

**Dynamic assessment of the effectiveness of digital game-based literacy training
in beginning readers**

Toivo Glatz^{1,2,3}, Wim Tops¹, Elisabeth Borleffs¹, Ulla Richardson⁴, Natasha Maurits^{2,5}, Annemie Desoete⁶ & Ben Maassen^{1,2,7}

¹Center for Language and Cognition (CLCG), Faculty of Arts, University of Groningen, Groningen, Netherlands

²Behavioural and Cognitive Neuroscience (BCN), University Medical Center Groningen (UMCG), Groningen, Netherlands

³Institute of Public Health, Charité – Universitätsmedizin Berlin, corporate member of Freie Universität Berlin and Humboldt-Universität zu Berlin, Berlin, Germany

⁴Centre for Applied Language Studies, University of Jyväskylä, Jyväskylä, Finland

⁵Department of Neurology, University Medical Center Groningen (UMCG), Groningen, Netherlands

⁶Department of Developmental Disorders, Ghent University, Ghent, Belgium

⁷Department of Neuroscience, University Medical Center Groningen (UMCG), Groningen, Netherlands

Corresponding author: Toivo Glatz^{1,2}

Email address: toivo.glatz@gmail.com

1 **Introduction**

2 Adequate early literacy instruction and well-developed literacy skills are indispensable for a child's
3 academic success and future career. It is therefore important to know how we can improve teaching
4 methods and accurately monitor reading progress. What tools do we have or need that help promote
5 reading skills, detect problems at an early stage and prevent struggling readers from developing more
6 serious literacy problems such as dyslexia? In this context, digital game-based learning shows
7 potential.

8
9 Developmental dyslexia or specific learning disorder in reading (DSM-5; APA, 2013; henceforth,
10 'dyslexia') is a developmental disorder characterized by persistent difficulties in word recognition
11 (reading) and/or spelling. These difficulties are not caused by a general cognitive delay or by a
12 hearing or vision impairment. Depending on a narrow or wider definition of poor reading
13 proficiency, this developmental disorder affects around 4 to 12% of children across languages (e.g.
14 Schulte-Körne, Deimel, Bartling, & Remschmidt, 1998; Schumacher, Hoffmann, Schmal, Schulte-
15 Körne, & Nothen, 2007). Language and orthography both play an important role in reading (Borleffs,
16 Maassen, Lyytinen, & Zwarts, 2017), with the prevalence of dyslexia differing across countries
17 depending on their characteristics (Bergmann & Wimmer, 2008; Ziegler & Goswami, 2005). As a
18 consequence of differences in the mapping of grapheme-phoneme correspondences, the
19 developmental trajectory and nature of the reading problems may also differ between languages with
20 regular and less transparent orthographies (Seymour, Aro, & Erskine, 2003; Bergmann & Wimmer,
21 2008; Ziegler & Goswami, 2005).

22
23 Dyslexia has been shown to be a disorder with a multifactorial aetiology in that it is associated with a
24 range of genetic, environmental, and cognitive risk factors rather than with a single cause
25 (Pennington, 2006). If one parent or sibling has dyslexia, the incidence rate rises to around 45%,

26 indicating a familial risk (for a review, see Snowling & Melby-Lervag, 2016). However, genetic
27 risks do not operate in isolation. They, for instance, interact with environmental factors such as
28 parental socioeconomic status (Mascheretti et al., 2013). There are certain early childhood cognitive
29 and behavioural precursors of future reading skills. The most prominent ones are letter knowledge as
30 a measure of grapheme-phoneme correspondences knowledge, phonological awareness, and rapid
31 automatized naming of familiar objects. Early below-average performance on tasks targeting these
32 skills increases the risk of dyslexia (van der Leij, Bergen, Zuijlen, Jong, Maurits, & Maassen, 2013;
33 Lyon, Shaywitz, & Shaywitz, 2003; Lyytinen et al., 2004; Lyytinen, Erskine, Kujala, Ojanen, &
34 Richardson, 2009). Phonological awareness and rapid automatized naming have been shown to
35 predict reading speed and accuracy across languages (Moll et al., 2014). However, certain cross-
36 linguistic variability exists with respect to the relative weight of each of the cognitive and
37 behavioural precursors of reading acquisition (Landerl et al., 2013; Ziegler et al., 2010): letter
38 knowledge being most predictive in Finnish with its extreme letter-sound consistency (Lyytinen et
39 al., 2009), rapid automatized naming being the best long-term predictor in German (Brem et al.,
40 2013), and letter knowledge, rapid automatized naming and phonological awareness being important
41 indicators in Dutch (van Bergen, Jong, Plakas, Maassen, & van der Leij, 2012). An important
42 distinction has to be made for letter knowledge, which just reflects the availability of letter-sound
43 associations that quickly reach ceiling and therefore has limited use for long-term predictions.
44 Adding time pressure within a speeded letter-speech sound identification task yields a measure
45 which is even more specifically related to the fluency of multimodal processing of audio-visual
46 information (Blomert, 2011; Hahn, Foxe, & Molholm, 2014). In addition to the availability of letter-
47 sound association this 'timed' letter knowledge also assesses the fluency in which these associations
48 can be retrieved from memory.

Commented [GM1]: There seems to be a jump in topic here that is difficult for the reader to follow. Provide a better segue.

Commented [GM2]: This is a reading skill. Do you mean predictors?

Formatted: Strikethrough

Commented [GM3]: Rephrase. This suggests that performance on tasks increases the risk of dyslexia, which is not what you mean, I presume

Commented [GM4]: I am not sure about this. Letter knowledge (ie the names of letters) is not the same as letter-sound knowledge. Some letter names provide hints about a letter-sound rule (T), but some do not (W).

Commented [GM5]: This paragraph is difficult for the reader to follow. It "jumps around" a lot from one topic to another without providing the reader with a logical flow of information. It is also very long and includes multiple topics.

I think I can see what you are trying to do here, but please revise so that the flow of your logic is simpler to follow for the reader.

One way you may do this is to use this paragraph to just focus on early predictors of later reading: letter knowledge, PA, and RAN - since that is where you want to "land" in the next paragraph anyway.

50 Given its multifactorial aetiology and early indicators such as poor letter knowledge, phonological
 51 awareness and rapid automatized naming, the question arises whether timely training of these skills
 52 might help remediate or even prevent reading difficulties. During the past decade, it has been shown
 53 that one promising way to deliver such a training is by computerized gaming (Chambers et al., 2008;
 54 Richardson and Lyytinen, 2014). While there is no commonly agreed definition for digital game-
 55 based learning or so called serious games (Susi, Johannesson, & Backlund, 2007) such tools have a
 56 range of interesting characteristics (De Freitas, 2006; Kiili, 2005; Prensky, 2001), as they: i) offer
 57 multimodal learning environment (Potocki, Ecalle, & Magnan, 2013), ii) provide immediate
 58 feedback for improved learning, iii) can adapt to individual learners depending on their responses, iv)
 59 are highly motivating for the players (e.g. Desoete, Praet, Van de Velde, De Craene, & Hantson,
 60 2016), and v) can monitor the development of accuracy and response times in task-relevant contexts
 61 and provide researchers with longitudinal data (e.g. Praet & Desoete, 2014; Puolakanaho & Latvala,
 62 2017).

64 GraphoGame is one such promising computerized training method targeting reading-related skills
 65 (for a review, see Richardson & Lyytinen, 2014). It is an adaptive, child-friendly digital learning
 66 environment that aims to help beginning readers. GraphoGame was originally designed for the very
 67 transparent orthography of Finnish as a tool to boost grapheme-phoneme correspondence knowledge
 68 in beginning readers. (REF) This is done by establishing accurate phonemic representations (X),
 69 connecting these to the orthographic stimuli (X), and establishing a fluent association between the
 70 two. More recent versions of the game also train phonological awareness, spelling, and reading
 71 fluency (Richardson & Lyytinen, 2014). A number of experimental studies using GraphoGame have
 72 been conducted in over 20 different languages, with mixed results. While most interventions led to
 73 improvements in letter knowledge, only a few demonstrated benefits for reading skills. A recent
 74 meta-analysis revealed that the average gain in reading performance across 19 GraphoGame studies

Commented [GM6]: predictors?

Formatted: Font colour: Red

Commented [GM7]: You have lost me here. Rephrase.

Commented [GM8]: which?

Formatted: Font colour: Red

Formatted: Font colour: Red

Formatted: Font colour: Red

Formatted: Font colour: Red

Commented [GM9]: what do you mean by interesting?

Formatted: Font colour: Red

Formatted: Font colour: Red

Commented [GM10]: why is this interesting?

Commented [GM11]: It is not clear why these features are interesting/important etc. Rewrite the sentence in red to make your point clearer.

Also, use PeerJ formatting for lists.

Commented [GM12]: You have not mentioned "promising" games immediate above. Rephrase to tie into information in previous paragraph more coherently.

Commented [GM13]: Provide reference to original version

Commented [GM14]: define what this is

Commented [GM15]: define what this is

Commented [GM16]: rephrase. Do you mean reliable? Strong?

Deleted: overall

Commented [GM17]: more advanced? other?

Commented [GM18]: what type of reading skills?

Deleted: terms of

Deleted: actual

Deleted:

Commented [GM19]: what type of reading performance?

79 was close to zero (McTigue, Solheim, Zimmer, Upstad, 2019). However, these studies differed in
80 various methodological aspects, including selection criteria for participants (provide range), age of
81 participants (5 to 10 years), the number and types of control groups (provide range), sample sizes
82 ($N = 10$ to 185), training hours (1 to 8 hours), training implementation (during school or at home,
83 with or without adult engagement), training period (1 to 28 weeks), and type of language (ranging
84 from transparent (Finnish) to opaque (???). An additional issue is that most standardised (reading)
85 tests are not designed to detect the subtle changes that occur over a few weeks of training. The large
86 variation in all these variables makes it challenging to grasp the characteristics that make up a
87 successful intervention. Unfortunately, many studies also lack relevant details, which further impede
88 comparison, be it in form of a detailed outline of the actual training material, the cognitive skills a
89 game is meant to train, or the training environment. Note that for sake of comparability we limit our
90 discussion to GraphoGame studies, but the same questions arise with many other literacy
91 interventions using digital game-based learning, such as ABRACADABRA (Savage et al., 2013;
92 Piquette, Savage, & Abrami, 2014) or Bouw! (Regtvoort, Zijlstra, & Leij, 2013; Zijlstra, Van
93 Bergen, Regtvoort, De Jong, & Van Der Leij, 2020).

94

95 In sum, the divergent outcomes of digital game-based learning studies in general, and GraphoGame
96 studies in particular, raise the question which circumstances impact the effectiveness of literacy
97 digital game-based learning. Deepening our understanding of how progress can be optimally
98 measured, how game content and training frequency and intensity modulate training outcomes, and
99 for which populations of children this approach works best might help us to improve the
100 effectiveness of digital game-based learning for early literacy. Therefore, in this work we focus on i)
101 the tools that are generally used to assess reading-related cognitive skills, ii) the characteristics of the
102 populations studied, and iii) training properties such as duration, intensity, and content.

103

Deleted: all

Deleted: also

Deleted: . Be it in terms of population properties like

Deleted:

Deleted: or

Deleted: per experimental group

Deleted: . Similarly, there are drastic differences in terms of

Deleted: , i.e. whether training took place in

Deleted: at

Deleted: during hours or after hours,

Deleted: how long the

Deleted: was

Deleted: ranging from one

Deleted: and the overall exposure to the game (one to eight hours on task

Deleted:). Linguistic properties form a third pillar of often discussed differences, with orthographies ranging from transparent to opaque and each language having its own unique mix of cognitive precursors of reading

Deleted: Therefore, training goals and content, as well as assessments of language skills will inevitably differ.

Deleted: regarding the latter

Commented [GM20]: You need to link this to differences between studies of GG. Are you trying to say that different studies used different standardised tests that vary in the degree to which they are sensitive to GG training effects? If not, rephrase to keep the discussion on point.

Deleted: of

Deleted: factors

Commented [GM21]: I presume you mean GG intervention? Clarify.

Commented [GM22]: At this point, the reader is wondering why your team would be interested in actually using GG. I think what is missing is a discussion of studies that ***have*** found an effect. This could be followed by an explanation that a number of studies have **failed** to find this effect (and cite them). And ***then*** move on to why there may be such mixed outcomes (ie the variables you have listed). Don't forget to cite as many GG studies as you can to recognise their contribution to the field.

Commented [GM23]: I would steer clear of this. You are opening up a can or worms that you cannot address in this study.

Formatted: Strikethrough

Formatted: Strikethrough

Commented [GM24]: I think you can remove this paragraph.

128 **The current study,**

129 Our main aim is to evaluate the effectiveness of a newly created Dutch version of GraphoGame.

130 More specifically, we want to investigate possible impacts of assessment tools, population

131 characteristics, and intervention properties on training outcomes. For this purpose, we invited entire

132 first grade classrooms to play GraphoGame-NL for up to seven weeks. We used traditional paper-

133 and-pencil as well as in-game tests to track their performance.

134 **Research Question 1:** *Does the evaluation of game effectiveness depend on the choice of assessment*

135 *tools used to track changes in performance?* First, we tested the effects of gaming on different

136 reading-related skills: i.e., direct assessment of word reading, on the one hand, and reading-related

137 skills, on the other hand. Second, we integrated comparisons of different assessment tools for the

138 same skill (see Methods). We hypothesized that tests of skills prerequisite for reading were more

139 sensitive for measuring training effectiveness than tests of reading fluency itself. Moreover, we

140 presumed that online measures such as response times would be overall more sensitive than offline

141 metrics (e.g., paper-and-pencil test).

142 **Research Question 2:** *How do population characteristics impact intervention effectiveness?* We

143 hypothesized that children who had below average performance in letter knowledge and

144 phonological awareness prior to the training would benefit most from GraphoGame-NL. We also

145 expected that certain subgroups of children perform below average at pre-test. These would

146 presumably be children at familial risk for dyslexia, younger children, or those speaking a foreign

147 language at home.

148 **Research Question 3:** *How do intervention properties contribute to training effectiveness?* From the

149 training phase itself, we acquired multiple metrics of child's gaming process: i.e., how many levels

150 were played (incl. levels that had to be repeated), the highest level reached by the child, and how

151 many targets and distractors they saw throughout the gaming. We hypothesized that these game

152 characteristics contributed to the variability in training outcomes.

Deleted: Present

Deleted: : research questions and hypotheses

Commented [GM25]: Provide 1-2 sentences that provides a nice segue between the previous paragraph (which will be two paragraphs above) and this new section. This will help the reader follow your story/logic for the study.

Commented [GM26]: Make sure these terms align/match the terms you use in the previous paragraph, otherwise the reader will not understand what these variables/factors actually are.

Commented [GM27]: Put this as your main aim. Then explain that you addressed this aim using GG with children in The Netherlands, and explain why you chose to do this from a scientific POV (ie what is special about the Dutch version that gives is special powers to address your aims)?

Formatted: Strikethrough

Commented [GM28]: This section has the same problem as the paragraph on page 3. It reads like a list of not quite connected ideas. Also, predictions are made without any justifications for those predictions or without referring to previous GG studies.

I suggest you rewrite this section restating your questions as aims. Then make predictions based on previous findings, which you mention briefly and clearly cite. If you cannot make any predictions, that is fine. Just explain why you cannot (ie no previous evidence to guide you).

155

156 **Materials and methods**

157 This research was approved by the ethical committee of the Faculty of Arts of the University of

158 Groningen, and the Faculty of Psychology of the University of Ghent. We aimed to follow the

159 CONSORT standards of randomized control trials (Schulz, Altman, & Moher, 2010).

160

161 **Participants**

162 Mainstream primary schools in the northern region (Groningen area) of the Netherlands (NL) and the

163 western region (Ghent area) of the Dutch-speaking part of Belgium (B) were contacted by phone or

164 letter and invited to join the study. Initial requirements for participation were i) schools had enough

165 computers with headphones for students to play GraphoGame on a daily basis, ii) classroom teachers

166 allowed students to play the game for 10-15 minutes per day for at least five weeks, iii) schools

167 allowed (who?) to administer behavioural tests to students at school during regular school hours, and

168 iv) teachers agreed to adhere to their allocated experimental condition. Eight schools were willing

169 and eligible to participate (three in the Netherlands, five in Belgium), which comprised 16 first-grade

170 classrooms (four in the Netherlands, 12 in Belgium) and 312 children (107 in the Netherlands, 205 in

171 Belgium).

172

173 Parents of participating children were asked for their written informed consent for the gaming and

174 additional behavioural tests. Parents of participating children were asked to complete a questionnaire

175 about their child's handedness, language(s)/dialect(s) spoken at home, reading history, families

176 reading history, problems in the child and/or in close family members, presence of confirmed

177 neurological problems, and medication. To enable as many children as possible to play the game

178 there were no initial exclusion criteria, and parents were also given the option to consent to the

179 gaming without additional behavioural assessments. A few families made use of that latter option

Commented [GM29]: Can you double check the formatting rules of PeerJ and make sure this manuscript aligns with those?

This section should be called Methods. I think there should be a separate section for ethics.

Formatted: Strikethrough

Commented [GM30]: To conduct the study? Or the report the study? And did you actually achieve that aim? (the meaning of this sentence is not clear).

I don't think you have followed reporting guidelines, which I would suggest you since this provides a clearer way of presenting your Methods. So, I would suggest finding the latest version of CONSORT reporting, and arrange the information in your Methods into the sections that

... [1]

Deleted: in which they were

Deleted: our

Deleted: to

Deleted: e

Deleted: that

Deleted: needed

Deleted: to have

Deleted: available allowing their first-grade

Deleted: the

Deleted: game

Deleted: expressed their willingness to motivate and e ... [2]

Commented [GM31]: when were the tests administ ... [3]

Deleted: additional

Deleted: could be carried out

Deleted: accepted to

Deleted: the (at that moment still unknown)

Deleted: their classroom would be assigned to

Deleted: We found e

Deleted: with a total of

Deleted: a potential of

Deleted: All

Deleted: p

Commented [GM32]: Put this info under Ethics

Deleted: Furthermore, they had

Deleted: inquiring

Deleted: ren

Formatted: Strikethrough

Formatted: Strikethrough

Deleted: in the child

Deleted: potential

Commented [GM33]: Put under Ethics

208 ($N = 4$) while some more did not consent to participation at all ($N = 26$). For final analysis we
209 excluded data of those children that were repeating the first grade ($N = 2$), one individual that was
210 one year older than its peers without repeating first grade, those children that were diagnosed with a
211 neurodevelopmental disorder ($N = 5$), those that missed an assessment ($N = 8$), and those who failed
212 to play at least 20 sessions (corresponding to four weeks of daily playing) or alternatively
213 accumulate at least 2.5 hours of game exposure ($N = 19$). The details of this sample are specified in
214 the CONSORT flow diagram (Figure 1).

216 To avoid that children would play the wrong game, they were assigned to an experimental condition
217 by classroom (i.e., cluster randomized) and not individually. Classrooms were semi-randomly
218 assigned to either play the reading version of GraphoGame-NL, a math version of GraphoGame
219 (active controls), or attend the normal school curriculum (passive controls). Importantly, even the
220 children of this passive group took part in two gaming sessions to complete the in-game assessments
221 at pre- and post-test. Where possible, a within-school design was set up: three schools participated
222 with three classrooms, so each classroom per school was randomly assigned to one of the three
223 experimental conditions. Another two schools joined with two classrooms, which were randomly
224 assigned to the reading and math conditions. The final three schools joined with one classroom and
225 each classroom was once more randomly assigned to one of the three conditions. See Table 1 for an
226 overview of the number of children per classroom and their experimental conditions, and Table 2 for
227 pre-test results.

229 Computerised training

230 Our research group specifically created a Dutch version GraphoGame for the present literacy study.
231 The content of which was selected from Veilig leren lezen (VLL; 'Learning to read safely';
232 Mommers, Verhoeven, & van der Linden, 1990) a widely used literacy teaching method in the

Commented [GM34]: This should be part of a Design section at start of Methods

Commented [GM35]: YOU need to indicate where to insert tables within the manuscript

233 Netherlands and a vocabulary achievement list for six-year olds (Schaerlaekens, Kohnstamm,
234 Lajaegere, & Vries, 1999). The game included 650 items, ranging from simple and complex
235 graphemes (e.g. ⟨n⟩, ⟨r⟩, ⟨ui⟩), to CV/VC syllables either representing separate words or occurring as
236 parts of existing words (e.g. vi / is), to monosyllabic words with CVC structure (e.g. vis, 'fish') or
237 targets with CCVC, CVCC or CCVCC consonant clusters (e.g. prijs, 'price'; zwart, 'black'). For a
238 detailed description of the tasks and materials used within the game, see Appendix 1. We excluded a
239 few infrequent complex graphemes (⟨ch⟩, ⟨sch⟩, ⟨aai⟩, ⟨auw⟩, ⟨eeuw⟩, ⟨ieuw⟩, ⟨oei⟩, ⟨ouw⟩) that are
240 not typically taught at the beginning of the first grade. We also created a limited number of
241 phonotactically legal pseudowords as minimal pairs using Wuggy (Keuleers & Brysbaert, 2010).
242 Five female students studying linguistics or speech-language pathology at the University of
243 Groningen spoke the auditory stimuli. Naïve native speakers of Dutch subsequently evaluated all
244 items with respect to their prototypicality and only the most prototypical items were then included in
245 the game, yielding one to four different spoken realizations per target. It needs to be noted here that
246 there are some systematic differences in pronunciation between the Dutch spoken in the Netherlands
247 and the Dutch variety or Flemish spoken in the northern part of Belgium (for phonetic distances
248 between Dutch dialects, see Nerbonne, Heeringa, Van den Hout, Van der Kooi, Otten, & Van de Vis,
249 1996). This should not be a big problem as Flemish children also have exposure to standard Dutch
250 through multimedia (movies, series, games, etc.).

252 A mathematics game specifically designed for this research was used as an active control condition.
253 Its framework was identical to that of the reading game, featuring a range of similar
254 reactive/interactive mini-games with varying graphics and task demands with the levels now
255 containing number/digit knowledge, counting, comparison of numbers and amounts, sorting of
256 adjacent or nonadjacent numbers in ascending or descending order, as well as simple addition and

Commented [GM36]: There is a lot of required information missing here about implementation of the program etc

257 subtraction. The range of numbers within the training goes from zero up to 20, thus mirroring the
258 classroom content of the first half of first grade.

259

260 **Assessment**

261 The following offline paper-and-pencil tests were used as outcome measures:

262

263 **Reading fluency:** Participants read out two custom lists of 45 words with a time limit of one minute
264 per list. List A contained potentially familiar or trained items (words that occurred in the game) and
265 list B untrained items (words that did not occur in the game or in any other assessment). Words were
266 selected from a vocabulary achievement list of six-year olds (Schaerlaekens et al., 1999) and
267 consisted of monosyllabic words ranging from two to five letters (mean and median length of 3.5 and
268 four letters respectively) with a frequency range of 0.3 to 36608 per million (mean and median
269 frequency of 1612 and 51 per million respectively). Based on results of 272 children, both lists
270 correlated strongly at $r = 0.93$, split-half reliability with Spearman-Brown correction was also very
271 high at 0.96, as was Cronbach's α at 0.96. Because children's performance did not statistically differ
272 between lists in any of our analyses, we took the average of both lists and z -transformed the result.

273

274 **Phonological Awareness 1:** All phonological awareness subtests of the CELF-4-NL (Kort,
275 Schittekatte, & Compaan, 2008) test battery were administered, including blending phonemes into
276 words, identification of final and middle phonemes in words, sentence segmentation (by clapping
277 words), final syllable deletion, word segmentation (by clapping syllables), syllable deletion of bi-
278 and trisyllabic words, and initial phoneme substitution. Reliability of this test as measured by
279 stratified α is very high at 0.91 (van den Bos & Lutje Spelberg, 2010) as is internal consistency
280 measured by Cronbach's α at 0.94 (D'hondt et al., 2008). For analyses, we used both z -transformed
281 raw scores as well as norm scores, acting both as dependent and independent variables.

Deleted: ¶

Formatted: None

Commented [GM37]: Check PeerJ formatting for all sections - particularly this one.

Commented [GM38]: Are these outcomes? And if so, are they primary outcomes or secondary outcomes? Again, see CONSORT.

Commented [GM39]: This is reading like a list rather than a typical Outcome section.

283

284 **Phonological Awareness 2:** All phonological awareness tasks of the Proef Fonologisch Bewustzijn
285 (PFB, Elen, 2006) were presented, including rhyming, word segmentation (with the number of
286 syllables being indicated by clapping), blending of phonemes, syllables or lexemes into a word, and
287 pseudoword repetition. No reliability measures are provided for this test by the author. We analysed
288 both, z-transformed raw scores as well as norm scores as dependent and independent variables.

289

290 To measure the children's progress in terms of accuracy and response times with tasks that we had
291 incorporated within the game itself, the following game-based tests are additional outcome measures.

292 A detailed description of these following tasks with screenshots can be found in Appendix 1.

293

294 **Timed letter-speech-sound identification:** Children heard a phoneme and had to select the
295 corresponding grapheme with a computer mouse on the screen as fast as they could. Simple and
296 complex graphemes were presented one by one with five to 10 distractors per trial. We tested 32
297 different graphemes in 42 trials distributed across four levels, each with a time limit of one minute.
298 The time limit meant that only the fastest children saw all 42 targets, while slower children have only
299 been able to see a fraction of that. Based on results of 270 children, internal consistency was high as
300 indicated by Cronbach α (0.87 and 0.93) as well as split-half reliability with Spearman-Brown
301 correction (0.87 and 0.91) for level and item analyses respectively. For data analyses, we considered
302 both, single trial accuracy (binary correct/incorrect) and response times as dependent variables, as
303 well as the absolute number of correctly named letters within four minutes as a covariate.

304

305 **Written lexical decision:** Children saw a word or pseudoword on screen and had to either accept it
306 as a real word or reject it as a pseudoword by clicking on a green checkbox or a red cross. This task
307 contained 16 words and 16 pseudowords and was split up into two levels of 16 items, each with a

308 three-minute time limit. For data analyses, we used single trial measures, i.e. considering accuracy
309 and response times for each target. Similar to the reading fluency task, monosyllabic words with two
310 to four characters (mean and median length 3.1 and three letters respectively) and a frequency range
311 of four to 24266 per million (mean and median frequencies of 2546 and 124 per million respectively)
312 were used. The pseudowords were created based on those 16 words with Wuggy (Keuleers, &
313 Brysbaert, 2010), in most cases with an edit distance of one grapheme (e.g. jas/jal or tijd/toed). The
314 pseudowords were therefore balanced in length and also featured a high neighbourhood density of
315 real words at an edit distance of one grapheme, ranging from eight to 36 neighbours (with a mean
316 and median of 21 neighbours). Based on results of 199 children, internal consistency was
317 questionable as indicated by Cronbach α (0.7 and 0.68) as well as split-half reliability with
318 Spearman-Brown correction (0.7 and 0.65) for level and item analyses respectively.

319
320 Finally, the following tests were administered as co-variates for the analyses:

321 **Abstract reasoning:** The analogies and categories subtests of the SON-R 6-40 (Tellegen & Laros,
322 2014) were used as an estimate of nonverbal fluid intelligence. Within the analogies subtest children
323 have to identify a pattern that changes one geometrical figure into another and apply this principle to
324 a new figure. The categories subtest presents three pictures with a common characteristic and
325 children have to pick two additional pictures (out of five), which also possess this characteristic.
326 Reliability of this test is generally high ranging from 0.87 to 0.95. Because norm scores are only
327 available for children aged six and older and we had a substantial number of children under the age
328 of six in our sample, the resulting raw scores of both subtests were averaged and z -transformed.

329
330 **Rapid Automatized Naming (objects and colours):** The test requires participants to name out loud
331 50 depicted objects and colours in five rows of 10 items as accurately and as quickly as possible (van
332 den Bos, 2003). We noted the time (in seconds) it took to name the entire list of 50 items. The

Commented [GM40]: You need to justify why you used these covariates. Since you do not discuss this kind of things in the Introduction, you may need to include this information at the end of the Introduction under the Aims section. Your choice will depend on the extra information you provide in the Introduction about previous studies of GG. You may find you can explain the need for co-variates as part of that information. Either way, it needs to go somewhere.

333 reliability of these subtests as indicated by stratified α is in the range of 0.89 to 0.91 (van den Bos &
334 Lutje Spelberg, 2010).

335

336 **Training properties:** We extracted six variables related to game progress and learning opportunity:
337 the number of played sessions, hours and levels, as well as the overall number of seen items and
338 given responses, and the maximum level that was reached at the end of the training. While the
339 reading and math game are both based on the same framework and mini-games, due to the
340 differences in cognitive load, the math game generally features less distractors and shorter levels
341 compared to the reading game. Even though both groups had the same exposure in terms of sessions
342 and hours on task, children in the math group were exposed to more levels, gave more responses, and
343 were overall exposed to less items on screen. Furthermore, the number of levels differed in both
344 games (265 in the reading and 178 in the math game).

345

346 **Procedure**

347 Pre-tests (T1) commenced in September, three to six weeks after the start of the new school term,
348 followed by a five to seven-week playing phase in October and November, with post-testing (T2)
349 being conducted in November and December. All tests mentioned above were administered twice,
350 except for abstract reasoning (T1 only), reading fluency (T2 only), and the in-game written lexical-
351 decision task (T2 only). In addition to parents giving written informed consent prior to the start of the
352 study, all children were asked for verbal consent before the assessments started. The behavioural
353 tests took up to an hour each and were administered by undergraduate and graduate students in
354 speech-language pathology or linguistics during school time.

355

356 The children in the two experimental groups played the respective games (reading or math) for 10 to
357 15 minutes every day during school hours, resulting in five to 10 minutes of effective playing time

Commented [GM41]: Is this a co-variate? If not, what section should it fall under?

Commented [GM42]: Some of this information needs to be shifted to a Design section and some to a Interventions section (see CONSORT)

on task. Children played individually on a computer or laptop wearing headphones. The supervision of the training sessions was carried out by the teachers and differed among schools depending on the numbers of computers, the curriculum, and other local circumstances. At some schools, all the children in a classroom played the respective games at the same time in a computer room, whereas at other locations children had to take turns using five to ten computers during classes. To ensure that the children understood the tasks and enjoyed playing the games, the teachers and student assistants asked them about their progress and encouraged them to give the game another try when the content became more difficult at least once a week. Whilst the intervention groups completed the game-based assessments during the first and last playing session, the passive controls did so at some point during September or October (T1) and November or December (T2).

368

369 Data analysis

370 Differences at T1 were checked using two-way ANOVAs including experimental condition (reading, 371 math, passive control) and country (Netherlands, Belgium) as predictors. Significant main effects or 372 interactions were then followed up with a *t*-test or Tukey HSD. We found several main effects of 373 country with the Dutch children consistently outperforming the Belgian children in letter knowledge, 374 phonological awareness and rapid automatized naming (for details see Results section). This finding 375 was unexpected and made our planned analyses obsolete. We thus decided to split up the analysis of 376 this study into three parts: i) the Belgian sample, to evaluate research questions one (assessment 377 tools) and two (population characteristics) in beginning readers with active and passive control 378 groups, ii) the Dutch sample, to evaluate research questions one and two in a population of more 379 advanced readers limited to an active control group, and iii) a combined sample of the reading and 380 math groups from both countries, to specifically evaluate research question three (intervention 381 properties) for a broad range of reading abilities.

Commented [GM43]: At this point, the reader has no idea that there was a planned analysis and so this is confusing.

In addition, it means that having an analysis section under the Methods is inappropriate.

Commented [GM44]: REFORMAT according to PeerJ lists format

Further statistical analyses were carried out using linear mixed effects regression in R 4.0.4 (R Core Team, 2021; Bates, Maechler, Bolker, & Walker, 2015) and non-linear mixed regression using generalized additive models (Wood, 2006). Separate models were fit for dependent variables of reading fluency, phonological awareness, letter speech sound identification (accuracy & response time) and written lexical decision (accuracy and response time) separately for both countries. Therefore, to answer our three research questions, a total of 16 mixed models were built. Due to the high number of models and effects, we will focus on those models which show effects relating to experimental conditions and research questions (see Table 3 for an overview of all models). To identify the best model based on the Akaike's Information Criterion (AIC; Akaike, 1974), we tested the stepwise inclusion of main effects and interactions for fixed effects (fitted with maximum likelihood), as well as random intercepts and slopes for the random effects structure (fitted with restricted maximum likelihood estimation). A predictor (or more generally an effect) was only kept if it reduced AIC by at least two, thus indicating a better model fit while penalizing the increase in complexity of the model. In case the best model ended up without covering our research questions, i.e. if gaming condition was not a term in the best model, we could infer that there is no effect of gaming condition. Regardless, driven by our research questions, in these cases we still added gaming condition as a main effect to the best model to be able to report measures of significance and effect size.

The following variables were considered during the model building: age at T1, gender, gaming condition, hours played, handedness, familial risk for dyslexia, language spoken at home (mono/multilingual), abstract reasoning, letter knowledge, log transformed rapid naming speed of colours and objects, and phonological awareness. Where available, we always tested for inclusion of raw and percentile/norm scores as predictor (e.g. for CELF phonological awareness we had both raw scores and percentiles, tested the inclusion of both, and picked the one that explained more variance). For in-game measures we also tested inclusion of previous trial response time, current trial response

Commented [GM45]: Sometimes you have gaps between paragraphs and sometimes not. I have tried to fix missing gaps but there are too many. Please get an understanding of the PeerJ format guidelines and stick to those throughout the manuscript

407 time (in case of accuracy models) and trial number as predictors to remove autocorrelation of
 408 observations. Due to their highly skewed nature, response times were always box-cox transformed
 409 (Sakia, 1992). For analyses of game exposure and learning opportunity we furthermore considered
 410 the variables of played sessions, hours and levels, as well as numbers of given responses and items
 411 seen over gameplay and the highest level that was reached at the end of the training. Apart from
 412 these fixed effects, we added random intercepts and slopes for variables such as items, subjects,
 413 classrooms, and schools. To facilitate interpretation, raw scores were centred and z -transformed
 414 where possible, so that the model coefficient β is identical to the effect size Cohen's d .
 415 Finally, for each resulting model we trimmed outliers based on residuals beyond ± 2 standard
 416 deviations of the model prediction and refitted the model to ensure that presented effects are not
 417 carried by outliers. Every model then underwent a model criticism to ensure that reported models
 418 fulfil test assumptions of independence of observations as well as a normal and homoscedastic
 419 distribution of residuals. Usually, model fit is evaluated by the squared correlation between the
 420 observed and the fitted values (R^2). For mixed-effects models, this method can only estimate the
 421 residual variance and thus ignores the random effects present in the model. Following the approach
 422 proposed by Nakagawa and Schielzeth (2013), a marginal and conditional R^2 was calculated. The
 423 former being an estimation of the fixed-effects structure alone, while the latter incorporates both
 424 fixed and random effects.
 425 Due to data trimming and cases of missing observations, most analyses were carried out on smaller
 426 subsets of the data, and exact sample sizes are reported at the corresponding positions of the results
 427 section. Most notably, because of data retrieval problems, the results of the in-game assessments at
 428 the post-test are not available for 66 children ($N_{\text{Read}} = 31$, $N_{\text{Math}} = 33$, $N_{\text{Passive}} = 2$).

429 Results

430 At T1 the Dutch children significantly outperformed their Belgian peers in terms of abstract
431 reasoning ($F_{1,242} = 15.74, p < .001$), letter knowledge ($F_{1,241} = 288.84, p < .001$), both phonological
432 awareness tests (CELF: $F_{1,238} = 59.68, p < .001$; PROEF: $F_{1,242} = 37.81, p < .001$) and both rapid
433 automatized naming measures (colours: $F_{1,242} = 31.40, p < .001$; objects: $F_{1,241} = 16.58, p < .001$; see
434 Table 2). For this reason, separate analyses were carried out within each country. For the Belgian
435 children there was a main effect of condition for rapid automatized naming colours at T1
436 ($F_{2,158} = 3.06, p = .050$) where the math group was significantly faster than the reading group (post-
437 hoc with TukeyHSD: $p = .039$). There was also a main effect of condition for letter knowledge
438 ($F_{2,157} = 3.22, p < .043$) where the passive group knew more letters than the reading group (post-hoc
439 with TukeyHSD: $p = .035$). For the Dutch sample, the math group had significantly more
440 multilingual children than the reading group (Fisher's exact test: $p = .018$) but otherwise the groups
441 did not differ in any other measures at T1.

442

443 Research Question 1: Assessment tools

444 To answer the first research question, we evaluated word reading fluency, phonological awareness
445 and the two in-game tests of letter-speech sound identification and written lexical decision. **Word**
446 **reading fluency** was assessed with two one-minute reading lists at T2. Whereas we did not find any
447 effects associated with gaming condition in the Dutch sample, there were effects in the Belgian
448 group (see Figure 2). Neither the reading ($\beta = 0.27, t = 1.60, p = .09$) nor the passive ($\beta = -0.27,$
449 $t = -1.63, p = .106$) group differed from the math group, but the reading group outperformed the
450 passive group ($\beta = 0.55, t = 3.18, p = .002$). In terms of effect sizes, these differences were small
451 ($d = 0.27$) for the passive and reading group compared to the math group, and medium-sized
452 ($d = 0.55$) when comparing the reading group to the passive group.

453 This best model was based on 150 Belgian children ($N_{\text{Passive}} = 48$, $N_{\text{Math}} = 52$, $N_{\text{Read}} = 50$, trimmed
 454 seven observations or 4.3% of data), controlled for the covariates letter knowledge at T1, CELF
 455 phonological awareness at T1, log transformed rapid automatized naming colours time at T1 and
 456 included random intercepts per school. The model had a conditional R^2 of 0.39 and a marginal R^2 of
 457 0.30.

458 **Phonological awareness** was measured using the nine subtests of the CELF and four subtests of the
 459 PFB. Neither for the Belgian nor for the Dutch sample did we find any effects related to gaming
 460 condition in either of the two tests.

461 **Letter-speech sound identification** assessment was embedded into the game itself. We found that
 462 the reading game boosted accuracy in the Belgian sample and a trend towards boosting response
 463 speed in the more advanced Dutch sample. For the Belgian sample, the best model predicting single
 464 trial accuracy at both testing sessions for 6756 trials of 101 children ($N_{\text{Passive}} = 47$, $N_{\text{Math}} = 21$,
 465 $N_{\text{Read}} = 33$) revealed an interaction of time \times condition (see Figure 3). At T1 the control group knew
 466 significantly more letters than the math ($\beta = 0.53$, $z = 2.37$, $p = .018$) and reading groups ($\beta = 0.51$,
 467 $z = 2.60$, $p = .009$). While we did find an effect of time for the math group ($\beta = 1.78$, $z = 12.69$,
 468 $p < .001$) this was smaller for the passive group ($\beta = -0.49$, $z = -3.02$, $p = .003$) and marginally bigger
 469 for the reading group ($\beta = 0.34$, $z = 1.88$, $p = .060$). The gain of the reading group therefore also far
 470 exceeded the passive group ($\beta = 0.83$, $z = 5.79$, $p < .001$). This best fitting model was controlling for
 471 level type, CELF phonological awareness at T1 and trial response time. The random effect structure
 472 consisted of random intercepts per subject and target, and random slopes for previous response time
 473 and CELF phonological awareness at T1 by subject. The model had a conditional R^2 of 0.47 and a
 474 marginal R^2 of 0.20.

475 For the Dutch sample the best model, which predicted box-cox transformed single trial response
 476 times based on 3646 trials of 75 children ($N_{\text{Math}} = 48$, $N_{\text{Read}} = 27$, trimmed 212 trials or 4.9% of data),
 477 also revealed a time \times condition interaction (see Figure 4). At T1 the two groups did not differ from

one another ($\beta = 0.01, t = 0.59, p = .597$) but we found a significant main effect of time ($\beta = 0.04, t = 7.76, p < .001$) and a marginally significant interaction of the two ($\beta = 0.01, t = 1.98, p = .052$) with the effect of speeding up from T1 to T2 being bigger for the reading than the math cohort. This best model controlled for PROEF phonological awareness at T1, age at T1, trial number and previous trial response time. The random effects structure consisted of intercepts per subject, class, target, and distractor order on screen, as well as random slopes for time by subject and random slopes for time by target.

Written Lexical Decision assessment was embedded into the game itself and did not show any differences between the gaming conditions in terms of accuracy or response times.

Research Question 2: Population characteristics

To answer the second research question, we were looking for possible interactions of gaming condition with pre-test scores, as well as age, gender, familial risk, abstract reasoning and home language environment. In almost all analyses, population characteristics explained unique variance as covariates and thus helped to describe more robust and generalizable intervention effects. However, we did not find any interactions of population characteristics and intervention type, which suggests that improvements are comparable across subpopulations. **Familial risk** for dyslexia was assessed by parental questionnaires inquiring about the occurrence of reading difficulties in first grade relatives. Familial risk was a relevant predictor for PROEF phonological awareness scores in the Belgian sample, reflected in slightly lower scores for children with a familial risk for dyslexia ($\beta = -0.23, t = -1.34, p = .183$) with a small effect size ($d = 0.23$). Otherwise, we found no evidence that familial risk for dyslexia had an effect on pre-test scores or response to intervention. **Age** was a relevant covariate in analyses of in-game response times of letter-speech sound identification ($\beta = 0.02, t = 2.22, p = .030$) in the Dutch sample and written lexical decision tasks ($\beta = -0.20, t = -3.17, p = .002$) in the Belgian sample. In both cases, younger children took, on average, longer to

503 reply than older children. No effect sizes for these in-game assessments could be computed from the
504 mixed models. **Gender** was a relevant covariate for letter-speech sound identification accuracy in the
505 Dutch sample ($\beta = -1.62, t = -4.59, p < .001$) and word reading fluency in the combined sample
506 ($F = 8.775, p < .001$). In both cases, girls outperformed their male peers by a significant margin.
507 Again, no effect sizes for these in-game assessments could be computed. **Nonverbal intelligence** as
508 measured by the SON-R 6-40 was a relevant co-variate for the CELF phonological awareness in the
509 Belgian sample ($\beta = 0.17, t = 2.65, p = .009$) with higher nonverbal intelligence resulting in higher
510 phonological awareness scores and a small effect size ($d = 0.17$). **Home language environment** and
511 **handedness** never came up as relevant predictors.

512

513 **Research Question 3: Intervention properties**

514 To answer the third research question, we were looking for possible interactions of gaming condition
515 with exposure measures. The inclusion of six game progress and achievement measures was tested in
516 all models fitted for research questions one and two, but none of them turned out being relevant,
517 suggesting that response to intervention is independent of training properties. One possible
518 explanation for this could be the exclusion of children who did not reach 20 playing sessions, which
519 reduces variance in exposure. We therefore re-included children who played less than 20 sessions
520 ($N = 15$), and combined children from both countries to increase statistical power, putting the
521 available sample size for this analysis at $N = 210$. The resulting groups did not differ in any test
522 scores at T1 (see Table 4 for sample characteristics).

523 The best model predicting z -transformed one-minute reading fluency scores at T2 for 196 children
524 ($N_{\text{Math}} = 95, N_{\text{Read}} = 101$, trimmed 14 observations or 6.7% of the data) described two nonlinear
525 interaction surfaces of CELF phonological awareness at T1 and the first principal component of
526 exposure per group. For the reading group this nonlinear interaction was significant ($F = 2.99$,
527 $p = .009$) while for the math group it was not ($F = 0.20, p = .653$; see Figure 5). Within the math

group (subplot A) there is an almost linear relation between phonological awareness skills at T1 and reading fluency as indicated by the vertical and equidistant topographic lines, whereas for the reading group (subplot B) there is a nonlinear interaction of these two variables. When phonological awareness is kept stable (e.g. at -1 or +1 z-scores) exposure modulates reading outcome, but only when exposure is above average. The difference between these two surfaces (subplot C) indicates that children with good phonological awareness skills and a lot of exposure to the reading game are more proficient readers with a medium-sized effect of Cohen's $d = 0.5$ than their peers from the math group. With additional main effects for country, as well as two nonlinear smooths for rapid automatized naming colours time at T2 and CELF phonological awareness at T1, this model explained 49.8% of variance in the reading fluency scores.

To investigate whether this effect within the reading group was carried by specific subpopulations we included gender and familial risk for dyslexia as covariates. This led to a similar model ($N = 98$, $N_{\text{male}} = 51$, $N_{\text{female}} = 47$, trimmed six observations or 5.8% of data) with two nonlinear interaction surfaces of CELF phonological awareness at T1 and game exposure for males ($F = 6.27$, $p < .001$) and females ($F = 8.28$, $p < .001$). Upon visualization (see Figure 5) it appears that the pattern of girls mirrors that seen for the reading game in general (subplot B vs. E). The difference between boys and girls who played the reading game, which was also significant ($F = 6.32$, $p < .001$, subplot F), shows a somewhat diffuse pattern without a clear interpretation. This model further included a nonlinear smooth for CELF phonological awareness at T1 and explained 80.1% of variance in reading fluency scores. Notably, the effects described above are only carried by a handful of children each and the model split by status of familial risk of dyslexia did not reveal additional effects, probably because of the limited number of children at familial risk for dyslexia in the current sample.

550 **Discussion**

551 In this study, we evaluated the effectiveness of a newly created version of GraphoGame for Dutch-
552 speaking beginning readers, employing active (math game) and passive (no game) control conditions
553 in 16 first-grade classrooms in the Netherlands and Flanders, the Dutch-speaking part of Belgium.
554 The main purpose of this game was to intensify exposure to relevant early reading materials and to
555 provide additional training for struggling beginning readers. We found large differences between the
556 two countries at both testing points irrespective of training, which led us to conduct separate analyses
557 within each country. In the Belgian sample, children who played the literacy game improved their
558 letter knowledge more than the other two groups (as measured by the accuracy in the timed letter-
559 speech sound identification task) and we observed somewhat faster reading fluency in this group at
560 post-test compared to the other two conditions with small to medium-sized effects. For the overall
561 more advanced Dutch sample, there was a trend towards faster responses in timed letter knowledge
562 for the reading group compared to the math group. Recombining both samples revealed that children
563 who played extensively and scored high on phonological awareness prior to training were more
564 fluent readers than could be expected based on other reading precursors and based on phonological
565 awareness alone. Beyond an overall evaluation of effectiveness, our secondary aim was to conduct
566 an exploration of further characteristics, possibly determining the effectiveness of digital game-based
567 learning in early literacy training within the framework of a single study. Thus, three factors
568 potentially contributing to the effectiveness of digital game-based learning were investigated: i)
569 assessment tools, ii) population characteristics, and iii) intervention properties.

570

571 **Research Question 1: Assessment tools**

572 We asked whether the evaluation of response to intervention partly depends on the choice of
573 assessment tools used to track changes in performance. We measured reading-related skills as well as
574 reading itself, by using different assessments for the same skill, and also combining online and

575 offline measures. To evaluate reading-related skills studies typically use assessments like word (e.g.,
576 EMT; Brus & Voeten, 1991) and pseudo word reading tests (e.g., De Klepel; van den Bos, Spelberg,
577 Scheepsma, & De Vries, 1994), word dictation (e.g. PI-dictee; Geelhoed & Reitsma, 1999) or tests
578 for vocabulary, phonological awareness, or rapid automatized naming (e.g. CELF-4-NL; Kort et al.,
579 2008). However, these paper-and-pencil tests are usually not designed to detect the subtle changes in
580 performance occurring during a few weeks of learning. If an improvement is not big enough, it may
581 simply get lost within variance related to sensitivity, specificity or test-retest reliability of a given
582 test. For example, the CELF-4-NL manual states reliabilities in the range of 0.71 to 0.86 for
583 subscales, and 0.88 to 0.92 for composite scores as well as a test-retest reliability over a 5-month
584 period of 0.75. These lead to a lack of **sensitivity to change**, which may be one of the main reasons
585 why shorter training studies fail to find significant improvements in reading-related skills even
586 though the games trained exactly these skills.

587 Two more issues need to be discussed in this context. Firstly, poor sensitivity to change could be
588 seen as a question of learning transfer: measuring (timed) letter knowledge before and after a training
589 of grapheme-phoneme correspondence can be considered a near training transfer, whereas evaluating
590 changes in reading fluency based on a combined training of grapheme-phoneme correspondences
591 and phonological awareness can be considered a far transfer, which may take longer and be overall
592 smaller (Froyen, Bonte, Atteveldt, & Blomert, 2009; Vaessen & Blomert, 2010). And indeed,
593 improvements in letter knowledge are almost unanimously reported in literacy digital game-based
594 learning research (Richardson & Lyytinen, 2014) as this skill is easily trainable and measurable,
595 while improvements in phonological awareness and reading fluency are rather the exception (e.g.
596 studies reported in McTigue et al., 2019; Carvalhais, Limpo, Richardson, & Castro, 2020; Lovio,
597 Halttunen, Lyytinen, Nääätänen, & Kujala, 2012; Ktisti, 2015). To assess phonological awareness in
598 the present study two paper-and-pencil tests were used which differed in nature. Whereas the CELF-
599 IV (Kort et al., 2008) tests nine different abilities one by one (e.g. segmentation, blending, phoneme

600 identification, deletion and replacement), the PROEF (Elen, 2006) constantly alternates between one
601 of four abilities (rime, segmentation, blending, pseudoword repetition). Our assumption was that the
602 PROEF would better reflect the automatization of phonological processing, but we did not find any
603 training effects with either test. While not being a far learning transfer per se, training of
604 phonological skills and/or measuring potential improvements appear to be difficult to achieve for
605 most short-term interventions.

606 Secondly, poor sensitivity to change is also a matter of how skills are assessed. For letter knowledge,
607 paper-and-pencil tests are typically administered without time pressure and reach ceiling within the
608 first few months of school (Blomert & Willems, 2010), thus losing predictive and evaluative power.
609 On top of accuracy, one can also measure response times and add time pressure, which within a
610 speeded letter-sound association task is even more specifically related to the fluency of multimodal
611 processing of audio-visual information (Blomert, 2011; Hahn et al., 2014). Whereas letter knowledge
612 just assesses the availability of letter-sound associations, letter-speech sound identification (here also
613 called 'timed' letter knowledge) additionally assesses the fluency with which these associations can
614 be retrieved from memory. As the game trains these exact skills, it can also be used to evaluate the
615 response to intervention. For this reason, we incorporated separate test units within the game to be
616 played at the start and the end of the training period to measure progress in reading-related skills.
617 These in-game assessments measure accuracy and response times at the item level, and thereby tap
618 into the domain of automatization to an extent which offline paper-and-pencil tests are not able to
619 capture.

620 However, collection of single trial data alone is not sufficient, the data should also be analysed as
621 such. In our case, conducting an ANOVA on aggregate data (i.e. count of correctly named letters
622 within four minutes) did not yield any group differences, whereas the use of mixed effects regression
623 of the single trial data did show that children who played the literacy game made more pronounced
624 progress in letter knowledge than their peers who played the math game or who did not play any

game. Notably, this finding was limited to the Belgian sample and not quite statistically significant between reading and math groups. The more advanced Dutch sample was already close to ceiling at pre-test, preventing gains in terms of accuracy. But there we also observed a trend towards faster responses in the reading group, stressing the importance of considering response times as an indicator of automatization once accuracy reaches ceiling. Although it can be argued that most children would eventually have attained the same skill levels without the game, our findings confirm that GraphoGame-NL can speed up acquisition of, and access to, grapheme-phoneme correspondences for children of all abilities. Moving on to another in-game measure, our attempt of assessing reading fluency in form of a written lexical decision task inside the game was not successful. The reliability measures were poor and there were only weak associations with other measured variables. The task was perhaps too difficult for starting readers and might only become a relevant measure for reading fluency at a later stage.

Regardless, our study demonstrates that the evaluation of game effectiveness depends on the choice of assessment tool and the statistical analysis. There are benefits of using in-game measures and we provide some starting points for future research. The use of in-game assessments does not remain without methodological issues though (Puolakanaho & Latvala, 2017). When children do not understand the task or find it too difficult or boring, they may just randomly click around to pass the level, achieving very fast response times but correspondingly low accuracy scores. This is easily controlled for in group-wide analyses by excluding the data of children performing below chance level, although it is often difficult to set a chance level because weighting the complexity and confusability of the many items presented simultaneously may pose a problem (Kujala, Richardson, & Lyytinen, 2010). To obtain useable data, in our experience it is essential to have an adult supervise the assessments in small groups of children. If performed individually and unsupervised, the assessments may generate little useable data because the tests may not have been performed as intended. Possibly, assessments can be repeated at certain intervals, but ultimately, it would be most

650 convenient to collect dynamic in-game data that considers the entire gameplay by continuously
651 tracking a child's progress, possibly even precluding the need for dedicated assessment levels.

652

653 **Research Question 2: Population characteristics**

654 Our question was whether population characteristics impact intervention effectiveness. We
655 hypothesized that poor performers would benefit most from GraphoGame, but we did not find any
656 evidence for that. Most effects relating to gaming condition were main effects, which indicates that
657 there were no systematic differences within the three experimental groups. The only exception, albeit
658 pointing the opposite direction, being the few children who performed above average in phonological
659 awareness skills at pre-test, who were comparatively faster readers conditional to more extensive
660 exposure to the reading game. We also anticipated that certain subgroups of children (like those at
661 familial risk or those speaking a different language at home) might perform worse at pre-test and also
662 exhibit a different outcome from exposure to the game, but we did not find evidence for that either.
663 Literacy interventions usually target poor performers, who are a generic group of children in whom
664 the underlying mechanism of reading-related difficulties may vary drastically. To account for this
665 variability, both reading-related performance and the presence/absence of familial risk for dyslexia
666 need to be taken into account as such children are more likely to share a common underlying deficit
667 accounting for their reading deficiencies (Snowling & Melby-Lervag, 2016; van Viersen, de Bree,
668 Zee, Maassen, van der Leij, de Jong, 2018). The questions whether poor performers and children at
669 familial risk of dyslexia can profit from digital game-based learning training, and whether or how
670 these groups differ from each other are of high clinical relevance.
671 Unfortunately, such a question remains difficult to answer if **inclusion criteria** vary across studies
672 and only certain children take part. Most studies use an inclusion criterion based on scores in
673 reading-related tests (e.g. Saine et al., 2010, 2011), the nomination by class teachers (e.g. Kyle et al.,
674 2013) or socioeconomic status (SES; e.g. Rosas, Escobar, Ramírez, Meneses, & Guajardo, 2017).

675 While the rationale for such inclusion criteria is clear, all these approaches pose certain difficulties.
676 In case of the test-based or SES-based approach, there is the question of finding the right cut-off
677 score. Furthermore, due to **regression to the mean**, children scoring at the lower end of the
678 population scale are more likely to perform closer to average at the next assessment (Morten &
679 Torgerson, 2004). On the other hand, teacher ratings may be subjective and based on the assessment
680 of skills unrelated to a child's reading abilities (Begeny, Krouse, Brown, & Mann, 2011).
681 To prevent such sampling bias, in the **present study** we invited all children from 16 classrooms to
682 play, independent of their performance on reading-related tasks and investigated the effect of pre-test
683 scores on training-induced skill improvement. Our approach was unintentionally strengthened further
684 because of the drastic pre-test differences between the Dutch and Belgian children in our sample.
685 These differences appear to stem from the different preschool systems, where Belgium has a stricter
686 separation of pre-school and school, often requiring a physical change of school around the age of
687 six, the Netherlands has a more gradual transition into formal instruction from four years of age
688 onwards within the same institution. This interpretation is also supported by the fact that similar
689 differences between these two neighbouring countries have been observed in early numeracy skills
690 (Torbeyns, Van den Noortgate, Ghesquière, Verschaffel, Van de Rijt, & Van Luit, 2002). Ultimately,
691 this gave even further spread to the cognitive measures in our sample and allowed us to evaluate the
692 impact of factors such as age, abstract reasoning, familial risk for dyslexia, gender,
693 language(s)/dialects spoken at home, and handedness more exhaustively than has been done in
694 previous literacy digital game-based learning research.
695 At first sight one could argue that, due to the absence of interactions of **pre-test scores** and outcome,
696 the intervention was equally effective for all children. However, when comparing results stratified by
697 country, it is apparent that the much weaker beginning readers in Belgium showed overall more
698 intervention effects (in letter knowledge and reading fluency), whereas for the more advanced Dutch
699 sample we found fewer effects (limited to grapheme-phoneme correspondence automation). This can

700 be taken as evidence that individual starting levels matter for intervention outcomes, which is in line
701 with most previous studies. Training poor performers at an early stage in their literacy development
702 usually yields group-wide benefits in easily trainable skills like letter knowledge (e.g. Brem et al.,
703 2010; Rosas et al., 2017), and in longer interventions also decoding and reading (e.g. Saine et al.,
704 2010; 2011). The opposite effect, that children with high pre-test scores have an increased benefit,
705 has also been reported before. Ruiz et al. (2017) found a small but significant advantage of early
706 readers who already scored high at pre-test in timed letter knowledge. The few studies who trained
707 entire classrooms (e.g. Jere-Folotiya et al., 2014; Koikkalainen 2015; Ronimus & Lyytinen, 2015)
708 did unfortunately not consider interaction terms in their analyses, thus providing no reference point
709 for comparisons. Regarding the general role of pre-test scores as predictors for intervention
710 outcomes, conventional reading interventions found that reading-related skills are actually poor
711 predictors for the response to intervention. Improvements were rather related to levels of short-term
712 memory and vocabulary (Byrne, Shankweiler, & Hine, 2008) - two variables which were not
713 measured in the present study and are not routinely collected and used to control for confounding in
714 analyses of reading interventions.

715 For effects relating to **familial risk** of dyslexia, we found that at-risk children had slightly lower
716 phonological skills, but the training effectiveness was not influenced by status of familial risk. The
717 former is somewhat surprising, given that other studies also reported weaker performance in other
718 reading precursors for children at familial risk (van Bergen et al., 2012; Lyytinen et al., 2004). So far
719 only two studies have specifically investigated the role of familial risk in GraphoGame effectiveness.
720 Whereas a study by Brem and colleagues (2010) did not find any distinct effects relating to familial
721 risk either, a study by Blomert and Willems (2010) found that risk children did not improve as much
722 as their peers. The authors concluded that familial risk is characterized by a letter-speech-sound
723 association and integration deficit, which the data from the present study does not support. Both
724 studies had shortcomings preventing the authors from drawing firm conclusions about the effects of

725 familial risk on GraphoGame effectiveness which merit mentioning. Including as few as 32 children
726 (14 risk, 18 no risk) across two experimental groups, the study by Brem and colleagues (2010) may
727 have suffered from a lack of power. In addition, playing GraphoGame followed by a math control
728 game or vice versa in a crossover design, the children spent systematically less time on the second
729 game. Blomert and Willems (2010) suggested that the absence of improvements in timed letter
730 knowledge, phonological awareness and reading skill in the risk children in their study might have
731 been due to the young (preschool) age of these (familial risk) children being exposed to reading
732 materials that were too difficult. The fact that the present study did not find any distinct training
733 effects attributable to status of familial risk may be due to the small number of at-risk children in
734 each condition (varying from seven to 18) or the rather weak self-report questionnaire asking for
735 reading failure in the close family, but without requesting proof of a formal diagnosis in first grade
736 relatives.

737 An interesting insight from **gender** effects in game-based learning in general is that previous gaming
738 experience may predict in-game achievement, which puts girls at a disadvantage (Nietfeld, Shores, &
739 Hoffmann, 2014). Ideally, studies should therefore control for gender or previous game experience in
740 their analyses, which is currently almost never done in the field (e.g. for studies reported in McTigue
741 et al., 2019). While creators of game-based learning tools should aim to build gender neutral and
742 inclusive games, considering that developmental dyslexia is diagnosed 1.5 to three times more often
743 in boys than in girls (Rutter et al., 2004), a slight male preference for game-based learning might
744 actually be an asset. In our sample, boys had significantly poorer letter knowledge and phonological
745 awareness skills compared to girls at the start of first grade. This appears to be the onset of a constant
746 difference which extends throughout school into adolescence, where girls outperform their male
747 peers in terms of reading (OECD, 2010; Ming Chui & McBride-Chang, 2006; Torppa, Eklund, van
748 Bergen, & Lyytinen, 2015). We also found that the observed benefits in terms of reading fluency
749 when phonological awareness and game exposure were high was mostly carried by girls. Thus, at the

group level, the boys and girls in our sample were in slightly different stages of reading acquisition. In sum, we feel that gender differences warrant further scrutiny in literacy digital game-based learning research, also given that boys generally play more games, show a stronger preference for game-based learning and are more open towards technology and computers than their female peers (Admiraal, Huizenga, Heemskerk, Kuiper, Volman, & ten Dam, 2013; Gwee, San Chee, & Tan, 2011; Bonanno & Kommers, 2007).

Research Question 3: Intervention properties

Finally, we asked how intervention properties contribute to training effectiveness and we hypothesized that characteristics from the gaming process itself might help explain variance in the intervention outcome. Our study provides only limited insights in this regard. Previous literacy digital game-based learning studies using GraphoGame usually relied on the number of gaming sessions or the time spent playing as a measure of training intensity. Only few communicate treatment fidelity measures such as attrition rates, which can be as high as 46% (Jere-Folotiya et al., 2014). Studies reporting positive effects used training durations ranging from one up to 28 weeks with an intensity of two to five training sessions per week (McTigue et al., 2019; Richardson & Lyytinen, 2014). Whether training duration and intensity act as independent variables modifying digital game-based learning outcomes or whether the overall exposure to the game (in hours) is a better predictor of training effectiveness remains an open empirical question. Furthermore, the ideal training duration and intensity may differ depending on population properties and training goals, which raises the obligation to investigate possible interactions of training and population properties. We therefore extracted additional game-exposure measures, such as the highest level that was reached, or total number of seen items which might capture the actual gameplay better than mere time on task. For example, even though all children played in the range of 20-30 sessions, the number of items seen within the training period had a much wider range from 5000 to 20000. This is

a result of speed and accuracy of children: responding faster will yield more levels, responses and seen items, while being less accurate results in being exposed to less items during the same period. Due to the adaptivity of modern games which constantly adjust the difficulty level to the individual learner, different children are therefore exposed to different content, making exposure comparisons difficult, even within the same study. However, individual response patterns do also vary over time depending on the complexity (simpler, more familiar content vs. more complex new information) of consecutive levels (Nja, 2019). Individually, these additional measures did not seem to be related to response to intervention in the present study, but rather reflect pre-test skills. This confirms that data extracted from in-game behaviour can be used for dynamic assessment (Koikkalainen et al., 2015; Puolakanaho & Latvala, 2017).

Possibly, the rather strict inclusion criterion of at least 20 playing sessions made the present sample too homogenous to find interactions with exposure. Upon re-inclusion of children who played less than 20 sessions and by fusing these exposure measures with a principal component analysis, we found that learning opportunity and phonological awareness modulated reading fluency when other reading pre-cursors were kept stable. Therefore, the time-course of development of phonological skills plays a crucial role for the benefits of GraphoGame-NL. Playing beyond mastery of grapheme-phoneme correspondences has little impact on reading fluency when phonological skills are poor, and we did not find evidence that the current game promotes phonological skills at all. This is problematic, as combined letter-sound training and phonological awareness training were found to be more successful in boosting reading and spelling skills than either of them in isolation (Schneider, Roth, & Ennemoser, 2000). We therefore suggest reducing the weekly playing intensity once letter knowledge accuracy reaches ceiling, and instead extend the overall training period. This might allow poor performers to get more out of the game, especially to give more time for maturation of phonological skills. Future studies should furthermore focus on identifying how to best train phonological skills with digital game-based learning. To achieve this goal, a more qualitative

approach should be taken to analyse games, their training content and training properties. Currently, this is not possible because for most games the underlying intentions, decisions, settings and materials are not sufficiently described and/or shared. In Appendix 1 we provide a detailed description of the games used in this research, as we believe this to be crucial in enabling future research to uncover the mechanisms of (more) successful interventions.

805

806 **Limitations**

As with all studies, we acknowledge several limitations in the design and procedure, which should be considered when interpreting the results and analyses presented above. The unexpectedly large pre-test differences forced us to split our sample by country, which led to smaller groups and reduced power compared to the study we initially conceived. Due to significant group differences at pre-test, we cannot rule out regression to the mean as a possible explanation for some effects described above. The analyses presented here also tested the inclusion of a wide range of measures as covariates in a conservative, yet exploratory fashion. We highly recommend replication of our results with other cohorts of Dutch and Flemish children. An additional weakness is that we only measured reading fluency at post-test. Due to an earlier pilot showing floor results and due to time constraints for testing at schools we decided not to collect such data at pre-test. As a result, we could not directly test interactions between reading fluency improvement and other factors, but by controlling reading fluency outcome for reading precursors at pre-test (letter knowledge, phonological awareness, rapid automatized naming and age), these results are nevertheless relevant and meaningful. Another issue arises from the fact that the teachers who participated were favourable, or at least open, towards the use of digital tools in their classrooms, and were furthermore not blinded to the experimental conditions, and thus knew their treatment allocation. This may have changed their teaching style in one way or another, which is something that is hard to control or correct for. To balance out the impact single classrooms may have on intervention effects, children should ideally be randomized

825 individually, i.e. one third of a classroom playing the reading game, one third playing a control game
826 and one third not playing. From our experience this is hard to implement in classrooms and it would
827 also negatively affect classroom atmosphere if some children were not allowed to play. Another
828 alternative could be to implement the playing at home, which would come with its own set of
829 challenges like how to ensure daily playing or prevent excessively long gaming sessions (Ronimus &
830 Lyytinen, 2015).

831 Finally, the math game may not have been the best control condition. Through ERP data collected
832 from a subset of the children in the present study it became apparent that playing the math game
833 might also contribute to the development of phonological awareness skills (Glatz, 2018). As
834 arithmetic representations are also phonological in nature (De Smedt & Boets, 2010; De Smedt,
835 Taylor, Archibald, & Ansari, 2010) both games ultimately promote careful listening and fast access
836 to phonological representations. Future research on computerized literacy training should therefore
837 try to make use of an active control condition where the improvements of video gaming can be
838 expected in the visual or motor domain (like described by Green & Bavelier, 2003) rather than in
839 verbal and/or auditory learning.

840

841 **Conclusion**

842 We conducted one of the first literacy digital game-based learning studies relying on single-trial data
843 from in-game tasks to evaluate its effectiveness. Playing GraphoGame-NL led to an increase in
844 mastery of grapheme-phoneme correspondences and to small to medium sized effects in reading
845 fluency. Demographic characteristics such as familial risk of dyslexia or languages/dialects spoken at
846 home had little impact on response to intervention and additional research investigating larger groups
847 of children at familial risk of dyslexia is needed. Follow-up studies will need to evaluate the longer-
848 term effects of such a brief computer-assisted literacy training in first graders learning to read the
849 semi-transparent Dutch orthography. It is unclear whether our findings are generalizable to more
850 opaque (e.g. English) or more transparent orthographies (e.g. Finnish and Greek). Studies employing
851 GraphoGame in Dutch are ongoing, with a focus on struggling readers and an exploration of new
852 learning materials and tasks.

853 **Acknowledgements**

854 We are very grateful to all the children, parents, and teachers who participated in our research. We
855 are indebted to Iivo Kapanen for helping to bring GraphoGame-NL to life, Sabien van Dyke and
856 Vanessa Janssens for test administration in Belgium, and Anastasia Glushko for feedback on early
857 drafts of this manuscript.

858

859 **References**

- 860 Admiraal, W., Huizenga, J., Heemskerk, I., Kuiper, E., Volman, M., & ten Dam, G. (2014). Gender-
861 inclusive game-based learning in secondary education. *International Journal of Inclusive*
862 *Education*, 18(11), 1208-1218.
- 863 Akaike, H. (1974). A new look at the statistical model identification. *IEEE transactions on automatic*
864 *control*, 19(6), 716-723.
- 865 American Psychiatric Association. (2013). Diagnostic and statistical manual of mental disorders
866 (DSM-5®). American Psychiatric Pub.
- 867 Bates, D., Maechler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects Models
868 Using lme4. *Journal of Statistical Software*, 67(1), 1-48. doi: 10.18637/jss.v067.i01
- 869 Begeny, J. C., Krouse, H. E., Brown, K. G., & Mann, C. M. (2011). Teacher judgments of students'
870 reading abilities across a continuum of rating methods and achievement measures. *School*
871 *Psychology Review*, 40(1), 23-38.
- 872 van Bergen, E., de Jong, P. F., Plakas, A., Maassen, B., & van der Leij, A. (2012). Child and parental
873 literacy levels within families with a history of dyslexia. *Journal of Child Psychology and*
874 *Psychiatry*, 53(1), 28-36. doi: 10.1111/j.1469-7610.2011.02418.x
- 875 Bergmann, J., & Wimmer, H. (2008). A dual-route perspective on poor reading in a regular
876 orthography: Evidence from phonological and orthographic lexical decisions. *Cognitive*
877 *Neuropsychology*, 25, 653-676. doi: 10.1080/02643290802221404
- 878 Blomert, L. (2011). The neural signature of orthographic-phonological binding in successful and
879 failing reading development. *Neuroimage*, 57(3), 695-703.
- 880 Blomert, L., & Willems, G. (2010). Is there a causal link from a phonological awareness deficit to
881 reading failure in children at familial risk for dyslexia?. *Dyslexia*, 16(4), 300-317. doi:
882 10.1002/dys.405
- 883 Bonanno, P., & Kommers, P. (2007, July). Exploring the Influence of Group characteristics on
884 Interactions during Collaborative Gaming. In *IADIS International Conference: e-Learning*
885 *held in Lisbon, Portugal*.
- 886 Borleffs, E., Maassen, B. A. M., Lyytinen, H., & Zwarts, F. (2017). Measuring orthographic
887 transparency and morphological-syllabic complexity in alphabetic orthographies: A narrative
888 review. *Reading and Writing*, 1-22, doi: 10.1007/s11145-017-9741-5
- 889 van den Bos, K. P., Spelberg, H., Scheepsmma, A., & De Vries, J. (1994). De Klepel. Vorm A en B.
890 Een test voor de leesvaardigheid van pseudowoorden. Verantwoording, handleiding,
891 diagnostiek en behandeling. Berkhout, Nijmegen, The Netherlands.
- 892 van den Bos, K. P. (2003). Snel Serieel Benoemen; Experimentele versie. [Rapid naming;
893 Experimental version]. Groningen: University of Groningen.
- 894 van den Bos, K. P., & Lutje Spelberg, H. C. (2010). CB&WL Continu Benoemen & Woorden Lezen.
895 Verantwoording [Continuous Naming & Reading Words. Technical Manual]. Amsterdam:
896 Boom Testuitgevers.
- 897 Brem, S., Bach, S., Kucian, K., Kujala, J. V., Guttorm, T. K., Martin, E., ... & Richardson, U. (2010).
898 Brain sensitivity to print emerges when children learn letter-speech sound correspondences.
899 *Proceedings of the National Academy of Sciences*, 107(17), 7939-7944. doi:
900 10.1073/pnas.0904402107
- 901 Brus, B., & Voeten, M. (1991). Een-minuut-test vorm A en B, schoolvorderingstest voor de
902 technische leesvaardigheid bestemd voor groep 4 tot en met 8 van het basisonderwijs.
903 Verantwoording en handleiding. Lisse: Swets & Zeitlinger.
- 904 Byrne, B., Shankweiler, D., & Hine, D. W. (2008). Reading development in children at risk for
905 dyslexia. In M. Mody, & E. R. Silliman (Eds.), *Brain, behavior and learning in language and*
906 *reading disorders* (pp. 240-270). NY: The Guilford Press.

- 907 Carvalhais, L., Limpo, T., Richardson, U., & Castro, S. L. (2020). Effects of the Portuguese
908 Graphogame on reading, spelling, and phonological awareness in second graders struggling
909 to read. *The Journal of Writing Research*, 12(1).
- 910 Chambers, B., Abrami, P., Tucker, B., Slavin, R. E., Madden, N. A., Cheung, A., & Gifford, R.
911 (2008). Computer-assisted tutoring in Success for All: Reading outcomes for first graders.
912 *Journal of Research on Educational Effectiveness*, 1(2), 120-137.
- 913 De Freitas, S. (2006). Learning in immersive worlds: A review of game-based learning.
- 914 Desoete, A., Praet, M., Van de Velde, C., De Craene, B., & Hantson, E. (2016) Enhancing
915 mathematical skills through interventions with virtual manipulatives. In Patricia S. Moyer-
916 Packenham (Eds.) *International Perspectives on Teaching and Learning Mathematics with*
917 *Virtual Manipulatives* (pp.171-187). Springer: Switzerland. doi: 10.1007/978-3-319-32718-1
- 918 D'hondt, M., Desoete, A., Schittekatte, M., Kort, W., Compaan, E., Neyt, F., ... & Surdiacourt, S.
919 (2008). De CELF-4-NL: een opvolger voor de TvK. *Signaal*, 65, 4-16.
- 920 Elen, R. (2006). Proef Fonologisch Bewustzijn (PFB): handleiding, materiaal, scoreformulieren.
921 *Vlaamse Vereniging voor Logopedisten*.
- 922 Froyen, D. J., Bonte, M. L., van Atteveldt, N., & Blomert, L. (2009). The long road to automation:
923 neurocognitive development of letter-speech sound processing. *Journal of Cognitive*
924 *Neuroscience*, 21(3), 567-580. doi: 10.1162/jocn.2009.21061
- 925 Geelhoed, J. W., & Reitsma, P. (1999). PI-dictee.
- 926 Glatz, T. (2018). Serious games as a level playing field for early literacy: A behavioural and
927 neurophysiological evaluation (Doctoral dissertation). Retrieved from the library of the
928 University of Groningen.
- 929 Green, C. S., & Bavelier, D. (2003). Action video game modifies visual selective attention. *Nature*,
930 423(6939), 534.
- 931 Gwee, S., San Chee, Y., & Tan, E. M. (2011). The role of gender in mobile game-based
932 learning. *International Journal of Mobile and Blended Learning (IJMBL)*, 3(4), 19-37.
- 933 Hahn, N., Foxe, J. J., & Molholm, S. (2014). Impairments of multisensory integration and cross-
934 sensory learning as pathways to dyslexia. *Neuroscience & Biobehavioral Reviews*, 47, 384-
935 392.
- 936 Jere-Folotiya, J., Chansa-Kabali, T., Munachaka, J. C., Sampa, F., Yalukanda, C., Westerholm, J., ...
937 & Lyytinen, H. (2014). The effect of using a mobile literacy game to improve literacy levels
938 of grade one students in Zambian schools. *Educational Technology Research and*
939 *Development*, 62(4), 417-436.
- 940 Keuleers, E., & Brysbaert, M. (2010). Wuggy: A multilingual pseudoword generator. *Behavior*
941 *research methods*, 42(3), 627-633. doi: 10.3758/BRM.42.3.627
- 942 Kiili, K. (2005). Digital game-based learning: Towards an experiential gaming model. *The Internet*
943 *and higher education*, 8(1), 13-24.
- 944 Koikkalainen, M. (2015). Computerized reading fluency assessment: Task validity and the strongest
945 discriminators of fluency skills among second-graders. [Master's thesis] University of
946 Jyväskylä. <http://urn.fi/URN:NBN:fi:ju-201510023300>
- 947 Kort, W., Schittekatte, M., & Compaan, E. (2008). *CELF-4-NL: clinical evaluation of language*
948 *fundamentals*. [Dutch version]. Pearson.
- 949 Ktisti, C. (2015). *Computer-based remediation for reading difficulties in a consistent orthography:*
950 *comparing the effects of two theory-driven programs* (Doctoral dissertation). Retrieved from
951 the library of the University of Cyprus.
- 952 Kujala, J. V., Richardson, U., & Lyytinen, H. (2010). A Bayesian-optimal principle for learner-
953 friendly adaptation in learning games. *Journal of Mathematical Psychology*, 54(2), 247-255.
- 954 Kyle, F., Kujala, J., Richardson, U., Lyytinen, H., & Goswami, U. (2013). Assessing the
955 effectiveness of two theoretically motivated computer-assisted reading interventions in the
956 United Kingdom: GG Rime and GG Phoneme. *Reading Research Quarterly*, 48(1), 61-76.

957 Landerl, K., Ramus, F., Moll, K., Lyytinen, H., Leppänen, P. H., Lohvansuu, K., ... & Kunze, S.
 958 (2013). Predictors of developmental dyslexia in European orthographies with varying
 959 complexity. *Journal of Child Psychology and Psychiatry*, 54(6), 686-694.
 960 van der Leij, A., Bergen, E., Zuijlen, T., Jong, P., Maurits, N., & Maassen, B. (2013). Precursors of
 961 developmental dyslexia: an overview of the longitudinal Dutch dyslexia programme study.
 962 *Dyslexia*, 19(4), 191-213. doi: 10.1002/dys.1463
 963 Lovio, R., Halttunen, A., Lyytinen, H., Näätänen, R., & Kujala, T. (2012). Reading skill and neural
 964 processing accuracy improvement after a 3-hour intervention in preschoolers with difficulties
 965 in reading-related skills. *Brain research*, 1448, 42-55. doi: 10.1016/j.brainres.2012.01.071
 966 Lyon, G. R., Shaywitz, S. E., & Shaywitz, B. A. (2003). A definition of dyslexia. *Annals of*
 967 *dyslexia*, 53(1), 1-14. doi: 10.1007/s11881-003-0001-9
 968 Lyytinen, H., Aro, M., Eklund, K., Erskine, J., Guttorm, T., Laakso, M. L., ... & Torppa, M. (2004).
 969 The development of children at familial risk for dyslexia: birth to early school age. *Annals of*
 970 *dyslexia*, 54(2), 184-220. doi: 10.1007/s11881-004-0010-3
 971 Lyytinen, H., Erskine, J., Kujala, J., Ojanen, E., & Richardson, U. (2009). In search of a science-
 972 based application: A learning tool for reading acquisition. *Scandinavian journal of*
 973 *psychology*, 50(6), 668-675. doi: 10.1111/j.1467-9450.2009.00791.x
 974 Mascheretti, S., Bureau, A., Battaglia, M., Simone, D., Quadrelli, E., Croteau, J., ... & Marino, C.
 975 (2013). An assessment of gene-by-environment interactions in developmental dyslexia-
 976 related phenotypes. *Genes, Brain and Behavior*, 12(1), 47-55. doi: 10.1111/gbb.12000
 977 McTigue, E. M., Solheim, O. J., Zimmer, W. K., & Uppstad, P. H. (2020). Critically reviewing
 978 GraphoGame across the world: Recommendations and cautions for research and
 979 implementation of computer-assisted instruction for word-reading acquisition. *Reading*
 980 *Research Quarterly*, 55(1), 45-73.
 981 Ming Chui, M. & McBride-Chang, C. (2006). Gender, context, and reading: A comparison of
 982 students in 43 countries. *Scientific Studies of Reading*, 10(4), 331-362.
 983 doi:10.1207/s1532799xssr1004_1.
 984 Moll, K., Ramus, F., Bartling, J., Bruder, J., Kunze, S., Neuhoff, N., ... & Tóth, D. (2014). Cognitive
 985 mechanisms underlying reading and spelling development in five European
 986 orthographies. *Learning and Instruction*, 29, 65-77. doi: 10.1016/j.learninstruc.2013.09.003
 987 Mommers, M. J. C., Verhoeven, L., & Van der Linden, S. (1990). Veilig leren lezen. *Zwijssen,*
 988 *Tilburg*.
 989 Morton, V., & Torgerson, D. J. (2003). Effect of regression to the mean on decision making in health
 990 care. *Bmj*, 326(7398), 1083-1084.
 991 Nakagawa, S., & Schielzeth, H. (2013). A general and simple method for obtaining R2 from
 992 generalized linear mixed-effects models. *Methods in Ecology and Evolution* 4(2), 133-142.
 993 10.1111/j.2041-210x.2012.00261.x
 994 Nerbonne, J., Heeringa, W., Van den Hout, E., Van der Kooi, P., Otten, S., & Van de Vis, W. (1996).
 995 Phonetic distance between Dutch dialects. In *CLIN VI: Proceedings of the sixth CLIN*
 996 *meeting* (pp. 185-202).
 997 Nietfeld, J. L., Shores, L. R., & Hoffmann, K. F. (2014). Self-regulation and gender within a game-
 998 based learning environment. *Journal of Educational Psychology*, 106(4), 961.
 999 Njå, M. (2019). Players' progression through GraphoGame, an early literacy game: influence of
 1000 game design and context of play. *Human Technology*, 15(2).
 1001 OECD (2010). *PISA 2009 Results: What Students Know and Can Do: Student Performance in*
 1002 *Reading, Mathematics and Science (Volume I)*, PISA, OECD Publishing.
 1003 doi:10.1787/9789264091450-en
 1004 Pennington, B. F. (2006). From single to multiple deficit models of developmental disorders.
 1005 *Cognition*, 101(2), 385-413.

- Piquette, N. A., Savage, R. S., & Abrami, P. C. (2014). A cluster randomized control field trial of the ABRACADABRA web-based reading technology: replication and extension of basic findings. *Frontiers in psychology*, 5, 1413.
- Potocki, A., Ecalte, J., & Magnan, A. (2013). Effects of computer-assisted comprehension training in less skilled comprehenders in second grade: A one-year follow-up study. *Computers & Education*, 63, 131-140. doi: 10.1016/j.compedu.2012.12.011
- Praet, M., & Desoete, A. (2014). Number line estimation from kindergarten to grade 2: a longitudinal study. *Learning and Instruction*, 33, 19-28. doi: 10.1016/j.learninstruc.2014.02.003
- Prensky, M. (2001). Fun, play and games: What makes games engaging. *Digital game-based learning*, 5(1), 5-31.
- Puolakanaho, A., & Latvala, J. M. (2017). Embedding Preschool Assessment Methods into Digital Learning Games to Predict Early Reading Skills. *Human Technology*, 13, 216-236.
- R Core Team (2021). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. retrieved from: <http://www.R-project.org/>
- Regtvoort, A., Zijlstra, H., & van der Leij, A. (2013). The Effectiveness of a 2-year Supplementary Tutor-assisted Computerized Intervention on the Reading Development of Beginning Readers at Risk for Reading Difficulties: A Randomized Controlled Trial. *Dyslexia*, 19(4), 256-280.
- Richardson, U., & Lyytinen, H. (2014). The GraphoGame method: the theoretical and methodological background of the technology-enhanced learning environment for learning to read. *Human Technology: An Interdisciplinary Journal on Humans in ICT Environments*, 10(1), 39-60. doi: 10.17011/ht/urn.201405281859
- Ronimus, M., & Lyytinen, H. (2015). Is school a better environment than home for digital game-based learning? The case of GraphoGame. *Human Technology: An Interdisciplinary Journal on Humans in ICT Environments*.
- Rosas, R., Escobar, J.P., Ramirez, M.P., Meneses, A., & Guajardo, A. (2017). Impact of a computer-based intervention in Chilean children at risk of manifesting reading difficulties. *Infancia y Aprendizaje*, 40(1), 158-188. <https://doi.org/10.1080/02103702.2016.1263451>
- Ruiz, J. P., Lassault, J., Sprenger-Charolles, L., Richardson, U., Lyytinen, H., & Ziegler, J. C. (2017). GraphoGame: un outil numerique pour enfant en difficultes d'apprentissage de la lecture. *ANAE Approche Neuropsychologique des Apprentissages chez l'Enfant*, 148, 333-343.
- Rutter, M., Caspi, A., Fergusson, D., Horwood, L. J., Goodman, R., Maughan, B., ... & Carroll, J. (2004). Sex differences in developmental reading disability: new findings from 4 epidemiological studies. *Jama*, 291(16), 2007-2012. doi: 10.1001/jama.291.16.2007
- Saine, N. L., Lerkkanen, M. K., Ahonen, T., Tolvanen, A., & Lyytinen, H. (2010). Predicting word-level reading fluency outcomes in three contrastive groups: Remedial and computer-assisted remedial reading intervention, and mainstream instruction. *Learning and Individual differences*, 20(5), 402-414. doi: 10.1016/j.lindif.2010.06.004
- Saine, N. L., Lerkkanen, M. K., Ahonen, T., Tolvanen, A., & Lyytinen, H. (2011). Computer-assisted remedial reading intervention for school beginners at risk for reading disability. *Child Development*, 82(3), 1013-1028. doi: 10.1111/j.1467-8624.2011.01580.x
- Sakia, R. M. (1992). The Box-Cox transformation technique: a review. *The statistician*, 169-178.
- Savage, R., Abrami, P. C., Piquette, N., Wood, E., Deleveaux, G., Sanghera-Sidhu, S., & Burgos, G. (2013). A (Pan-Canadian) cluster randomized control effectiveness trial of the ABRACADABRA web-based literacy program. *Journal of Educational Psychology*, 105(2), 310.
- Schaerlaekens, A. M., Kohnstamm, G. A., Lejaegere, M., & Vries, A. K. (1999). *Streeflijst woordenschat voor zesjarigen: gebaseerd op nieuw onderzoek in Nederland en België*. Swets & Zeitlinger.

Schneider, W., Roth, E., & Ennemoser, M. (2000). Training phonological skills and letter knowledge in children at risk for dyslexia: a comparison of three kindergarten intervention programs. *Journal of Educational Psychology*, 92(2), 284.

Schulte-Körne, G., Deimel, W., Bartling, J., & Remschmidt, H., (1998). Auditory processing and dyslexia: evidence for a specific speech processing deficit. *Neuroreport*, 9, 337-340. doi: 10.1097/00001756-199801260-00029

Schulz, K. F., Altman, D. G., & Moher, D. (2010). CONSORT 2010 statement: updated guidelines for reporting parallel group randomised trials. *BMC medicine*, 8(1), 18.

Schumacher, J., Hoffmann, P., Schmä, C., Schulte-Körne, G., & Nöthen, M. M. (2007). Genetics of dyslexia: the evolving landscape. *Journal of medical genetics*, 44(5), 289-297. doi: 10.1136/jmg.2006.046516

Seymour, P. H., Aro, M., & Erskine, J. M. (2003). Foundation literacy acquisition in European orthographies. *British Journal of psychology*, 94(2), 143-174. doi: 10.1348/000712603321661859

De Smedt, B., & Boets, B. (2010). Phonological processing and arithmetic fact retrieval: evidence from developmental dyslexia. *Neuropsychologia*, 48(14), 3973-3981.

De Smedt, B., Taylor, J., Archibald, L., & Ansari, D. (2010). How is phonological processing related to individual differences in children's arithmetic skills?. *Developmental Science*, 13(3), 508-520.

Snowling, M. J., & Melby-Lervåg, M. (2016). Oral Language Deficits in Familial Dyslexia: A Meta-Analysis and Review. *Psychological bulletin*, 142(5), 498-545. doi: 10.1037/bul0000037

Susi, T., Johannesson, M., & Backlund, P. (2007). Serious games: An overview.

Torppa, M., Eklund, K., van Bergen, E., & Lyytinen, H. (2015). Late-emerging and resolving dyslexia: A follow-up study from age 3 to 14. *Journal of Abnormal Child Psychology*, 43(7), 1389-1401.

Tellegen, P. J., & Laros, J. A. (2014). *SON-R 6-40. Snijders-Oomen non-verbal intelligence test*. Göttingen, Germany: Hogrefe

Torbeyns, J., Van den Noortgate, W., Ghesquière, P., Verschaffel, L., Van de Rij, B. A., & Van Luit, J. E. (2002). Development of early numeracy in 5-to 7-year-old children: A comparison between Flanders and The Netherlands. *Educational Research and Evaluation*, 8(3), 249-275.

Vaessen, A., Bertrand, D., Tóth, D., Csépe, V., Fásca, L., Reis, A., & Blomert, L. (2010). Cognitive development of fluent word reading does not qualitatively differ between transparent and opaque orthographies. *Journal of Educational Psychology*, 102(4), 827.

Vaessen, A., & Blomert, L. (2010). Long-term cognitive dynamics of fluent reading development. *Journal of experimental child psychology*, 105(3), 213-231. ISO 690

van Viersen, S., de Bree, E. H., Zee, M., Maassen, B., van der Leij, A., & de Jong, P. F. (2018). Pathways into literacy: The role of early oral language abilities and family risk for dyslexia. *Psychological Science*, 29(3), 418-428.

Wood, S.N. (2006) Generalized Additive Models: An Introduction with R. Chapman and Hall/CRC.

Ziegler, J. C., Bertrand, D., Tóth, D., Csépe, V., Reis, A., Fásca, L., . . . Blomert, L. (2010). Orthographic depth and its impact on universal predictors of reading: A cross-language investigation. *Psychological Science*, 21(4), 551-559.

Ziegler, J. C., & Goswami, U. (2005). Reading acquisition, developmental dyslexia, and skilled reading across languages: a psycholinguistic grain size theory. *Psychological bulletin*, 131(1), 3.

Zijlstra, H., van Bergen, E., Regtvoort, A., de Jong, P. F., & van der Leij, A. (2020, July 13). Prevention of Reading Difficulties in Children With and Without Familial Risk: Short- and Long-Term Effects of an Early Intervention. *Journal of Educational Psychology*. Advance online publication. <http://dx.doi.org/10.1037/edu0000489>

Page 7: [1] Commented [GM30] Genevieve McArthur 10/07/2021 13:52:00

To conduct the study? Or the report the study? And did you actually achieve that aim? (the meaning of this sentence is not clear).

I don't think you have followed reporting guidelines, which I would suggest you since this provides a clearer way of presenting your Methods. So, I would suggest finding the latest version of CONSORT reporting, and arrange the information in your Methods into the sections that CONSORT suggests - including a Design section. See here for latest CONSORT requirements: <http://www.consort-statement.org/checklists/view/32--consort-2010/66-title>

Page 7: [2] Deleted Genevieve McArthur 10/07/2021 13:38:00

Page 7: [3] Commented [GM31] Genevieve McArthur 10/07/2021 13:39:00

when were the tests administered? Pre and post?