136 their ability to predict eye movement patterns reflecting changes in comprehension demands.

137 **Individual differences in adult readers' eye movements**

138 We turn our focus now to individual differences in adult readers' eye movements.
139 Skilled adult readers typically read more quickly, make fewer and shorter fixations, longer
140 saccades and fewer regressions than less skilled readers [28,29]. However, reading skill is not
141 directly related to reading rate, and a speed-accuracy trade-off means that faster reading
142 eventually leads to lower levels of comprehension [30]. There is much variability within
143 groups of skilled adult readers with fixation durations varying between approximately 50-600
144 ms and saccade lengths between 1-20 letter spaces [31,32]. Skilled readers also vary in how
145 they respond to features of a text. Ashby et al. [29] found that poor adult readers (identified
146 using NDRT reading comprehension and vocabulary tests) benefitted more from highly
147 constraining sentential contexts compared to skilled readers. Similarly, Bisanz et al. [33]
148 reported a complex relationship between reading ability and reading times in line with
149 Stanovich's [34] interactive-compensatory model which stated that poor readers, who had
150 below average bottom-up processing skills, would rely more heavily on contextual cues when
151 they were available. Bisanz et al. [33] showed that poor readers actually read some sentences
152 more quickly than skilled readers. It has been suggested that some readers might use a 'risky'
153 reading strategy where they read more quickly and make fewer refixations than other readers
154 [35].

155 **Intra-individual differences**

156 In addition to the differences observed between readers, intra-individual differences
157 (variability within the same reader) can also influence reading behaviours. It has been well
158 established that task demands can influence the way that readers process a text: skilled
159 readers are able to adjust their reading behaviours (and pace) to the demands of the task [36]

160 and are able to read thoroughly or superficially when needed [37]. Aaronson and Ferres
161 [38,39] noted that skilled readers are more likely to use a ‘recall strategy’ focussed on
162 structural aspects of a text when a reading task involves direct recall of words/sentences, but
163 when the task involves true/false questions, their focus is driven by the meaning of the text
164 using a ‘comprehension strategy’. This research was influential as it gave clear evidence that
165 skilled readers had some autonomy over how deeply they processed a text.

166 It has been noted that when texts are more difficult, a more ‘careful’ strategy might be
167 used where, in comparison to a risky reading strategy [35], readers tend to make more
168 refixations, have smaller average fixation durations and smaller saccade amplitudes [40].
169 Researchers have investigated whether these strategies can be observed for identical
170 sentences when different comprehension demands are placed upon them. Radach et al. [16]
171 investigated differences in eye movement behaviours related to the specific reading task as
172 well as different text formats. Participants took part in one of two tasks: comprehension,
173 where participants were asked detailed questions about the text; and a word verification task
174 where participants had to indicate which word had appeared in the sentence from some given
175 options. Radach et al. [16] also compared eye movement measures within these groups for
176 identical sentences that were either embedded within a passage or were presented alone.
177 Researchers concluded that top-down processes influenced by the task (comprehension vs
178 word identification) and format of the text (sentences vs paragraphs) clearly impacted the eye
179 movement record. Word-viewing times were significantly longer on comprehension tasks and
180 more fixations were made on a word in this task than in the verification task, indicating more
181 careful reading when reading for comprehension. Passages were read more quickly on the
182 first-pass but featured more rereading than sentences.

183 Similarly, Wotschack and Kliegl [41] investigated the effect of easy ‘verification’
184 questions (after 27 % of sentences) compared to ‘hard’ comprehension questions about
185 sentence meaning (following 100 % of sentences) and found that the more difficult questions
186 were associated with more careful reading as indicated by more rereading and more
187 regressions. However, they found that accuracy was high in both conditions and questioned
188 the strength of their manipulation. In response, Weiss et al. [42] aimed for a stronger
189 difficulty manipulation and investigated ‘easy’ lexical verification questions versus ‘difficult’
190 comprehension questions that required resolving some syntactic ambiguity. For example, a
191 sentence containing a subjective relative clause such as ‘The chef that distracted the waiter sifted
192 the flour onto the counter’, was followed by an easy question: ‘Did a chef do something?’ Or a
193 difficult question: ‘Did the waiter distract the chef?’ They did see differences in accuracy
194 between the difficult (83 %) and easy (97 %) conditions, and also found that participants
195 made more regressions and spent more time rereading texts in the difficult condition but that
196 no disruptions were seen in first pass fixation times. Weiss et al. [42] concluded that inflated
197 differences happened at the end of passages even when the ambiguity occurred earlier in the
198 sentence. Accuracy was not predicted by the magnitude of the disruption, suggesting that the
199 increased processing time was a ‘checking mechanism’ rather than additional information
200 processing.

201 Christianson et al. [43] reached a conclusion similar to Weiss et al. [42] in a study that
202 investigated rereading behaviours in garden-path sentences (where an ambiguity in the
203 sentence meaning is revealed fairly late in the sentence e.g. The babysitter who was
204 purchased a gift card thanked the parents) vs. local coherence structures (where ambiguities
205 were resolved earlier, e.g. The parents thanked the babysitter who was purchased a gift card).
206 They found that rereading behaviours were more consistent with confirmatory rereading
207 (checking) than revisionary rereading (for understanding) because rereading was not

208 consistently predicted by critical regions in the sentence structure, and rereading behaviours
209 were not predictors of offline comprehension accuracy.

210 Recent investigations have looked more closely at rereading behaviours and have
211 started to examine individual differences in rereading. A study by Andrews and Veldre [44]
212 investigated ‘wrap-up’ effects in tasks with different comprehension loads in relation to
213 individual differences in reading proficiency (measured by vocabulary, reading
214 comprehension reading rate (NDRT [13]), spelling dictation and spelling recognition [45]).
215 Wrap-up effects [46] are where longer reading times are observed at clause and sentence
216 boundaries, where readers integrate information before moving forward in a text [47,48].
217 Wrap-up times have been associated with the goals of the reading task, for example in a study
218 by Stine-Morrow et al. [49] where differences in wrap-up predicted recall but not
219 comprehension success. Importantly, Andrews and Veldre assessed readers’ individual
220 differences in spelling, reading comprehension (NDRT), vocabulary and reading rate
221 alongside manipulating how often comprehension questions occurred (after 25 % of passages
222 or 100 % of passages). They found that comprehension load had little effect on wrap-up,
223 however it did lead to shallower (more risky) reading strategies when comprehension
224 demands were low, with longer passage reading times, more refixations and regressions, but
225 no differences in average fixation times or forward saccade lengths. Andrews and Veldre [44]
226 found that the better readers (as identified via a composite score of the individual differences
227 measures that have been shown to provide a good assessment of lexical quality [50-54])
228 generally read passages more quickly, made fewer and shorter fixations, longer forward
229 saccades and marginally fewer regressions than poorer readers. They did not find that reading
230 proficiency composite scores interacted with the effect of comprehension load on eye
231 movement measures, but they noted that accurate comprehension was associated with more

232 consistent reading behaviour, where readers did not adjust their reading strategy much in
233 response to comprehension load.

234 Reading strategies may of course be adapted over time during an experiment. For
235 example, readers may read through early trials more slowly when they have higher
236 comprehension demands, until they are familiar with the format of the questions in the
237 experimental block, after which they may adjust their reading rate to speed up processing
238 time. This rate of adaptation may be modulated by individual differences, whereby better
239 comprehenders might be able to increase their reading rate to one that is optimal/preferred
240 more quickly over trials than less skilled comprehenders. Therefore, besides examining the
241 differences between predictions based on two comprehension tests, a second goal of current
242 study was to determine whether individuals alter their reading strategies in response to
243 comprehension demands gradually as trials progress. We were interested to see if individual
244 differences in reading ability predicted differences in the rate of adaptation to different
245 comprehension demands as well as whether discrepancies occurred between the two
246 measures of reading comprehension that we included. Following Radach et al. [16], identical
247 reading materials were used between conditions in the current study to directly compare the
248 influence of comprehension demands placed on the reader via differences in the difficulty of
249 questions that followed.

250 **Predictions**

251 We expected high scores on the comprehension tests to predict faster passage reading
252 times as faster sentence reading times were associated with higher scores on the
253 comprehension subtests from the WIAT-II [12] and the NDRT [13] in Lee, Godwin et al.
254 [10] and Lee, Pagán et al. [11]. Note however that the format of our experimental materials in
255 the current study (paragraphs) was different in comparison to our previous investigations

256 (sentences). Longer reading times and more rereading have been observed for passages
257 compared to sentences [16]. Similarly, since comprehension was included as part of the
258 composite measure of reading proficiency in Andrews and Veldre [44], who found that higher
259 reading proficiency predicted faster passage reading times, shorter average fixation durations,
260 longer forward saccades and a greater number of regressions than low proficiency, we
261 expected similar patterns to emerge for our comprehension scores.

262 We expected that higher offline comprehension scores would predict faster passage
263 reading times, shorter average fixation durations, longer forward saccades and fewer
264 regressions. Higher comprehension demands were expected to increase the number of
265 fixations and the time that participants spent reading the passages. We anticipated that all
266 readers would adapt their reading strategy to become more efficient (they would make fewer
267 fixations, longer saccades, shorter fixations and read passages more quickly), but that there
268 might be individual differences observed in the rate of adaptation or ceiling levels in saccade
269 lengths that poorer readers could reach, since poorer readers have been shown to have shorter
270 rightwards perceptual spans (in languages read left to right) than better readers [55].
271 Similarly, as poorer readers usually exhibit slower reading times and longer fixations than
272 skilled readers [28,29,44,10,11] we anticipated floor effects for poor readers' minimum
273 passage reading times, fixation durations and the number of fixations. Since the intended
274 population was skilled adult readers, it was likely that accuracy would be high across tasks
275 (as was observed in [44,41]). Therefore, because comprehension accuracy is often higher for
276 narratives than expository texts [17], expository passages were used in the current study to
277 maximise the likelihood of variability in accuracy scores.

278 We note that the NDRT exclusively uses expository texts to measure comprehension,
279 therefore it may be more similar in format to the passages used in this experiment. As noted

280 by Ready et al. [20] and Coleman et al. [21], the NDRT may also feature a high degree of
281 overlap with general knowledge or world knowledge, which has been found to be associated
282 with expository text comprehension. Therefore, it would not be surprising if the NDRT
283 predicts higher comprehension accuracy across conditions, than the WIAT-II, which features
284 more varied reading formats. We also anticipated that the WIAT-II may be more noisy in its
285 predictions due to some performance anxiety induced by the experimenter's presence.

286 **Method**

287 **Participants**

288 Participants were 91 students and staff from the University of Southampton over the
289 age of 18 (11 Males, $M = 20.27$ years, range = 18 – 45 years). An additional 9 participants
290 took part in the study, but their data were removed from the final dataset due to poor overall
291 accuracy on the comprehension questions in the eye tracking task (below 60 % where chance
292 level was 50 %). Participants were all native English speakers with normal or corrected to
293 normal vision and no known reading difficulties. Participants received course credits or £25
294 for completing the study. Recruitment took place from 29/10/2021 to 10/06/2022. This study
295 was approved by the University of Southampton Ethics and Research Governance Board.

296 **Apparatus**

297 Paragraphs and questions were presented on a 21-inch CRT monitor, with a refresh rate
298 of 120 Hz and a resolution of 1024 x 768 at a viewing distance of 60 cm. Passages were
299 presented in Courier New, size 14 font on a grey background; three characters equated to
300 about 1° of visual angle. Although reading was binocular, eye movements were recorded
301 from the right eye only using an EyeLink 1000 tracker [56]. Forehead and chin rests were

302 used to minimize head movements. The spatial resolution of the eye tracker was 0.05°, and
303 the sampling rate was 1000 hz.

304 Participants used a 14-inch Dell Laptop Computer to complete the NDRT
305 comprehension test administered using an online web browser running Qualtrics. For
306 copyright issues, whenever we ran a participant using the online version, we voided a
307 purchased paper version. Participants were required to select answers using a mouse. During
308 WIAT-II comprehension test researchers used the testing flip pad and scoring sheets included
309 in the test pack.

310 **Materials**

311 Forty experimental paragraphs ($M = 138.33$ words, $SD = 19.28$) were adapted from
312 freely available online practice comprehension tests [57]. Two conditions were created for
313 each paragraph, one with lower comprehension demands where participants were asked
314 ‘What is the passage about?’ and were given two short options that consisted of a word or
315 phrase (e.g., Synaesthesia/Claustrophobia). One option was directly related to the passage
316 and the other was unrelated. In the higher comprehension demands condition participants
317 were asked, ‘What is the main idea of the passage?’ and two longer and more detailed options
318 were presented from which participants were asked to select an answer (e.g., People with
319 synaesthesia experience a fusing of different senses/People with synaesthesia may hear a
320 sound when they touch an object). In this condition, both answers were related to the passage,
321 but one provided a better evaluation of the passage meaning. Questions were similarly
322 phrased but differences were presented by the type of options available, and level of detail
323 needed to select a correct answer. The original questions from the online practice materials
324 were the ‘higher demands’ questions, a ‘lower demands’ alternative was then created for each

325 of them. Paragraph naturalness and comprehension question difficulties were independently
326 rated by participants who did not take part in subsequent testing. Passages were rated on a
327 scale from 0 (very unnatural) to 100 (very natural) ($M = 63.04$, $SD = 5.31$) to ensure there
328 were no outliers in the readability of the text and questions were rated on a scale from 0 to 100
329 as more difficult ($M = 23.71$, $SD = 4.45$) than low comprehension demand questions ($M = 19.97$, SD
330 $= 3.67$), $t(49) = -8.57$, $p < .001$. Two counterbalanced lists were then created so that each
331 participant viewed 20 of each question type but did not view the same paragraph twice. The
332 paragraphs occupied 10 - 13 lines on the screen ($M = 859.95$ characters including spaces,
333 $Max = 1159$ characters).

334 **Design and procedure**

335 Testing took place over two sessions with a minimum of two days in between them.
336 During the first session participants were given an information sheet and were asked to sign a
337 consent form and completed two eye tracking tasks (the first eye tracking task was for a
338 separate experiment, where participants read 60 single sentences and lasted approximately 30
339 minutes), followed by the experimenter administered WIAT-II comprehension test and some
340 other cognitive tasks belonging to an unrelated experiment (Rapid Automated Naming and
341 the pseudoword decoding and word reading subtests of the WIAT-II. These tasks took
342 approximately 15 minutes to complete). The same experimenter administered this task to all
343 participants to control for as much experimental variation between participants as possible. A
344 script was read from the test materials to ensure that instructions were identical for all
345 participants. Participants read passages (short narratives and information texts) aloud or
346 silently and were asked literal and inferential comprehension questions by the experimenter,
347 participants gave spoken responses which the experimenter transcribed.

348 For the eye tracking task, participants were asked to sit comfortably at the computer,
349 resting their chin on a chinrest and were then guided through the set up and 9-point
350 calibration of the eye tracker by the researcher. Participants were then required to direct their
351 gaze to a fixation cross presented in the upper left portion of the screen. Once participants
352 fixated upon the cross sentences were presented and always began at the location marked by
353 the fixation cross. Participants were asked to read the paragraphs and answer questions
354 presented on the screen using the keyboard to respond. Participants either answered questions
355 with longer, and more detailed options from which to select an answer (higher
356 comprehension demands) or with shorter, simpler options (lower comprehension demands)
357 depending on the condition. The same participants completed both conditions over two
358 sessions. Eye tracking sessions took place on two separate days. At the start of the blocks,
359 participants read five practice paragraphs with questions matching the type for the current
360 condition. Practice questions were followed by 20 experimental paragraphs that were each
361 followed by a comprehension question. Block order was counterbalanced so that participants
362 who read paragraphs and answered questions with higher comprehension demands in session
363 1, then read paragraphs and answered questions with lower comprehension demands in
364 session 2 and vice versa. Block order was randomly assigned and within each block trial
365 order was randomised. Participants could take breaks when needed.

366 During the second session participants took part in the second part of the eye tracking
367 task (featuring the condition that they had not yet completed). Participants were then asked to
368 complete the NDRT comprehension test, and some other online tasks for a separate study (the
369 vocabulary subtest of the NDRT, a test of vocabulary knowledge, spelling dictation and
370 spelling recognition tasks, an Author Recognition test, and a backwards digit span task in a
371 randomised order. These tasks took approximately 40 minutes to complete). During the
372 NDRT participants silently read up to 7 passages and answered 5 - 8 MCQ questions about

373 them with 4 available options. Questions appeared below the passages on the same screen.
374 On the first passage participants were stopped after 1 minute and were asked to record the
375 number corresponding to the line that they had been reading to measure their reading rate.
376 Testing automatically stopped after 10 minutes and answers were recorded. Testing for the
377 NDRT comprehension test followed a half-timed procedure, in which the standard time limit
378 for completing this test was reduced by half. Participants were not aware of this reduced time
379 limit. This procedure has been shown to generate a more normal distribution for university
380 student readers (like those who took part in the current study) than the standard time limits in
381 an investigation by Andrews et al. [50]. To measure test reliability, Cronbach's alpha was
382 calculated for both the WIAT-II comprehension test ($\alpha = .62$) and the NDRT comprehension
383 test ($\alpha = .75$). We note that even though these estimates are still considered acceptable, they
384 are lower than those reported for normed data (NDRT = .89 to .98 [12] and WIAT = .98 [13]
385). Given that we focus on a university sample of readers, this might suggest these
386 comprehension tests are less reliable for this population of readers. Furthermore, we note that
387 by using a timed version of the NDRT comprehension test, reliability will be lower for
388 instances where participants answered fewer questions in the given time.

389 **Results**

390 Overall accuracy on comprehension questions was high but not at ceiling level ($M =$
391 79.42% , $SD = 10.53$). Reading comprehension scores on WIAT-II were calculated and
392 normed following guidance from the experimenter manual ($M = 110.47$, $SD = 9.37$, range =
393 $71 - 124$). NDRT comprehension scores were calculated based on raw scores due to the half-
394 timed aspect of the task. NDRT comprehension scores ($M = 57.64$, $SD = 11.02$, range = $20 -$
395 74) were weakly correlated with the WIAT-II comprehension scores ($r = 0.22$, $p = .039$).
396 Both WIAT-II Comprehension and NDRT Comprehension were weakly correlated with

397 overall accuracy on the eye tracking task (WIAT-II $r = .22$, $p < .001$; NDRT $r = 0.11$, $p <$
398 $.001$). Scores on both tests were standardised for further analyses.

399 **Data cleaning**

400 Eye tracking trials identified by the experimenter as having issues with tracker loss or
401 featuring excessive blinking were removed prior to the analysis. Fixations shorter than 80 ms
402 that landed within one character of the previous or next fixation were merged. Then, of the
403 remaining fixations, those shorter than 80 ms and longer than 800 ms were removed. Practice
404 trials were also removed. Due to an error in the programming of the experiment, texts were
405 presented with a justified alignment which meant that word level data would have confounds
406 between word length and visual extent. For this reason, word level measures such as
407 regressions and refixations were not included in these analyses.

408 The following global eye tracking measures were calculated for each trial; Number of
409 Fixations (total number of fixations made on a trial); Average Fixation Duration (mean
410 duration in ms of all fixations in a trial); Forward Saccade Length (the distance in degrees of
411 visual angle between one fixation and the next); and Total Passage Reading Time (total time
412 in ms spent reading the passage in a trial). Trials where total passage reading times fell
413 outside of 2.5 standard deviations from the mean for each participant were removed as
414 outliers (1.31 % of data removed). Data were then removed for each eye movement measure
415 per participant that fell outside of 2.5 standard deviations from the mean (Number of
416 Fixations (0.59 % data removed); Average Fixation Duration (0.88 % data removed);
417 Forward Saccade Length (1.16 % data removed). Descriptive statistics per condition for these
418 measures were calculated across participants and are displayed in Table 1.

419 **Table 1. Descriptive Statistics for Eye Movement Measures**

	Condition	Min	Max	Mean	SD
Number of Fixations	Low	42.00	285.00	138.38	33.43
	High	65.00	269.00	143.64	33.95
Average Fixation Duration (ms)	Low	138.53	285.94	206.14	24.06
	High	144.30	279.17	206.79	23.80
Average Forward Saccade Length (visual degrees)	Low	3.38	9.66	6.07	0.98
	High	3.41	9.25	5.99	0.99
Total Passage Reading Time (ms)	Low	8388.00	58115.00	28854.00	8484.70
	High	11232.00	63869.00	30138.17	8868.98

420 Descriptive statistics are based on participant means per condition.

421 **Linear mixed models**

422 Eye movement measures were analysed using the lme4 package (version 1.1-31 [58]) in
 423 R (version 4.2.2 [59]). Data were checked for normality and were not transformed for
 424 modelling as their distribution closely resembled a normal distribution. Binomial Generalized
 425 Linear Mixed Models were used to model accuracy data. The following model building
 426 strategy was followed. Models featured all fixed effects of interest: the main effect of
 427 experimental condition (lower vs higher comprehension demands), either the NDRT or
 428 WIAT-II comprehension test scores and the trial number and all the interactions. To ensure
 429 the maximal model was achieved, we started with a full random structure (all random slopes
 430 were included for subjects and items) and performed stepwise trimming of this structure until
 431 the model converged [60]. Slopes were first trimmed from the random effects structure where
 432 perfect correlations were indicated and subsequently factors that explained the smallest
 433 amount of variance until the model converged.

434 **Number of fixations**

435 Models shown in Table 2 and 3 indicated that overall, more fixations were made on
436 paragraphs where comprehension demands of the questions were high compared to when
437 they were low. The number of fixations decreased slightly over trials, however, a significant
438 three-way interaction between trial number, condition and scores on the NDRT
439 comprehension test revealed a more complex pattern based on individual differences (Table
440 2). Fig 1 shows that high scorers on the NDRT comprehension test reduced the number of
441 fixations further into the experimental session (analyses were based on continuous
442 comprehension scores but are presented in 3 panels for the mean +/- 1SD in figures to clearly
443 demonstrate the 3-way interaction). They also made more fixations in the difficult condition
444 and these two factors did not interact. A different pattern emerged for the low scorers on the
445 NDRT comprehension test. Low scorers did make more fixations on a paragraph at the
446 beginning of the experiment than on trials nearing the end of the experiment, but when
447 comprehension demands were low, this decrease was not as steep. This pattern may indicate
448 that less skilled comprehenders were nearing floor effects where they were close to the
449 minimum number of fixations that they could accommodate whilst still reading for
450 comprehension when comprehension demands were low.

451 No significant interactions or individual differences were observed for scores on the
452 WIAT-II comprehension test, and in this model a trial by condition interaction was
453 marginally significant (Table 3).

454 **Table 2. LMM for Number of Fixations predicted by NDRT Comprehension and**
 455 **Interactions with Trial Number and Condition**

	β	95 % CI	t	df	p
Intercept	153.47	[147.01, 159.92]	46.60	140.91	< .001 ***
Trial Number	-0.75	[-0.85, -0.64]	-14.38	3300.23	< .001 ***
Condition	9.64	[5.20, 14.07]	4.26	335.04	< .001 ***
NDRT Comprehension	-1.10	[-6.65, 4.46]	-0.39	104.33	.699
Trial Number \times Condition	-0.22	[-0.42, -0.01]	-2.09	3302.67	.037 *
Trial Number \times NDRT Comprehension	-0.16	[-0.26, -0.06]	-3.13	3299.89	.002 **
Condition \times NDRT Comprehension	-4.23	[-8.60, 0.14]	-1.90	346.14	.059.
Trial Number \times Condition \times NDRT Comprehension	0.25	[0.05, 0.45]	2.41	3300.46	.016 *

456 The baseline of the condition term is lower comprehension demands. Estimates represent the
 457 change when going from lower to higher comprehension demands.

458 **Table 3. LMM for Number of Fixations predicted by WIAT-II Comprehension and**
 459 **Interactions with Trial Number and Condition**

	β	95 % CI	t	df	p
Intercept	153.17	[146.69, 159.64]	46.35	140.84	< .001 ***
Trial Number	-0.77	[-0.87, -0.67]	-14.77	3302.63	< .001 ***
Condition	9.22	[4.78, 13.65]	4.07	351.96	< .001 ***
WIAT-II Comprehension	1.90	[-4.39, 8.20]	0.59	104.53	.555
Trial Number \times Condition	-0.18	[-0.39, 0.02]	-1.77	3303.08	.076.
Trial Number \times WIAT-II Comprehension	0.07	[-0.04, 0.18]	1.21	3300.95	.226
Condition \times WIAT-II Comprehension	-0.18	[-5.14, 4.78]	-0.07	351.61	.943
Trial Number \times Condition \times WIAT-II Comprehension	-0.06	[-0.29, 0.17]	-0.51	3300.98	.611

460 The baseline of the condition term is lower comprehension demands. Estimates represent the
 461 change when going from lower to higher comprehension demands.

462 **Fig 1. A Three-Way Interaction between NDRT Comprehension Scores, Condition**
 463 **(higher vs lower comprehension demands) and Trial Number on the Number of**
 464 **Fixations made when Reading a Paragraph.** Shaded areas represent 95 % confidence
 465 intervals.

466 -FIG1 HERE-

467 **Average fixation duration**

468 Tables 4 and 5 present models for average fixation durations. These models indicated
 469 that overall, average fixation durations increased slightly over trials. A three-way interaction
 470 between trial number, condition and scores on the NDRT comprehension test was observed
 471 (Table 4). Fig 2 shows that for low scorers on the NDRT average fixation durations increased

472 from early to late trials in the experiment. For the high scorers a different pattern emerges
 473 depending on the comprehension demands with average fixation time going up when
 474 comprehension demands are high and going down when they are low.

475 WIAT-II comprehension scores were not significant predictors of average fixation
 476 durations (Table 5), though an interaction of trial by condition was marginally significant.

477 **Table 4. LMM for Average Fixation Durations predicted by NDRT Comprehension and**
 478 **Interactions with Trial Number and Condition**

	β	95 % CI	t	df	p
Intercept	205.48	[200.82, 210.14]	86.42	98.98	< .001 ***
Trial Number	0.10	[0.05, 0.15]	4.00	3300.34	< .001 ***
Condition	-0.41	[-2.71, 1.89]	-0.35	280.96	.729
NDRT Comprehension	-2.18	[-6.75, 2.38]	-0.94	94.13	.351
Trial Number \times Condition	0.07	[-0.03, 0.17]	1.43	3301.82	.154
Trial Number \times NDRT Comprehension	-0.07	[-0.12, -0.03]	-2.97	3299.18	.003 **
Condition \times NDRT Comprehension	-2.62	[-4.89, -0.34]	-2.26	278.34	.025 *
Trial Number \times Condition \times NDRT Comprehension	0.12	[0.02, 0.21]	2.31	3296.20	.021 *

479 The baseline of the condition term is lower comprehension demands. Estimates represent the
 480 change when going from lower to higher comprehension demands.

481 **Table 5. LMM for Average Fixation Durations predicted by WIAT-II Comprehension**
 482 **and Interactions with Trial Number and Condition**

	β	95 % CI	t	df	p
Intercept	205.58	[200.90, 210.26]	86.08	98.93	< .001 ***
Trial Number	0.09	[0.04, 0.14]	3.64	3297.93	< .001 ***
Condition	-0.65	[-2.96, 1.66]	-0.55	280.78	.580
WIAT-II Comprehension	-3.27	[-8.44, 1.90]	-1.24	94.18	.218
Trial Number \times Condition	0.09	[-0.01, 0.19]	1.75	3298.81	.080.
Trial Number \times WIAT-II Comprehension	0.01	[-0.04, 0.07]	0.38	3297.93	.703
Condition \times WIAT-II Comprehension	-0.02	[-2.60, 2.56]	-0.01	278.01	.989
Trial Number \times Condition \times WIAT-II Comprehension	-0.05	[-0.16, 0.06]	-0.87	3298.43	.383

483 The baseline of the condition term is lower comprehension demands. Estimates represent the
 484 change when going from lower to higher comprehension demands.

485 **Fig 2. A Three-Way Interaction between NDRT Comprehension Scores, Condition**
 486 **(higher vs lower comprehension demands) and Trial Number on Average Fixation**
 487 **Durations when Reading a Paragraph.** Shaded areas represent 95 % confidence intervals.

488 -FIG2 HERE-

489 **Average forward saccade length**

490 Models for average forward saccade lengths are displayed in Tables 6 and 7. In both
 491 models, longer forward saccades were observed for passages with lower comprehension
 492 demands than for identical passages with higher comprehension demands. A slight increase
 493 in forward saccade length over trials was also predicted by both models. Table 6 shows that a
 494 three-way interaction between trials, conditions and NDRT comprehension scores was

495 significant though numerically small. Differences can be seen in Fig 3 where those who
 496 scored highly on the NDRT made slightly longer forward saccades when comprehension
 497 demands were low compared to when comprehension demands were high, but in both
 498 comprehension demand conditions forward saccade lengths became longer further in the
 499 experiment. Low scorers also made longer forward saccades when comprehension demands
 500 were low than when they were high, however the average length of their forward saccades
 501 only increased over time when comprehension demands were high. When comprehension
 502 demands were low, these readers did not make longer forward saccades over trials, with
 503 comparable average forward saccade lengths across all trials.

504 No significant effects of individual differences in WIAT-II comprehension test scores
 505 were observed for average forward saccade lengths (Table 7).

506 **Table 6. LMM for Average Forward Saccade Length predicted by NDRT**
 507 **Comprehension and Interactions with Trial Number and Condition**

	β	95 % CI	t	df	p
Intercept	5.89	[5.70, 6.08]	60.79	104.80	< .001 ***
Trial Number	0.01	[0.01, 0.01]	7.37	3279.81	< .001 ***
Condition	-0.14	[-0.24, -0.04]	-2.72	340.04	.007 **
NDRT Comprehension	0.02	[-0.16, 0.21]	0.26	96.07	.797
Trial Number \times Condition	0.00	[0.00, 0.01]	1.68	3289.38	.092.
Trial Number \times NDRT Comprehension	0.01	[0.00, 0.01]	5.19	2762.05	< .001 ***
Condition \times NDRT Comprehension	0.10	[0.00, 0.20]	1.94	334.43	.054.
Trial Number \times Condition \times NDRT Comprehension	-0.01	[-0.01, 0.00]	-2.16	2633.99	.031 *

508 Note. The baseline of the condition term is lower comprehension demands. Estimates
 509 represent the change when going from lower to higher comprehension demands.

510 **Table 7. LMM for Average Forward Saccade Length predicted by WIAT-II**

511 **Comprehension and Interactions with Trial Number and Condition**

	β	95 % CI	t	df	p
Intercept	5.88	[5.69, 6.07]	60.24	104.60	< .001 ***
Trial Number	0.01	[0.01, 0.01]	7.99	3284.76	< .001 ***
Condition	-0.14	[-0.24, -0.04]	-2.65	351.51	.009 **
WIAT-II Comprehension	0.06	[-0.15, 0.27]	0.54	96.75	.591
Trial Number \times Condition	0.00	[0.00, 0.01]	1.41	3287.86	.159
Trial Number \times WIAT-II Comprehension	0.00	[0.00, 0.00]	-0.55	2314.64	.580
Condition \times WIAT-II Comprehension	0.01	[-0.10, 0.12]	0.19	353.67	.851
Trial Number \times Condition \times WIAT-II Comprehension	0.00	[0.00, 0.01]	1.46	2661.39	.144

512 Note. The baseline of the condition term is lower comprehension demands. Estimates

513 represent the change when going from lower to higher comprehension demands.

514

515 **Fig 3. A Three-Way Interaction between NDRT Comprehension Scores, Condition**

516 **(higher vs lower comprehension demands) and Trial Number on the Average Forward**

517 **Saccade Length when Reading a Paragraph.** Shaded areas represent 95 % confidence

518 intervals.

519 -FIG3 HERE-

520 **Total passage reading times**

521 Models for total passage reading times are presented in Tables 8 and 9. In both models,

522 passages in trials that occurred later in the experiment for both conditions were read more

523 quickly than earlier passages. Passages were also read more quickly when comprehension

524 demands were low compared to when comprehension demands were high (this was

525 significant in both models). The model presented in Table 8 also revealed that total passage
526 reading times were influenced by a three-way interaction between trials, conditions and
527 NDRT comprehension scores. Fig 4 shows that high scorers consistently read passages more
528 quickly towards the end of the experimental conditions than at the beginning, and read
529 passages with lower comprehension demands more quickly than passages with higher
530 comprehension demands. High scorers also read more quickly than low scorers by the end of
531 the experiment in both conditions.

532 Low scorers on the NDRT comprehension scores displayed a different pattern, where
533 their reading times were longer when comprehension demands were high compared to low at
534 the beginning of the experiment and decreased over trials. However, when comprehension
535 demands were low a potential floor effect was observed for low scorers where only a small
536 decrease in reading times across trials was seen for passages.

537 No other significant effects were observed in the model including the WAIT-II
538 comprehension scores (Table 9).

539 **Table 8. LMM for Total Passage Reading Times predicted by NDRT Comprehension**
 540 **and Interactions with Trial Number and Condition**

	β	95 % CI	t	df	p
Intercept	31996.48	[30331.38, 33661.57]	37.66	132.60	< .001 ***
Trial Number	-144.55	[-168.41, -120.70]	-11.88	3321.32	< .001 ***
Condition	2152.00	[1081.86, 3222.15]	3.94	304.94	< .001 ***
NDRT Comprehension	-727.71	[-2219.17, 763.74]	-0.96	100.48	.341
Trial Number \times Condition	-42.15	[-89.90, 5.59]	-1.73	3323.40	.084.
Trial Number \times NDRT Comprehension	-39.11	[-62.68, -15.55]	-3.25	3320.14	.001 **
Condition \times NDRT Comprehension	-1270.52	[-2324.79, -216.25]	-2.36	311.73	.019
Trial Number \times Condition \times NDRT Comprehension	66.07	[19.03, 113.10]	2.75	3321.15	.006 **

541 The baseline of the condition term is lower comprehension demands. Estimates represent the
 542 change when going from lower to higher comprehension demands.

543 **Table 9. LMM for Total Passage Reading Times predicted by WIAT-II Comprehension**
 544 **and Interactions with Trial Number and Condition**

	β	95 % CI	t	df	p
Intercept	31935.5 2	[30246.02, 33625.02]	37.05	131.55	< .001 ***
Trial Number	-150.52	[-174.47, -126.58]	-12.32	3321.70	< .001 ***
Condition	2059.34	[986.54, 3132.14]	3.76	308.20	< .001 ***
WIAT-II Comprehension	-70.65	[-1782.12, 1640.82]	-0.08	100.25	.936
Trial Number \times Condition	-34.52	[-82.43, 13.40]	-1.41	3322.80	.158
Trial Number \times WIAT-II Comprehension	18.14	[-8.47, 44.76]	1.34	3320.94	.182
Condition \times WIAT-II Comprehension	-291.17	[-1483.74, 901.41]	-0.48	314.56	.633
Trial Number \times Condition \times WIAT-II Comprehension	-9.23	[-62.48, 44.03]	-0.34	3321.66	.734

545 The baseline of the condition term is lower comprehension demands. Estimates represent the
 546 change when going from lower to higher comprehension demands.

547 **Fig 4. A Three-Way Interaction between NDRT Comprehension Scores, Condition**
 548 **(higher vs lower Comprehension demands) and Trial Number on Total Passage**
 549 **Reading Times.** Shaded areas represent 95 % confidence intervals.

550 -FIG4 HERE-

551
 552 **Accuracy**

553 Neither model showed significant differences in accuracy for high (M = 84 %, SD =
 554 10.69) compared to lower comprehension demands (M = 73.78 %, SD = 7.65), or across
 555 trials (Tables 10 and 11). In terms of individual differences in accuracy, one interaction
 556 between trials, conditions and WIAT-II comprehension scores was found to be marginally

557 significant (Table 11). The pattern observed suggested that when comprehension demands
 558 were high, high scorers on this test became more accurate over time, whereas low scorers
 559 became less accurate in later trials. However, these trends were marginal.

560 **Table 10. Binomial GLMM for Accuracy predicted by NDRT Comprehension and**
 561 **Interactions with Trial Number and Condition**

	β	95 % CI	z	p
Intercept	2.65	[1.54, 3.76]	4.66	< .001 ***
Trial Number	0.01	[-0.02, 0.03]	0.56	.576
Condition	-1.57	[-3.97, 0.84]	-1.28	.201
NDRT Comprehension	0.05	[-0.40, 0.49]	0.20	.841
Trial Number \times Condition	0.02	[-0.03, 0.07]	0.66	.510
Trial Number \times NDRT Comprehension	0.01	[-0.01, 0.04]	0.96	.338
Condition \times NDRT Comprehension	-0.09	[-0.95, 0.76]	-0.21	.832
Trial Number \times Condition \times NDRT Comprehension	0.00	[-0.06, 0.05]	-0.16	.873

562 The baseline of the condition term is lower comprehension demands. Estimates represent the
 563 change when going from lower to higher comprehension demands.

564 **Table 0.11. Binomial GLMM for Accuracy predicted by WIAT-II Comprehension and**
 565 **Interactions with Trial Number and Condition**

	β	95 % CI	z	p
Intercept	2.58	[1.49, 3.67]	4.63	< .001 ***
Trial Number	0.01	[-0.02, 0.04]	0.79	.430
Condition	-1.48	[-3.85, 0.88]	-1.23	.220
WIAT-II Comprehension	0.06	[-0.42, 0.53]	0.23	.815
Trial Number \times Condition	0.02	[-0.04, 0.07]	0.59	.558
Trial Number \times WIAT-II Comprehension	0.01	[-0.02, 0.04]	0.84	.401
Condition \times WIAT-II Comprehension	-0.59	[-1.51, 0.33]	-1.26	.209
Trial Number \times Condition \times WIAT-II Comprehension	0.05	[-0.01, 0.11]	1.69	.090 .

566 The baseline of the condition term is lower comprehension demands. Estimates represent the
 567 change when going from lower to higher comprehension demands.

568 **Discussion**

569 The current study investigated two offline reading comprehension tests (the NDRT and
 570 the WIAT-II) as predictors of individual differences in skilled readers' eye movements during
 571 paragraph reading. Eye movement patterns were investigated under higher and lower
 572 comprehension demands and across trials. Parallel sets of analyses were conducted for each
 573 test to determine whether individual differences in offline comprehension tests predicted
 574 patterns in eye movement behaviour that was reflective of changes in comprehension
 575 demands, and whether readers adapted to comprehension demands over time. The main aim
 576 was to determine whether discrepancies arose between the two tests that claim to measure
 577 reading comprehension [10,11], and a secondary aim was to investigate whether individual
 578 differences could be observed in the way that skilled readers adapted their reading strategies

579 over time and in response to comprehension demands. First, we will focus on the overall
580 patterns in the data across global eye movement measures, then on individual differences that
581 were observed, and finally, we will discuss the two offline comprehension tests and
582 differences in the predictive power associated with them for skilled readers.

583 **Overall patterns**

584 Overall, within an experimental block, paragraphs in later trials were read more quickly
585 than in earlier trials. Reading strategies appeared to become more efficient, or perhaps more
586 ‘risky’ [35] over time, with fewer fixations and increasing saccade lengths. Future
587 investigations would need to include analyses of regressions to determine whether readers do
588 use a riskier reading strategy in later trials since this may be more clearly observed though
589 rereading behaviours. Participants were not more accurate on the comprehension questions in
590 any one condition, or over time in the experiment. Though the increased difficulty in the
591 higher comprehension demands condition was confirmed by a pre-test, it may be that for our
592 skilled readers, the higher demands were not enough to reduce their accuracy. Indeed, the
593 pattern observed in the means suggested that participants had higher levels of accuracy when
594 comprehension demands were high, which would be compatible with previous observations
595 by Andrews and Veldre [44]. However, this difference was not significant in the analyses.
596 We did, however, observe differences in eye movement patterns in response to the higher and
597 lower comprehension demands. Readers were able to adjust their reading behaviours to the
598 comprehension demands [36] and were able to read more thoroughly when comprehension
599 demands were high and more superficially when comprehension demands were low [37].
600 Passages with higher comprehension demands were read more slowly, and featured more
601 fixations and shorter saccades than passages with lower comprehension demands.

602 **Individual differences**

603 Passage reading has the potential to introduce more variance in eye movement data
604 compared to sentence reading simply due to the increase in processing demands, and the
605 potential for allowing individual differences to be expressed in more varied ways. Slower
606 reading and more rereading is often observed during passages compared to sentences [16].
607 Although our data do not echo Andrews and Veldre's [44] observations of shorter passage
608 reading times, shorter average fixation durations and longer saccades directly related to
609 individual differences in reading proficiency, their findings were based on a composite
610 measure which included vocabulary, reading comprehension, reading rate and spelling, rather
611 than comprehension alone. It may be that the direct effects of individual differences on
612 fixation time measures observed by Andrews and Veldre [44] are better explained by other
613 measures included in their composite score (e.g., spelling or vocabulary). In our analyses,
614 individual differences as measured by offline comprehension measures seem to predict the
615 response to higher versus lower comprehension demands in the way that readers adapt over
616 time.

617 Analysis of eye movements in relation to NDRT comprehension scores presented a
618 clear picture of individual differences in response to comprehension demands. When reading
619 behaviours were measured across trials, there were observable individual differences in the
620 way that readers adapted their behaviour in response to comprehension demands. Differences
621 between readers were smaller at the beginning of the experimental blocks and became larger
622 in later trials where high scorers read more quickly, made fewer fixations and longer saccades
623 than low scorers. High scorers read passages with higher comprehension demands more
624 slowly, with more fixations and shorter saccades than passages with lower comprehension
625 demands, but the changes over time for higher and lower comprehension demands were

626 comparable. In contrast, low scorers adapted their reading behaviours at a slower rate and
627 approached a threshold for the fastest reading times, lowest number and duration of fixations,
628 and the largest saccade lengths they were able to accommodate whilst reading for
629 comprehension, even when comprehension demands were low. This evidence that less skilled
630 comprehenders have a lower limit to how quickly they can read for comprehension than
631 highly skilled comprehenders complements the general finding that less skilled readers often
632 read more slowly and make longer fixations than more skilled readers [28,29,44,10,11].

633 **Offline comprehension measures**

634 Critically, this pattern of results was highly dependent on which offline measure of
635 comprehension was used to measure comprehension. Analyses of the same participants' eye
636 movement data in relation to their scores on the WIAT-II comprehension test did not predict
637 differences in eye movement patterns for different comprehension demands. Earlier, we
638 described some differences in the format of each test that could indicate differences in the
639 underlying skills measured by them. We return to these now to consider possible reasons why
640 the NDRT revealed patterns in our data that the WIAT-II did not.

641 Higher comprehension accuracy is often observed for questions following narratives
642 than expository texts [17]. Therefore, expository passages were selected for the current study
643 to ensure that the materials were appropriate for skilled reading and to maximise the
644 likelihood of finding variation in accuracy scores within this population. Potentially as a
645 result of this choice, accuracy scores were not close to ceiling levels in the current study. The
646 NDRT includes expository texts that are more similar to the current study materials than the
647 WIAT-II comprehension test. Therefore, it is reasonable to suggest that comprehension based
648 on similar test materials will account for a comparatively larger proportion of variance in

649 reading behaviour. It has also been suggested that the NDRT is closely related to general
650 knowledge [20,21]. If the NDRT comprehension measure is highly related to general
651 knowledge, we would expect to see higher levels of comprehension accuracy on our
652 experimental questions for participants who score highly on the NDRT, but this was not
653 observed.

654 In contrast, since the WIAT-II comprehension test includes some items that must be
655 read aloud, it may feature some overlap with working memory processes [22]. However, our
656 previous investigations of eye movement behaviours in sentence reading included the WIAT-
657 II and a test of working memory (a backwards digit span task) amongst other reading skill
658 predictors [10,11]. These investigations did not suggest that there was much overlap between
659 working memory and the WIAT-II comprehension test as they did not load together in
660 principal components analyses [10,11]. We highlighted some aspects of the WIAT-II
661 comprehension subtest that may mean it also has less power to discriminate between skilled
662 adult readers than the NDRT. First, narrative test comprehension is often higher than
663 expository texts, which may indicate that portions of the WIAT-II comprehension test are not
664 difficult enough to allow much variance within skilled readers. We also noted in the
665 introduction that the reading aloud parts of the test may not be as informative about
666 individual differences in adults as it is for children since adults rely on reading aloud less
667 often [25], though further research would be needed to confirm this. In addition, the face-to-
668 face aspect of the WIAT-II may lead to noisier data for adults where participants might
669 experience performance anxiety.

670 It is also important that we acknowledge the potential impact that the low reliability of
671 the WIAT-II may have had on results in this study. Low reliability in the WIAT-II may be the
672 underlying reason for a weak correlation with the NDRT and null findings when predicting

673 eye movement measures. Future research will need to explore whether the low reliability of
674 the WIAT-II that we observed is due to the specific population of skilled readers our study
675 examined. Regardless, we maintain that when used to predict individual differences in eye
676 movement patterns in adult readers that researchers should be cautious when selecting an
677 appropriate test to use.

678 **Limitations**

679 As is very common for participant samples which are mostly based on Psychology
680 Undergraduate students, the current sample from the University of Southampton featured a
681 high proportion of females, which may limit the generalisation to male participants.

682 We also note that the NDRT was not administered with the standard time limit. A
683 precedent has been set by Andrews et al [51] for administering a shortened version of the
684 NDRT for researchers examining individual differences in skilled readers' eye movement
685 patterns since it increases the variance between skilled readers. This choice might limit the
686 comparability to research that uses the NDRT with the standard time limit. However, since
687 our focus was on skilled readers and the NDRT shortened time limit is increasing in
688 popularity in the field of examining skilled reading, we feel the choice for the shortened
689 version of the NDRT was justified. In addition, reliability estimates for both comprehension
690 tests based on our data were somewhat lower than the estimates given by each test manual
691 [12, 13], therefore we note that these tests may have comparatively reduced reliability for
692 university level populations.

693 All questions in the experimental conditions had two options from which participants
694 were required to select an answer. Such limited response options may have limited the
695 capacity to find differences in accuracy in our data. However, we note that this would not
696 limit findings drawn from the eye movement record.

697 **Conclusion**

698 Overall, it appears that the NDRT comprehension test (notably when following a half-
699 timed procedure) is more sensitive to differences in eye movement behaviours in response to
700 higher and lower comprehension demands observed between skilled adult readers compared
701 to the WIAT-II comprehension test. Individual differences captured by the half-timed version
702 of the NDRT have been previously shown to be sensitive to individual differences in skilled
703 readers eye movements [50]. The current study extends this and suggests it can be used to
704 predict differences in eye movement behaviours across trials in response to varying
705 comprehension demands. We highlight the importance of careful test selection when
706 measuring eye movement behaviour in skilled adult readers and advise that comprehension
707 tests should not be used interchangeably, because they *jingle* [14] and that researchers should
708 exercise caution when selecting a reading comprehension test for future research. We echo
709 advice from Flake and Fried [15] who call for transparency when reporting test selection
710 processes and urge researchers to select comprehension tests that are clearly based on the
711 theoretical concepts that the researcher wishes to assess.

712 **Acknowledgements**

713 Special thanks to Karolina Vakulya for her work on data collection for this research.

714 **References**

- 715 1. Kintsch W. *Comprehension: A paradigm for cognition*. Cambridge University Press;
716 1998.
- 717 2. Perfetti C, Stafura J. Word knowledge in a theory of reading comprehension. *Sci Stud*
718 *Read*. 2014;18(1):22-37. doi:10.1080/10888438.2013.827687
- 719 3. Van Den Broek P, Young M, Tzeng Y, Linderholm T. The landscape model of
720 reading. In: van Oostendorp H, Goldman SR, eds. *The construction of mental*
721 *representations during reading* (pp. 71-98). Mahwah: Erlbaum; 1999.
- 722 4. Gough PB, Tunmer WE. Decoding, reading, and reading disability. *RASE: Remedial*
723 *Spec Educ*. 1986;7(1):6–10. doi:10.1177/074193258600700104
- 724 5. Kintsch W, Rawson KA. Comprehension. In: Snowling MJ, Hulme C, editors. *The*
725 *science of reading: A handbook* (pp. 209-226). Blackwell Publishing; 2005.
726 doi:10.1002/9780470757642.ch12
- 727 6. Cutting LE, Scarborough HS. Prediction of reading comprehension: Relative
728 contributions of word recognition, language proficiency, and other cognitive skills can
729 depend on how comprehension is measured. *Sci Stud Read*. 2006;10:277-299.
730 doi:10.1207/s1532799xssr1003_5
- 731 7. Keenan JM, Betjemann RS, Olson RK. Reading comprehension tests vary in the skills
732 they assess: Differential dependence on decoding and oral comprehension. *Sci Stud*
733 *Read*. 2008;12(3):281-300. doi:10.1080/10888430802132279
- 734 8. Kendeou P, Papadopoulos TC, Spanoudis G. Processing demands of reading
735 comprehension tests in young readers. *Learn Instr*. 2012;22(5):354-367.
736 doi:10.1016/j.learninstruc.2012.02.001

- 737 9. Mézière DC, Yu L, Reichle ED, von der Malsburg T, McArthur G. Using Eye-
738 Tracking Measures to Predict Reading Comprehension. *Read Res Q.* 2023.
739 doi:10.1002/rrq.498
- 740 10. Lee CE, Godwin HJ, Blythe HI, Drieghe D. Individual differences in skilled reading
741 and the word frequency effect. (Manuscript submitted for publication).
- 742 11. Lee CE, Pagán A, Godwin HJ, Drieghe D. Individual differences and the transposed
743 letter effect during reading. (Manuscript submitted for publication).
- 744 12. Wechsler D. Wechsler Individual Achievement Test 2nd Edition (WIAT-II II).
745 London: The Psychological Corp; 2005.
- 746 13. Brown JA, Fishco VV, Hanna G. Nelson-denny reading test: Manual for scoring and
747 interpretation, forms G & H. Riverside Publishing; 1993.
- 748 14. Thorndike EL. An introduction to the theory of mental and social measurements.
749 Science Press, New York: Teacher's College, Columbia University; 1904.
750 doi:10.1037/13283-000
- 751 15. Flake JK, Fried EI. Measurement schmeasurement: Questionable measurement
752 practices and how to avoid them. *Adv Methods Pract Psychol Sci.* 2020;3(4):456-465.
753 doi:10.1177/2515245920952393
- 754 16. Radach R, Huestegge L, Reilly R. The role of global top-down factors in local eye-
755 movement control in reading. *Psychol Res.* 2008;72:675-688. doi:10.1007/s00426-
756 008-0173-3
- 757 17. Best RM, Floyd RR, McNamara DS. Differential competencies contributing to
758 children's comprehension of narrative and expository texts. *Read Psychol.*
759 2008;29(2):137-164. doi:10.1080/02702710801963951

- 760 18. McNamara DS, Ozuru Y, Floyd RG. Comprehension challenges in the fourth grade:
761 The roles of text cohesion, text genre, and readers' prior knowledge. *Int Electr J Elem*
762 *Educ.* 2011;4(1):229-257.
- 763 19. Zwaan RA. Effect of genre expectations on text comprehension. *J Exp Psychol Learn*
764 *Mem Cogn.* 1994;20(4):920-933. doi:10.1037/0278-7393.20.4.920
- 765 20. Ready RE, Chaudhry MF, Schatz KC, Strazzullo S. "Passageless" administration of
766 the Nelson-Denny reading comprehension test: Associations with IQ and reading
767 skills. *J Learn Disabil.* 2013;46(4):377-384. doi:10.1177/0022219412468160
- 768 21. Coleman C, Lindstrom J, Nelson J, Lindstrom W, Gregg KN. Passageless
769 comprehension on the Nelson-Denny reading test: well above chance for university
770 students. *J Learn Disabil.* 2010;43(3):244-9. doi:10.1177/0022219409345017.
- 771 22. Laubrock J, Kliegl R. The eye-voice span during reading aloud. *Front Psychol.*
772 2015;6(1432). doi:10.3389/fpsyg.2015.01432
- 773 23. Hale AD, Skinner CH, Williams J, Hawkins R, Neddenriep CE, Dizer J. Comparing
774 comprehension following silent and aloud reading across elementary and secondary
775 students: Implication for curriculum-based measurement. *Behav Anal Today.* 2007;8.
776 doi:10.1037/h0100101
- 777 24. García JR, Cain K. Decoding and reading comprehension: A meta-analysis to identify
778 which reader and assessment characteristics influence the strength of the relationship
779 in English. *Rev Educ Res.* 2014;84(1):74-111. doi:10.3102/0034654313499616
- 780 25. Duncan S, Freeman M. Adults reading aloud: A survey of contemporary practices in
781 Britain. *Br J Educ Stud.* 2020;68(1):97-123. doi:10.1080/00071005.2019.1610555
- 782 26. Gambrell LB, Heathington BS. Adult disabled readers' metacognitive awareness
783 about reading tasks and strategies. *J Read Behav.* 1981;13(3):215-222.
784 doi:10.1080/10862968109547409

- 785 27. Liversedge SP, Findlay JM. Saccadic eye movements and cognition. *Trends Cogn*
786 *Sci.* 2000;4(1):6-14. doi:10.1016/s1364-6613(99)01418-7
- 787 28. Rayner K. Eye movements in reading and information processing: 20 years of
788 research. *Psychol Bull.* 1998;124(3):372. doi:10.1037/0033-2909.124.3.372
- 789 29. Ashby J, Rayner K, Clifton C. Eye movements of highly skilled and average readers:
790 Differential effects of frequency and predictability. *Q J Exp Psychol A.*
791 2005;58(6):1065–1086. doi:10.1080/02724980443000476
- 792 30. Rayner K, Schotter ER, Masson MEJ, Potter MC, Treiman R. So much to read, so
793 little time: How do we read, and can speed reading help? *Psychol Sci Public Interest.*
794 2016;17(Supplement 1):4-34. doi:10.1177/1529100615623267
- 795 31. Rayner K. Eye movements and attention in reading, scene perception, and visual
796 search. *Q J Exp Psychol.* 2009;62(8):1457-1506. doi:10.1080/17470210902816461
- 797 32. Andrews S. Individual differences in skilled visual word recognition and reading: The
798 role of lexical quality. In Adelman J, editor. *Visual Word Recognition* (pp. 151-172).
799 Hove: Psychology Press; 2012.
- 800 33. Bisanz GL, Das JP, Varnhagen CK, Henderson HR. Structural components of reading
801 time and recall for sentences in narratives: exploring changes with age and reading
802 ability. *J Educ Psychol.* 1992;84(1):103–114. doi:10.1037/0022-0663.84.1.103
- 803 34. Stanovich KE. Towards an interactive-compensatory model of individual differences
804 in the development of reading fluency. *Read Res Q.* 1980;16:32-71.
805 doi:10.2307/747348
- 806 35. O'Regan JK. Optimal viewing position in words and the strategy-tactics theory of eye
807 movements in reading. In: Rayner K, editor. *Eye movements and visual cognition:*
808 *Scene perception and reading* (pp. 333-354). New York: Springer-Verlag; 1992.
809 doi:10.3758/BF03213829

- 810 36. Tinker M. Recent studies of eye movements in reading. *Psychol Bull.* 1958;55:215-
811 231.
- 812 37. Heller D. Eye movements in reading. In: Groner R, Fraisse P, editors. *Cognition and*
813 *eye movements* (pp. 139-154). Amsterdam: North Holland; 1982.
- 814 38. Aaronson D, Ferres S. Reading strategies for children and adults: Some empirical
815 evidence. *J Verbal Learning Verbal Behav.* 1984;23(2):189–220. doi:10.1016/S0022-
816 5371(84)90137-3
- 817 39. Aaronson D, Ferres S. Reading strategies for children and adults: A quantitative
818 model. *Psychol Rev.* 1986;93(1):89–112. doi:10.1037/0033-295X.93.1.89
- 819 40. Inhoff AW, Radach R. Definition and computation of oculomotor measures in the
820 study of cognitive processes. In: Underwood G, editor. *Eye guidance in reading and*
821 *scene perception* (pp. 29-53). Elsevier Science Ltd; 1998. doi:10.1016/B978-
822 008043361-5/50003-1
- 823 41. Wotschack C, Kliegl R. Reading strategy modulates parafoveal-on-foveal effects in
824 sentence reading. *Q J Exp Psychol.* 2013;66:548-562.
825 doi:10.1080/17470218.2011.625094140S.
- 826 42. Weiss AF, Kretschmar F, Schlesewsky M, Bornkessel-Schlesewsky I, Staub A.
827 Comprehension demands modulate rereading, but not first-pass reading behavior. *Q J*
828 *Exp Psychol.* 2018;71:198-210. doi:10.1080/17470218.2017.1307862
- 829 43. Christianson K, Luke SG, Hussey EK, Wochna KL. Why reread? Evidence from
830 garden-path and local coherence structures. *Q J Exp Psychol.* 2017;70:1380–1405.
831 doi:10.1080/17470218.2016.1186200
- 832 44. Andrews S, Veldre A. Wrapping up sentence comprehension: The role of task
833 demands and individual differences. *Sci Stud Read.* 2021;25(2):123–140.
834 doi:10.1080/10888438.2020.1817028

- 835 45. Andrews S, Hersch J. Lexical precision in skilled readers: Individual differences in
836 masked neighbor priming. *J Exp Psychol Gen.* 2010;139(2):299–318.
837 doi:10.1037/a0018366
- 838 46. Just MA, Carpenter PA. A theory of reading: From eye fixations to comprehension.
839 *Psychol Rev.* 1980;87(4):329-354. doi:10.1037/0033-295X.87.4.329
- 840 47. Aaronson D, Scarborough HS. Performance theories for sentence coding: Some
841 quantitative evidence. *J Exp Psychol Hum Percept Perform.* 1976;2:56–70.
842 doi:10.1037/0096-1523.2.1.56
- 843 48. Rayner K, Kambe G, Duffy SA. The effect of clause wrap-up on eye movements
844 during reading. *Q J Exp Psychol Sect A.* 2000;53(4):1061-1080.
845 doi:10.1080/713755934
- 846 49. Stine-Morrow EAL, Milinder L, Pullara O, Herman B. Patterns of resource allocation
847 are reliable among younger and older readers. *Psychol Aging.* 2001;16:69-84.
848 doi:10.1037/0882-7974.16.1.69
- 849 50. Andrews S, Veldre A, Clarke IE. Measuring lexical quality: The role of spelling
850 ability. *Behav Res Methods.* 2020;52:2257-2282. doi:10.3758/s13428-020-01387-3
- 851 51. Andrews S. Individual differences among skilled readers: The role of lexical quality.
852 In Pollatsek A, Treiman R, editors. *The Oxford Handbook of Reading* (pp. 129–148).
853 Oxford University Press; 2015.
- 854 52. Veldre A, Andrews S. Parafoveal lexical activation depends on skilled reading
855 proficiency. *J Exp Psychol Learn Mem Cogn.* 2015;41(2):586-595.
856 doi:10.1037/xlm0000039
- 857 53. Veldre A, Andrews S. Parafoveal preview benefit is modulated by the precision of
858 skilled readers' lexical representations. *J Exp Psychol Hum Percept Perform.*
859 2015;41(1):219-232. doi:10.1037/xhp0000017

- 860 54. Perfetti C. Reading ability: Lexical quality to comprehension. *Sci Stud Read.*
861 2007;11(4):357-383. doi:10.1080/10888430701530730
- 862 55. Veldre A, Andrews S. Lexical quality and eye movements: Individual differences in
863 the perceptual span of skilled adult readers. *Q J Exp Psychol.* 2014;67(4):703-727.
864 doi:10.1080/17470218.2013.826258
- 865 56. SR Research. *Eyelink User Manual 1.0.12.* Mississauga, Ontario, Canada; 2017.
- 866 57. Determine the main idea of the passage. [Internet]. 2021 Oct 27. Available from:
867 <http://www.uk.ixl.com>
- 868 58. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using
869 *lme4.* *J Stat Softw.* 2015;67(1):1–48. doi:10.18637/jss.v067.i01.
- 870 59. R Core Team. *R: A language and environment for statistical computing.* R Foundation
871 for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>; 2022.
- 872 60. Barr DJ, Levy R, Scheepers C, Tily HJ. Random effects structure for confirmatory
873 hypothesis testing: Keep it maximal. *J Mem Lang.* 2013;68:255-278.
874 doi:10.1016/j.jml.2012.11.001