**Title: Young EFL learners' processing of multimodal input: Examining learners' eye movements**

Tragant Mestres, Elsa & Pellicer-Sánchez, Ana

There is a tradition among English as a second and foreign language (ESL/EFL) teachers to use authentic materials written for children, in a number of different ways: for example, telling or reading stories, having learners read storybooks or comics silently, or watching cartoons (Shin, 2014). These materials, even if lexically challenging (Webb & Macalister, 2013), are likely to provide a rich source of input either in the classroom context or as out-of-class material, and may complement the more carefully planned input students are likely to encounter in ESL/EFL textbooks. In part, this is because the rich visual element that often accompanies these children's texts, and the pattern and plot repetition they often include, may aid comprehension and promote language learning. In addition, materials written for children have been found to be motivating for young language learners (Ghosn, 2002; González & Taronna 2012).

The present study examines the ways in which young EFL learners engage with two types of authentic materials: a cartoon with second language (L2) subtitles, and a storybook with audio support. While both of these materials offer language learners

multimodal input through images, audio support and text /captions, there are also important differences between the two that cannot go unnoticed. One of these differences is the fact that the non-verbal input is provided through dynamic clips in motion in cartoons in contrast to the bi-dimensional static visual input in storybooks. It may be the case that the former has the added value of being engaging to the child's eye and of making the plot more explicit and easy to understand. However, it is also possible that the attention children pay to the dynamic images in the video is detrimental to the attention that is left for them to attend to the processing of verbal input. In addition to the distinct features of the non-verbal input in the two materials, there are clear differences in the presentation of the text in the two modalities. The speed and dynamism in the presentation of subtitles in the video condition is evidently different from the more static presentation of the text in the storybook condition. Significant differences between the processing of the text and images in these two types of multimodal materials are therefore expected. However, a description of how young learners engage with these types of multimodal materials is yet to be provided. Thus, the work that we set to carry out aims at examining and describing how young language learners process text when accompanied by dynamic vs. static images and how they split their attention between the verbal and non-verbal sources of input in two types of multimodal materials. Empirical research has shown that both static and full-motion images can facilitate comprehension (Chang & Read, 2006, 2007; Herron, Julia & Cole, 1995) but we still do not know how the dynamic nature of the non-verbal input in movies affects young learners' distribution of attentional resources. Having a clearer picture of learners' processing in these two types of multimodal materials would allow us to provide a better evaluation of the potential of these materials for language learning and to predict potential detrimental effects to learning and comprehension.

**Background**

*Reading while listening*

In illustrated storybooks with audio support, the verbal input is provided both through aural and written modes. The combination of written and aural input is believed to support the reading process, leading to improved comprehension, learning gains and more positive attitudes. Simultaneous written and aural input, referred to here as 'reading while listening' and also known as 'assisted reading', was originally conceived in the 1970s to support children with reading difficulties (Carbo, 1978). In the late seventies and eighties, several shared book reading programmes, which involved children's exposure to printed texts that were simultaneously read aloud by the teacher, were implemented for the acquisition of L2 literacy in elementary schools (Elley, 1991). More recently, researchers have started to use this instructional technique with older L2 learners, often in combination with extensive reading and/or repeated reading, an instructional intervention involving multiple successive readings of a passage (see for example, Han & Chen, 2010).

In L2 research several studies have compared 'reading while listening '(one type of bimodal input) with 'reading only' and/or 'listening only' (two types of unimodal input). Most of these studies (i.e., Chang, 2009; Chang & Millett 2015; Taguchi, Takayasu-Maass & Gorsuch, 2004) have focused on how these modes of input contribute to improvement in comprehension and fluency with the general conclusion that the effectiveness of reading while listening largely depends on the amount of reading materials over time. So, in studies like Chang and Millett (2014) where students were exposed to ten graded readers over a 13-week period, reading while listening produced superior results to just reading or listening. In contrast, in studies that involved

lower weekly amounts of input like in Chang (2012), where students read 15 audio graded readers over a period of 26 weeks, the gains of the reading while listening group were much smaller. The intensity of the intervention also seems to play a role when it comes to vocabulary learning in the context of reading while listening. Whereas Brown, Waring and Donkaewbua (2008) found that most words were not learned no matter the input mode in a laboratory study where students were exposed to three stories on one sitting, Webb and Chang (2012) did find significant differences in vocabulary learning between the reading while listening group and the reading only group. In this case the intervention involved 28 texts over a period of 14 weeks. In an attempt to explain the superiority of the reading while listening group, these authors provide several hypotheses like the fact that students in that group may have had more time left to focus on vocabulary given that their reading speed was determined by the oral text. Nevertheless, this and other comparative studies cited in this paragraph can only speculate in their explanation since they look at the outcomes of processing written/oral text. In contrast, our study attempts to depict the process, what students do when they read and listen simultaneously.

In the type of storybook with audio support that we examine in the present study, another important source of input is the visual component, i.e., the static images. Children books and graded readers very often include images to support the reading process. Information processing theories suggest that this combination of verbal and non-verbal input supports language learning. According to Paivio's Dual Coding Theory (1986), the two different input modes are processed through two different systems, the verbal system and the imagery system, and the simultaneous activation of the two systems in multimodal materials such as illustrated storybooks leads to higher learning gains. This theory was later extended by Mayer and Sims (1994) in the context

of multimedia learning, who claim that meaningful learning takes place when learners build three types of connections: the visual representational connections, the verbal representational connections, and the referential connections between visual and verbal representations. In accordance with the contiguity effect, they explain that learners are able to form better referential connections when verbal and visual materials are presented contiguously than when presented separately. Previous research has shown that the use of static images to support the reading process leads to increased comprehension in both the first language (L1) (e.g., Gambrell & Jawitz, 1993; Hall, Bailey & Tillman, 1997) and the L2 context (e.g., Elley & Mangubhai, 1983; Omaggio, 1979), improved vocabulary learning (e.g., Bisson et al., 2015) and positive perceptions by learners (e.g., Tang, 1992). However, very little is known as yet about how learners make use of the visual input while reading in this type of multimodal materials.

The study of eye movements may provide some insights into the processing of text and images. Eye-tracking provides a detailed record of online processing behavior. Eye-tracking research has provided useful insights about typical eye-movement behavior when processing verbal and non-verbal stimuli (Rayner, 1998). In the context of L2 reading, eye-tracking research has shed light into our understanding of how L2 learners process written words during naturalistic reading and of the different factors that affect this process, such as frequency (Whitford & Titone, 2012), cognate status (e.g. Duyck, Van Assche, Drighe, & Hartsuiker, 2007) and orthographic neighborhood (e.g. Dirix, Cop, Drieghe, Duyck & Hartsuiker, 2017). Second language acquisition researchers have also recently started to use eye-tracking to explore the relationship between patterns of processing and learning of several linguistic features, including grammatical structures (e.g. Winke, 2013), and vocabulary (e.g. Elgort, Brysbaert, Stevens & Van Assche, 2017; Godfroid, Boers & Housen, 2013; XXX, 2016).

Research using this technique has also contributed to our understanding of how texts are processed in different conditions. Texts are processed more slowly in oral reading (i.e., reading while listening or reading aloud) compared to silent reading, and fixations (i.e., the period of time when the eyes remain still) tend to be longer. According to Rayner (2009), mean fixation duration in oral reading ranges from 275-325 ms and from 225-250ms in silent reading. Interesting insights about the processing of images have also emerged from eye-tracking research. Fixations on images (260-330ms) tend to be longer than in silent reading, since during scene perception useful information is gained from a fairly wide field of view (Rayner, 2009). However, very few eye-tracking studies have looked at how learners process text and images simultaneously in reading. In the L1 context, previous studies with pre-literate children have found that little attention is paid to the text when compared to the illustrations in shared reading experiences (e.g., Evans & Saint-Aubin, 2005; Justice, Skibbe & Canning, 2005) and that attention to print increases with reading proficiency (e.g., Roy-Charland, Saint-Aubin & Evans, 2007). A few other studies have examined how young L1 readers process the text and images when reading to learn content and the effect that adding non-verbal support has on learning (e.g. Mason et al. 2013; Mason & Tornatora, 2015). To the best of our knowledge, similar research looking at the simultaneous processing of text and images in the L2 context, and in particular with young learners reading for comprehension, is yet to be conducted.

*Watching subtitled videos*

Most of the studies exploring the effect of subtitles on comprehension are based on short videos. Overall, the comparison between the performance of learners under captioned (i.e., subtitles and soundtrack in the same language) and uncaptioned conditions have generally found an advantage for the captioned condition, providing

evidence for the positive effect of captions to support the comprehension process (e.g., Guillory, 1998; Markham & Peter, 2003; Montero Perez, Peters & Desmet, 2014). Fewer studies have looked at the use of captions with longer viewing conditions. For example, comparing the comprehension of ten TV episodes under captioned and uncaptioned conditions Rodgers and Webb (2017) found that the advantage of the captioned group was only evident in some of the episodes, and that factors like difficulty and position in the viewing process also played a role. Particularly relevant for the present investigation are recent studies using eye-tracking to examine L2 learners' processing of the different sources of input in captioned and subtitled videos. Bisson, van Heuven, Conklin and Tunney (2014), for example, examined participants' processing of different types of subtitle conditions, i.e., standard (foreign language (FL) soundtrack and native language (NL) subtitles), reversed (NL soundtrack and FL subtitles), and intralingual (FL soundtrack and FL subtitles), while watching 25 minutes of a FL movie. Examining the fixations made in the image and subtitle areas, they found that both areas were processed in all conditions; participants read the subtitles regardless of the subtitle condition, but more regular reading was exhibited when the soundtrack was in an unknown FL. Learners' processing of L2 subtitled videos and the time spent on L2 captions has also been found to be affected by L1 background (Winke, Gass & Sydorenko, 2013). Montero Perez, Peters and Desmet (2015) examined the learning of L2 vocabulary from subtitled movies, assessing learners' processing and learning of L2 words in two types of captioning (i.e., full captioning and keyword captioning) and under two test announcement conditions (i.e., the intentional condition in which participants were informed of the upcoming vocabulary post-tests, and the incidental condition in which they were not informed). Analyses of the fixation measures on the target words showed a significant interaction between type of captioning and test

announcement. Significant correlations were found between reading times and word learning for learners in the full captioning intentional group, with longer reading times being associated with higher recognition scores.

The processing of subtitles by adults and children was compared in d'Ydewalle and de Bruycker (2007), who examined children's (Grade 5 and 6) and adults' eye movements while watching a FL movie in two subtitling conditions, i.e., standard subtitling (FL soundtrack and NL subtitles) and reversed subtitling (NL soundtrack and FL subtitles). The results showed irregular reading patterns in both conditions in both groups. Regarding age differences, the results showed that children took longer to shift attention to the subtitle at its onset and had longer fixations and shorter saccades in the text.

**The present study**

Overall, the studies reviewed in the above section show that multimodal exposure seems to be beneficial for language learning, but there is no clear picture as yet of how children's reading behavior changes when exposed to different multimodal materials. In addition, little is known about the amount of attention received by the different input sources. To shed light on these areas, the present study intends to answer the following two questions using eye-tracking methodology:

(1) What is EFL young learners' reading behavior in the storybook and the video formats?

(2) How is attention split between the processing of text and visual input in the storybook and video formats?

**Methodology**

*Participants*

Students from two semi-private schools in Barcelona were initially selected to participate in this study (henceforth, schools A and B). The schools were comparable in terms of the students' family profiles and the importance given to English in their curriculum. In both schools more time was devoted to English than is mandatory, and English instruction started at the age of three.

A group of 36 Grade 5 students participated in the study aged either 10 or 11 years old (a mean age 10.74). There were 21 girls and 15 boys (23 from school A and 13 from school B). The students were selected by the teachers from seven intact classes (five from school A and two from school B) following the researchers' instructions to choose only students without reading or learning difficulties. Their scores in the Peabody Picture Vocabulary Test (Dunn & Dunn, 2007) ranged from 31 to 131 ($M$=91.5, $SD$=28.14). While the majority of these students (N=22) had some experience of book reading in English beyond the classroom, only one of them reported having experience of FL subtitled videos.

Barcelona is a bilingual city, and all the students in the study were competent users of both Catalan and Spanish. In the case of the participants in this study, 52% of them spoke Catalan at home, 16% spoke Spanish, and 26% spoke both. There were two students who spoke Catalan and Italian at home (6%). As regards the educational background of the families, most of the fathers at the two schools (94% at school A and 85% at school B) and the mothers (100% at school A and 85% at school B) were university graduates.

*Design*

The experiment followed a within-subjects design in which all the participants were exposed to the two formats under study: storybook with audio support and video with FL subtitles. The presentation of the two formats was counterbalanced. Students in Condition 1 (n=19) read while listening to the first part of the story and then watched the video of the second part with subtitles. Students in Condition 2 (n=17) watched the video first and continued with reading while listening. Participants were randomly assigned to conditions 1 and 2 and comprehension of the story in the two formats was measured in order to control for potential differences in comprehension.

*Instruments*

*Reading Materials*. The present study uses an episode from the series "Charlie and Lola", an authentic series for young children (3 to 7 years old) created by Lauren Child. The episode, "We honestly can look after your dog", was selected since it was available both as a video (with FL subtitles) (BBC, 2006) (https://www.youtube.com/watch?v=Nw5iqOHvB8w) and as a picture book with CD (Child, 2005), the two formats examined in the study. One characteristic feature of the book collection is the layout and typestyle of the text, which shows considerable variation from page to page. As for the video collection, the subtitles are intended for the Deaf and Hard-of-Hearing, which means that important non-dialogue audio sound effects and speaker identification are included in the subtitles.

The video and book formats are comparable in terms of the length of the soundtrack[i] and the pace of reading (150 words per minute in both formats)[ii]. As regards the text, an analysis of the vocabulary profile of the two formats with Lextutor (Cobb, 2016) showed that they were very similar in terms of the frequency profile, type/token ratio and lexical density (See Table 1 for a comparison of the video and book formats). The percentages of open word classes are comparable too as well as the average number

of characters per word and the syllable count. The average number of words per sentence is higher in the case of the book (6.9) in comparison to that of the video (5.7) mainly due to the fact that many of the sentences in the book started with a reporting clause which was absent in the video format (See Appendix A for sample excerpts from the two formats). The Flesch Reading Ease score for the text in both formats was classified as 'very easy to read' (storybook format= 64; video format= 97.9). The Text Readability Consensus Calculator (www.readabilityformulas.com) was used to further explore the similar level of difficulty and readability of the text in both formats. This calculator uses seven popular readability formulas (Flesch Reading Ease score, Gunning Fog, Flesch-Kincaid Grade Level, The Coleman-Liau Index, The SMOG Index, Automated Readability Index, Linsear Write Formula) to calculate the average grade level, reading age, and text difficulty of a text. The results of this calculation showed that for the texts in both formats the grade level was 2 and the reading level 'very easy to read'.

[Place Table 1 about here]

In order to prepare the stimuli for the eye-tracking session, the episode was divided into part 1 and part 2 in the two formats. The optimal place to split the story was chosen taking into account the length of the story and the plot. Table 2 shows the main features of the text/audio of the two conditions. These materials were used to design the eye-tracking experiment with Tobii Pro Studio (version 3.4.2).

[Place Table 2 about here]

*Comprehension test.* A 16-item test was used to check that the comprehension of the story was comparable when reading the book and when watching the video. The test, which was previously piloted, included questions that required either short (i.e., Who is

the owner of the dog?) or longer answers (i.e., Why does Marv ask the dog to sit?). Since certain images from the video were more explicit about the plot of the story than the illustrations from the book, special care was taken not to ask for information from those more revealing video images. The questions were worded and administered in Catalan, most students' L1 and the language students are taught in at school, because we wanted students to answer in a language they had full command of. Appendix B includes the complete set of questions.

*Procedure*

The experimental session was conducted individually in a quiet room on the school premises. It started with the administration of the Peabody Picture Vocabulary Test. Then participants were informed that they would watch a story in English, and they were shown the cover of the DVD to introduce them to the two characters in the story, Charlie and Lola. They were not informed that they would visualize one part of the story in book format and the other in video format. They were told that they would be asked a few questions about the story in Catalan after the experiment.

The eye-tracker was then set up and a 5-point calibration of the equipment was conducted. Participants' eye-movements were recorded using Tobii T120 (Tobii, www.tobii.com), a remote, desktop eye-tracker, with the camera and infrared light integrated in the monitor. It has a sampling rate of 120 Hz, which is considered adequate for the examination of fixations to larger regions of interest (XXX, 2016). It has a typical accuracy of 0.5° (measured in ideal conditions) and 0.2° resolution. Recording was done binocularly and data quality was checked (minimum recording accuracy 70%) (for a detailed discussion of types of recording and effects on data quality see XXX, 2018). Data from the left eye was included in the analyses. The stimuli were displayed on a 24" screen using Tobii Pro Studio (version 3.4.2).

In both conditions, participants visualized the full story and heard the soundtrack (with headphones) once via computer. In the book part, pages changed automatically in synchrony with the soundtrack. After viewing the story, participants were asked to answer the comprehension questions orally, and their answers were written down by the researcher. The whole procedure took about 30-40 minutes.

*Analyses*

In the process of analyzing the eye-movement data, areas of interest (AOIs) were first created for the selected pages in the storybook and the corresponding subtitles/images in the video. The text/subtitle areas and the image areas constituted the two types of AOIs in the study. Unlike most video studies where AOIs are created for each subtitle, in this study the AOIs in the video condition included groups of subtitles/scenes that correspond to the text in the selected pages from the book. This was done to allow for a better comparison of the results. In the case of the video, all the image and subtitle AOIs had the same size and position and took up the whole width of screen. In the case of the book, the size, shape and position of the image and text AOIs were different for each page because of the variable layout. Book pages were discarded when the lines of the text were not horizontal, when the position of the text and image was different from most other pages, or when the font size of several words on a page was much larger than most of the text in the book. As for the video, text was also excluded when visual aspects of the video co-occurred with the subtitle area. In order to identify co-occurrences, a group of five students who were not participating in the study were asked to watch the video without subtitles, and their eye-movements were recorded. After the appropriate deletions, 12 pages from the book and the corresponding 58 subtitles from the video were included in the analysis (see Tables 3 and 4).

[Place Table 3 about here]

[Place Table 4 about here]

The eye-tracking measures used in the present study are based on fixations, which can be described as the periods of time between saccades when the eyes remain fairly still. We explored fixations to the text and the image AOIs. Three late eye-tracking measures were examined: *average fixation duration* (i.e., the mean of the duration of each individual fixation within an AOI), *total fixation duration* (i.e., the sum of all fixation durations made within an AOI), and *fixation count* (i.e., the total number of fixations made within an AOI). Repeated measures t-tests were run to answer research questions 1 and 2, and the level of significance was set at .05.

Prior to data analysis, eye-movement data were inspected for outliers. Data from two participants were excluded from the analysis because the quality of the recording was below 70%. The final sample thus comprised 34 participants (condition 1 n=19; condition 2 n=15). Fixations shorter than 70 ms were also discarded and short fixations were not merged, meaning the loss of 19% of the data points (4627 of the total 23661 data points). This percentage of data loss is slightly higher than those reported in most reading studies. The more frequent blinking and head movement that is common with children could lead to unreliable contact with the eye-tracker which then results in shorter fixations (Was, Smith & Johnson 2013)[iii].

**Results**

Comprehension of the story was first checked to make sure that it was similar in the two viewing formats. Independent sample t-tests showed that there were no significant differences between the comprehension scores for the storybook and video formats (see

Table 5). This is so both for part 1 of the story [$t(32)=1.2$, $p=.24$] and part 2 [$t(32)=1.39$, $p=.68$], the effect size being small in the two cases (.05 and .04 respectively). Total scores (part 1 + part 2) of participants in Condition 1 (*M*=6,21, *SD*=2.23) and 2 (*M*=6.27, *SD*=2.71) were not significantly different either [$t(32)=0.07$, $p=.95$], and the effect size was small (.01); any differences observed in the eye-movement patterns could therefore be attributed to the different formats under examination and not to differences in comprehension.

[Place Table 5 about here]

**Research question 1: reading behavior**

Learners' reading behavior with regard to the text AOIs in the storybook format and the subtitle AOIs in the video format were first investigated (see the descriptive measures in Table 6). As explained in the Analysis section, the reading of the text area of each page in the storybook and the reading of the corresponding set of subtitles was examined. Since the number of words and the duration of the book and video formats were not exactly the same, the number of fixations and the total fixation duration during text/subtitle reading were normalized: the number of fixations were divided by the number of words, and the total fixation durations were divided by presentation time.

The examination of the normalized number of fixations in the two formats (see Table 7) showed that the ratios are very similar and amount to over one fixation per word (1.11 in the case of the book and 1.19 in the case of the video). In reading the storybook, the average fixation duration (247.37ms) is within the range of values provided for silent reading (225-250ms) by skilled readers of English and lower than those provided for oral reading (275-325ms) (Rayner, 2009). In the case of subtitles in the video format, the average fixation duration (214.97ms) is shorter than in silent

reading and shorter than when children read standard subtitles (247ms) (D'Ydewalle &
De Bruycker, 2007).

[Place Table 6 about here]

[Place Table 7 about here]

A further analysis was conducted to explore how students' English proficiency related
to their eye movements by using Pearson Product-moment correlation coefficient.
Regarding book reading, correlations between the vocabulary scores and the three
dependent variables (average fixation duration, total fixation duration and fixation
count) were negative but non-significant. Regarding subtitle reading, correlations were
also negative and there was one significant correlation with average fixation duration
[$r=-.39$, $n=34$, $p<.05$] but not with the other two measures.

**Research question 2: processing of text vs. image**

In order to analyze the amount of time devoted to reading the text/subtitles in relation to
the images, a percentage was calculated based on the total duration of all fixations (text
and image AOIs). The examination of the mean percentage time devoted to the
text/subtitle areas showed that learners spent a higher percentage of the viewing time
reading the text/subtitles than looking at the images/video in the two formats (see Table
8).

[Place Table 8 about here]

Further analysis was conducted to explore how participants processed text vs
image by classifying them into three groups: those who had spent less than 50% of their
time reading the text/subtitles, those who had spent between 50-70% of their time, and a

16

third group with those students who had read the most (more than 70% of their time). Figure 1 shows that almost all students (31 out of 34) in the book format spent most of their time reading (>70%). In the case of watching the video a third of the participants did so (12 out of 34), while the rest were divided between attending to the subtitles for 50-70% of their time (15 out of 34) and for less than half of their time (7 out of 34).

[Place figure 1 about here]

Some other interesting differences between the processing of the text and images were also observed. The results in Table 7 show that when reading a book, average fixation durations are longer when reading than when looking at the illustrations (247.37ms vs. 205.17ms). When watching the video the opposite was the case, with much shorter average fixation durations when reading the subtitles than when watching the images (214.97ms vs 306.38ms). Differences were significant in the case of both the book [$t(33)$=4.49, $p$<.000] and the video [$t(33)$=9.75, $p$<.000], and effect sizes were large (0.38 and 0.74 respectively).

The results also show a higher number of fixations when reading the text than when looking at the illustrations (34.20 vs. 5.79). A similar pattern was found in the case of subtitle reading vs watching the images (35.96 vs. 16.58). Similarly, descriptive statistics suggested a longer total fixation duration in the text/subtitle AOIs than in the image AOIs in both formats. However, these results should be treated with caution, as the sizes of the text/subtitle AOIs were smaller than the image AOIs. Importantly, these differences in the processing of the text and images are driven by the different nature of these two types of input (i.e verbal vs. non-verbal).

**Discussion**

The present study has examined two types of multimodal material that are quite different in nature (both in terms of the images and the presentation of the text) and this should call for some caution when interpreting the results. The inclusion of these two types of materials in the present paper was in part motivated by the fact that they are two types of multimodal materials that are often used in the L2 classroom as activities to improve reading and listening comprehension and teachers often find themselves in the dilemma of choosing between books or cartoons as complementary sources of input in children's EFL lessons. In addition, in the ELT literature, reading while listening to books and extensive viewing to TV are often referred to as alternative options to extensive reading (i.e., Siyanova-Chanturia & Webb, 2016) and in fact Renandya and Jacobs (2016) use the same term ('extensive listening') to refer to these two activities. As explained in the introduction, it is expected that the evident distinctive features of these two types of multimodal materials will lead to processing differences. However, a better understanding of how audio-storybooks and subtitled videos are processed is still needed in order to evaluate how these two formats can improve L2 comprehension and language learning. Thus, the present study attempted to contribute to this better understanding.

With regard to the first research question about learners' reading behavior, the results show that the probability that an individual word would be fixated was similar in the two formats. If we compare our video data with d'Ydewalle and De Bruycker's subtitle data reported in children (2007), we obtain a higher word fixation probability (1.12) than when Dutch-speaking children were watching excerpts from a movie in Swedish with reversed subtitles (0.54), and a slightly higher probability than when these children were watching the movie with standard subtitles (0.92). These differences could be explained by the fact that the Dutch children did not know any Swedish

whereas our students had been learning English for several years and were watching the video under intralingual subtitling conditions.

Further investigations of the book and video formats were conducted taking into account the duration of the fixations. In storybook reading, average fixation durations were in line with those identified for silent reading (Rayner, 2009), whereas in the video format they were shorter than in silent reading. Comparison with video data from the study by Bisson et al. (2014) with adults shows that normalized total fixation durations are very similar (0.43 in their study and 0.46 in ours). This indicates that the reading behavior of subtitles seems to be similar in both age groups even though the standard deviation was higher in our study (0.13 compared with 0.03 in their study). Our examination of the relationship between learners' processing of the text and their level of English proficiency (as measured by the Peabody Picture Vocabulary test) points towards an interesting relationship between these two variables. In the processing of subtitles in the video condition, learners' average fixation durations were negatively correlated with proficiency levels, with higher scores in the proficiency test associated with shorter average fixation durations. This would suggest that, as expected, learners with a higher proficiency would show a faster and more fluent reading. However, this relationship failed to reach the significance level in the storybook condition and with the other processing measures. Future studies with learners of a wider range of proficiencies should further explore this relationship.

As for the second research question regarding the processing of the text and images in the two formats, our results show that learners process both input sources (text and image AOIs) in both modalities, in contrast to the results of studies with pre-literate children (Evans & Saint-Aubin, 2005; Justice, Skibbe & Canning, 2005) but in line with findings by Bisson et al. (2014) with adult learners. Thus, younger learners'

behavior is similar to that of adult learners with regard to both sources of input. The learners in this study tended to pay considerable attention to the images in the video, as evidenced by the long percentage times. It seems that the dynamic nature of the visual input in the video condition does not distract learners' attention from the text entirely, but it does make learners frequently attend to the visual component in this format. The patterns found in the processing of images in both conditions are closely connected, and partially a consequence of the patterns of text processing found as they all relate to the same underlying differences between the two conditions being examined.

The examination of the average fixation durations for the text and image areas yields interesting results. In the video format, fixations were on average shorter when reading the subtitles than when watching the video (214.97ms and 306.38ms respectively). This is in line with the average fixation durations for adult L1 readers, who showed longer fixations in visual scene perception (330ms) than in oral reading (275ms) (Rayner, 1998). However, the opposite pattern was found in the book condition, with longer average fixation durations in the text than in the images. The more regular and careful reading shown by young learners in the storybook modality may account for this difference.

The results concerning the percentage of time devoted to the text in relation to the image in the book format shows that students spent proportionally longer time processing text and less time processing the images. In the case of the video, the trend is less pronounced but the proportions from our data (62.6) are considerably higher than those reported for the children (47) in d'Ydewalle and De Bruycker's study (2007). However, this difference could also be due to the fact that the procedures applied to calculate these proportions were not the same. The fact that those children were reading standard subtitles may have enabled them to split their attention more evenly between

the subtitles and the images than our students who were reading in a foreign language. In addition, the instructions that our students received at the start of the experiment, in which they were explicitly asked to read from the screen, might also have influenced their reading behavior.

Further inspection of the percentage times in the video data shows high individual variability. This variability may indicate that students react differently to the verbal and visual information in the video format. This difference may partially be due to previous experience with this type of material. As explained in the methodology section, most learners in this study reported having little experience with subtitled viewing.

In any case, our data on normalized fixation durations and percentage time on text/subtitles confirm that students spend a considerable amount of time processing the text both when reading the book and when watching the video, and that learners do attend to the verbal input in both formats. Nevertheless, the book format seems to show more regular reading patterns whereas the video format presents high variability among learners.

The present study has focused on the online processing of the different sources of information in two types of multimodal materials but it has not compared the attention given to the written and aural modes of verbal input. Nor has it examined whether the processing differences observed are related to story comprehension. These are issues that remain to be explored. It would also be interesting to look into the possible relationships with the students' perceptual styles. It could be the case that students favoring a visual style spend different percentage times on subtitles from those favoring a verbal learning style.

Even though the administration of the audiobook condition via computer may be considered a less authentic reading experience in certain contexts, the use of authentic materials in the present study increases its ecological validity. However, this also obliged us to apply a very conservative approach in the creation and selection of AOIs for the analysis of eye-movement data, and to reject a large number of pages and subtitles because of the unconventional layout of the book and the overlapping of action in the subtitle areas. Future studies conducted with specifically designed materials would allow researchers to control for the position and specific features of the input, which are factors that influence eye-movements.

**Conclusions**

The results of the study provide evidence that young learners process both sources of information (i.e., written verbal information and visual information) in multimodal materials, and that they spend a longer time processing the text than the visual component in both formats, and challenge the assumption that the engaging nature of the visual information may distract learners' attention from the text. The study has also shown that a more regular reading pattern is observed in the storybook condition with audio support, which is potentially caused by the non-dynamic nature of the images in this condition. Overall, the study has shed light on our understanding of how young learners engage with different types of multimodal materials. Future investigations are necessary in order to determine whether the processing patterns observed in this study are reflected in differences in comprehension or learning. Finally, the study shows that the use of eye-tracking can provide a richer picture of multimodal learning and opens up an important new avenue for L2 research.

**References**

BBC. (2006). We do Promise Honestly we can Look after your Dog. In *Charlie and Lola One*. BBC Children's DVD.

Bisson, M-J., Van Heuven, W., Conklin K., & Tunney, R. (2014). Processing of Native and Foreign Language Subtitles in Films: An Eye Tracking Study. *Applied Psycholinguistics, 35*, 399-418. doi: 10.1017/S0142716412000434.

Bisson, M-J., Van Heuven, W., Conklin, K., & Tunney R. (2015). The Role of Verbal and Pictorial Information in Multi-Modal Incidental Acquisition of Foreign Language Vocabulary. *Quarterly Journal of Experimental Psychology, 68,* 306–26. doi: 10.1080/17470218.2014.979211.

Brown, R., Waring, R., & Donkaewbua S. (2008). Incidental Vocabulary Acquisition from Reading, Reading-While-Listening, and Listening to Stories. *Reading in a Foreign Language*, *20(2)*, 136-163.

Carbo, M. (1978). Teaching Reading with Talking Books. *The Reading Teacher, 32,* 267-273.

Carrol, G., & Conklin K. (2017). Cross language priming extends to formulaic units: evidence from eye-tracking suggests that this idea "has legs". *Bilingualism: Language and Cognition (Special Issue on Cross Language Priming), 20(2)*, 299-317. doi:10.1017/S1366728915000103.

XXX. (2016). Using Eye-Tracking in Applied Linguistics and Second Language Research. *Second Language Research, 32*(3), 453-467. doi: 10.1177/0267658316637401.

XXX. (2018). *Eye-tracking: A Guide for Applied Linguistics Research*. Cambridge: Cambridge University Press.

Chang, C.-S. (2009). Gains to L2 Listeners from Reading While Listening vs. Listening only in Comprehending Short Stories. *System, 37*, 652-663. doi: 10.1016/j.system.2009.09.009.

Chang, C.-S. (2012). Improving reading rate activities for EFL students: timed reading and repeated oral reading. *Reading in a Foreign Language,* 24(1), 56-83.

Chang, C.-S., & Millett S. (2014). The Effect of Extensive Listening on Developing L2 Listeing Fluency: Some Hard Evidence. *ELT Journal, 68(1)*, 31-40. doi: 10.1093/elt/cct052.

Chang, C.-S., & S. Millett S. (2015). Improving Reading Rates and Comprehension through Audio-Assisted Extensive Reading for Beginner Learners. *System*, *52*, 91-102. doi: 10.1016/j.system.2015.05.003.

Chang, C.-S., & Read J. (2006). The effects of listening support on the listening performance of EFL learners. *TESOL Quarterly 40*, 375–397.
doi: 10.2307/40264527.

Chang, C.-S., & Read J. (2007). Support for foreign language listeners: Its effectiveness and limitations. *RELC Journal*, *38*, 373–394. doi: 10.1177/0033688207085853.

Child, L. (2005). *We Honestly can Look after your Dog*. London: Puffin.

Cobb,T. Lextutor Vocabprofile [accessed January 2016 from http://www.lextutor.ca/vp/], an adaptation of Heatley, Nation & Coxhead's (2002) Range.

Dirix, N., Cop, U., Drieghe, D., Duyck, W., & Hartsuiker, R. J. (2017). Cross-lingual neighborhood effects in generalized lexical decision and natural reading. *Journal of Experimental Psychology: Learning, Memory and Cognition, 43(6)*, 887-915. doi: 10.1037/xlm0000352.

Duyck, W., Van Assche, E., Drighe, D., & Hartsuiker R. (2007). Visual word recognition by bilinguals in a sentence context: Evidence for nonselective lexical access. *Journal of Experimental Psychology: Learning, Memory, and Cognition 33(4)*, 663-679. doi: 10.1037/0278-7393.33.4.663.

Dunn, Ll., & Dunn D. (2007). Peabody Picture Vocabulary Test (4th ed.). Pearson.

D'Ydewalle, G., & de Bruycker W. (2007). Eye Movements of Children and Adults while Reading Television Subtitles. *European Psychologist, 12(3)*, 196-205. doi: 10.1027/1016-9040.12.3.196.

Elgort, I., Brysbaert, M., Stevens, M., & Van Assche E. (2017). Contextual word learning during reading in a second language: An eye movement study. *Studies in Second Language Acquisition,* 1-26. doi:10.1017/S0272263117000109

Elley, W. (1991). Acquiring literacy in a second language: The effect of book-based programs. *Language Learning*, 41(3), 375-410. Doi: 10.1111/j.1467-1770.1991.tb00611.x.

Elley, W. B., & Mangubhai F. (1983). The Impact of Reading on Second Language Learning. *Reading Research Quarterly, 19(1)*, 53-67. doi: 10.2307/747337.

Evans, M. A., & Saint-Aubin J. (2005). What Children are Looking at during Shared Storybook Reading: Evidence from Eye Movement Monitoring. *Psychological Science, 16 (11)*, 913-920. doi: 10.1111/j.1467-9280.2005.01636.x.

Gambrell, L. B., & Jawitz P. B. (1993). Mental Imagery, Text Illustrations, and Children's Story Comprehension and Recall. *Reading Research Quarterly,* 28(3), 264-276. doi: 10.2307/747998.

Ghosn, I. K. (2002). Four Good Reasons to Use Literature in Primary School ELT. *ELT Journal, 56,* 172-179. doi: 10.1093/elt/56.2.172.

Godfroid, A., Boers, F., & Housen A. (2013). An eye for words: Gauging the role of attention in incidental L2 vocabulary acquisition by means of eye tracking. *Studies in Second Language Acquisition, 35*, 483-517. Doi: 10.1017/S0272263113000119.

González, M., & Taronna A. (2012). *New Trends in Early Foreign Language Learning*. Newcaslte upon Tyne: Camp Scholars Publishing.

Guillory, H. G. (1998). The Effects of Keyword Captions to Authentic French Video on Learner Comprehension. *CALICO Journal, 15(1–3)*, 89–108. doi: 10.1558/cj.v15i1-3.89-108.

Hall, V. C., Bailey, J., & Tillman C. (1997). Can Student Generated-Illustrations be Worth Ten Thousand Words? *Journal of Educational Psychology, 89(4)*, 677-681. doi: 10.1037/0022-0663.89.4.677.

Han, Z., & Chen A. C. (2010). Repeated-reading-based instructional strategy and vocabulary acquisition: A case study of a heritage speaker of Chinese. *Reading in a Foreign Language, 22(2)*, 242–262.

Herron, C., Julia, E. B., & Cole S. P. (1995). A comparison study of two advance organizers for introducing beginning foreign language students to video. *The Modern Language Journal, 79*, 387–395. doi: 10.2307/329353.

Justice, L. M., Skibbe, L., Canning, A., & Langkford C. (2005). Pre-Schoolers, Print and Storybooks: An Observational Study Using Eye Movement Analysis. *Journal of Research in Reading, 28(3)*, 229–243. doi: 10.1111/j.1467-9817.2005.00267.x.

Markham, P., & Peter L. A. (2003). The Influence of English Language and Spanish Language Captions on Foreign Language Listening/Reading Comprehension. *Journal of Educational Technology Systems, 31(3)*, 331–341. doi: 10.2190/BHUH-420B-FE23-ALA0.

Mason, L., Pluchino, P., Tornatora, M. C., & Ariasi N. (2013). An eye-tracking study of

    learning science text with concrete and abstract illustrations. *Journal of*

    *Experimental Education, 81(3)*, 356-384. doi: 10.1080/00220973.2012.727885.

Mason, L., & Tornatora M. C. (2015). Integrative processing of verbal and graphical

    information during re-reading predicts learning from illustrated text: An eye

    movement study. *Reading and Writing*, 28, 851-872. doi: 10.1007/s11145-015-

    9552-5.

Mayer, R. E., & Sims V. K. (1994). For Whom is a Picture Worth a Thousand Words?

    Extensions of a Dual-Coding Theory of Multimedia Learning. *Journal of*

    *Educational Psychology, 86(3)*, 389-401. doi: 10.1037/0022-0663.86.3.389.

Montero Perez, M., Peters, E., & Desmet P. (2014). Is Less More? Effectiveness and

    Perceived Usefulness of Keyword and Full Captioned Video for L2 Listening

    Comprehension. *ReCALL, 26(01)*, 21–43. Doi: 10.1017/S0958344013000256.

Montero Perez, M., Peters, E., & Desmet P. (2015). Enhancing Vocabulary Learning

    through Captioned Video: An Eye-Tracking Study. *The Modern Language Journal,*

    *99*, 308–28. doi: 10.1111/modl.12215.

Niehorster, D.C., Cornelissen, T.H.W., Holmqvist, K., Hooge, I., & Hessels R. (2018).

    What to expect from your remote eye-tracker when participants are unrestrained.

    *Behavior Research Methods & Instrumentation, 50(1)*, 213-227. doi:

    10.3758/s13428-017-0863-0.

Omaggio, A. C. (1979). Pictures and Second Language Comprehension: Do they Help?

    *Foreign Language Annals, 12(2)*, 107-116. doi: 10.1111/j.1944-

    9720.1979.tb00153.x.

Paivio, A. (1986). *Mental Representations: A Dual Coding Approach*. Oxford: Oxford

    University Press.

XXXXX. (2016). Incidental vocabulary acquisition from and while reading: An eye-tracking study. *Studies in Second Language Acquisition*, *38*, 97-130. doi:10.1017/S0272263115000224.

Rayner, K. (1998). Eye movements in reading and information processing: 20 years of research. *Psychological Bulletin, 124(3)*, 372–422. doi: 10.1037/0033-2909.124.3.372.

Rayner, K. (2009). Eye Movements and Attention in Reading, Scene Perception, and Visual Search. *The Quarterly Journal of Experimental Psychology,* 62(8), 1457-1506. doi: 10.1080/17470210902816461.

Renandya, W.A., & Jacobs G. M. (2016). Extensive reading and listening in the language classroom. In W.A. Renandya & H.P. Widodo (Eds.), *English Language Teaching Today: Linking Theory and Practice* (pp.97-110). Switzerland: Springer International Publishing.

Rodgers, M.P.H., & Webb S. (2017). The Effects of Captions on EFL Learners' Comprehension of English-Language Television Programs. *CALICO Journal, 33(01)*, 20-38. doi: 10.1558/cj.29522.

Roy-Charland, A., Saint-Aubin, J., & Evans M. A.. (2007). Eye Movements in Shared Book Reading with Children from Kindergarten to Grade 4. *Reading and Writing, 20*, 909–931. doi: 10.1007/s11145-007-9059-9.

Shin, J. K. (2014). Teaching Young Learners in English as a Second/Foreign Language Setting. In M. Celce.Murcia, D. M. Briton, & M. A. Snow (Eds.), *Teaching English as a Second or Foreign Language* (pp.550-567). Boston: National Geographic Learning.

Siyanova-Chanturia A., & Webb S. (2016). Teaching vocabulary in the EFL context. In W.A. Renandya and H.P. Widodo (Eds.), *English Language Teaching Today:*

*Linking Theory and Practice* (pp. 227-239). Switzerland: Springer International Publishing.

Taguchi, E., Takayasu-Maass, M., & Gorsuch G. (2004). Developing reading fluency in EFL: How assisted repeated reading and extensive reading affect fluency development. *Reading in a Foreing Language, 16(2),* 70-96.

Tang, G. (1992). The Effect of Graphic Representation of Knowledge Structures on ESL Reading Comprehension. *Studies in Second Language Acquisition, 14*, 177-195. doi: 10.1017/S0272263100010810.

Wass, S.V., Smith, T.J., & Johnson M. H. (2013). Parsing eye-tracking data of variable quality to provide accurate fixation duration estimates in infants and adults. *Behavior Research Methods, 45(1),* 229-250. Doi:10.3758/s13428-012-0245-6.

Webb, S., & Chang A. (2012). Vocabulary learninig through assisted and unassisted repeated reading. *The Canadian Modern Language Review, 68(3)*, 267-290. doi: 10.3138/cmlr.1204.1.

Webb, S., & Macalister, J. (2013), Is text written for children useful for L2 extensive eeading? *TESOL Quarterly*, 47, 300-322. doi:10.1002/tesq.70

Whitford, V., & Titone D. (2012). Second language experience modulates first- and second language word frequency effects: Evidence from eye movement measures of natural paragraph reading. *Psychonomic Bulletin and Review, 19*, 73–80. doi: 10.3758/s13423-011-0179-5.

Winke, P. M. (2013). The effects of input enhancement on grammar learning and comprehension. *Studies in Second Language Acquisition, 35(2)*, 323-352. doi: 10.1017/S0272263112000903.

Winke, P., Gass, S., & Sydorenko T. (2013). Factors Influencing the Use of Captions by

Foreign Language Learners: An Eye-Tracking Study. *The Modern Language

Journal, 97*, 254–75. doi: 10.1111/j.1540-4781.2013.01432.x.

**Appendix A.** Excerpts from two book pages and corresponding subtitles (differences highlighted in italics)

| Book: page 5 (38 words) | Video: 4 subtitles (29 words) |
|---|---|
| *Lola says,* "You ask." <br><br> *Lotta says,* "No you ask." <br><br> *So Lola says,* "Marv, can we look after Sizzles?" <br><br> *Marv says*, "Lola, do you know about dogs?" <br><br> *Lola says*, "Yes I do. Everything." <br><br> *And Lotta says*, "So do I." | *You*. You ask. No, you ask. *Go on, ask.* <br><br> Um, Marv, can we look after Sizzles? <br><br> Lola, do you know about dogs? <br><br> Yes, I do. Everything. So do I. |

| Book: page 10 (31 words) | Video: 4 subtitles (30 words) |
|---|---|
| *So Marv says, "OK*. But you do know that there are lots of rules if you want to look after Sizzles. <br><br> No chocolates. <br><br> Or cakes. <br><br> *And* no sweets of any kind. | But you do know that there are lots <br><br> of rules if you want to look after Sizzles. <br><br> *(Marv)* No chocolates. Or cakes. <br><br> *Or* no sweets of any kind. |

**Appendix B.** Comprehension questions (originally administered in Catalan)

Warm up questions:

> What is the relationship between Lola and Charlie? Are they brother and sister,
> cousins or friends?
>
> And between Lola and Lotta?
>
> And between Charlie and March?

1. Do you know whose the dog is?
2. There is an image where we can see the dog doing acrobatics, what do Lola and Lotta want to demonstrate when they imagine the dog doing acrobatics?
3. Who says that the dog can walk on two legs: the girls, Marv or Charlie?
4. At one point, Marv asks the dog to sit. Why does he ask him to sit?
5. Who proposes playing football, Charlie or Marv?
6. What problem does the other boy see?
7. Later, we can see an image of just the dog stuck to a lead. What does this image mean?
8. Before the boys go, what do Lola and Lotta promise to do?
9. What do the girls say about the dog when they are sitting at the bench?
10. According to Lola, what makes dogs happy?
11. What does Lola tell Lotta so that she lends her the lead?
12. Who knows more about dogs: Lola, Lotta or the two of them to the same extent?
13. The dog goes away because they quarrel over the lead. But why does Lotta think the dog has escaped?

14. When Marv and Charlie are back from playing football, there is something the girls do not explain to them. What is it?

15. Later Charlie shows the dog tag to Lola and Lotta. Why didn't the girls check this before?

16. Do the girls admit that they did not know about the dog tag or do they pretend they knew?

Table 1.  Length and lexical profile of video and book

|  | Book | Video |
|---|---|---|
| Length | 7'30" | 10' |
|  | 30 pages | 169 subtitles |
| Tokens | 850 | 984 |
| Lexical analysis | K1 94.1% | K1 91.2% |
|  | Type/token ratio .25 | Type/token ratio .26 |
|  | Lexical density .55% | Lexical density .54% |
|  | Adverbs 9 % | Adverbs 10% |
|  | Nouns 29% | Nouns 29% |
|  | Verbs 25% | Verbs 20% |
|  | Characters per word 3.8 | Characters per word 3.8 |
|  | Syllables per word 1.3 | Syllables per word 1.2 |
|  | Words per sentence 5.6 | Words per sentence 6.9 |

Table 2. Text characteristics by condition

|  | Condition 1 | | Condition 2 | |
| --- | --- | --- | --- | --- |
| Part 1 | Book | 393 tokens | Video | 442 tokens |
|  |  | 13 pages |  | 74 subtitles |
|  |  | 4'03'' |  | 4'32'' |
| Part 2 | Video | 516 tokens | Book | 457 tokens |
|  |  | 90 subtitles |  | 17 pages |
|  |  | 4'35'' |  | 3'20'' |
| Total |  | 909 tokens |  | 899 tokens |
|  |  | 7:52 min |  | 8:38 min |

Table 3. Selected material from the book

|  | Part 1 | Part 2 | total |
|---|---|---|---|
| All pages | 13 | 17 | 30 |
| Excluded pages | 7 | 11 | 17 |
| $N$ pages final (num. words and duration) | 6 (193) (84'05'') | 6 (160) (88'83'') | 12 (353) (172'88'') |
| $M$ num. words per page | 32.17 | 26.67 | 29.45 |
| $M$ duration per page | 14'01'' | 14'81'' | 14'41'' |

Table 4. Selected material from the video

|  | Part 1 | Part 2 | total |
|---|---|---|---|
| All subtitles | 84 | 85 | 169 |
| Excluded subtitles | 55 | 56 | 111 |
| *N* subtitles final | 29 | 29 | 58 |
| (num. words and duration) | (197) (97'29'') | (175) (104'52'') | (372) (201'83'') |
| *n* 1 line of text | 12 | 13 | 25 |
| *n* 2 lines of text | 17 | 16 | 33 |
| *M*  num. words per subtitle | 6.79 | 6.03 | 6.4 |
| *M* duration per subtitle | 3'35'' | 3'6'' | 3'48'' |

Table 5. Comprehension of the story by format

|                  | Book (n=19) M (SD) | Video (n=15) M (SD) |
|------------------|--------------------|---------------------|
| Part 1 (max. 8)  | 3.26 (1.66)        | 4 (1.93)            |
| Part 2 (max. 8)  | 2.27 (1.16)        | 2.95 (1.39)         |

Table 6. Mean (standard deviation) for the number of fixations (per page /group of

subtitles), average fixation duration and total fixation duration in text/subtitle area and

image area

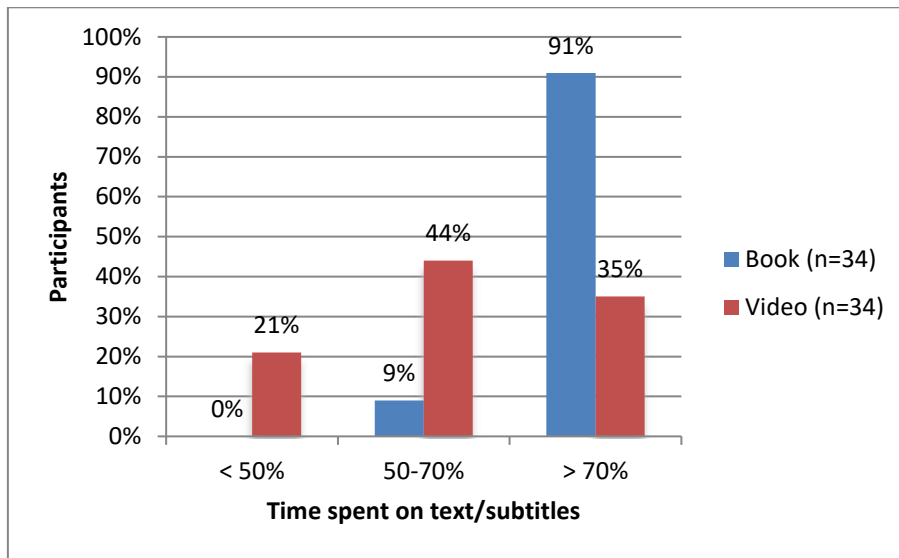| | Text/subtitle area | | | Image area | | |
|---|---|---|---|---|---|---|
| | Num. fix | Avg. fix. dur. (ms) | Tot. fix dur.(ms) | Num. fix | Avg. fix. dur. (ms) | Tot. fix dur. (ms) |
| Format | *M(SD)* | *M(SD)* | *M(SD)* | *M(SD)* | *M(SD)* | *M(SD)* |
| Book (n=34) | 34.20 (5.03) | 247 (50) | 8,445 (2,297) | 5.79 (3.25) | 205 (47) | 1,259 (880) |
| Video (n=34) | 35.96 (8.33) | 215 (33) | 7,758 (2,332) | 16.58 (6.44) | 306 (51) | 4,755 (1,835) |

Table 7. Mean (standard deviation) normalized number of fixations and normalized total
fixation duration in text/subtitle area

| Format | Normalized N of fixations text/subtitles | | | Normalized total fixation duration Text/subtitles | | |
|---|---|---|---|---|---|---|
| | *M(SD)* | Min. | Max. | *M(SD)* | Min. | Max. |
| Book (n=34) | 1.11 | 0.79 | 1.45 | 0.61 | 0.37 | 0.78 |
| | (0.16) | | | (.12) | | |
| Video (n=34) | 1.19 | 0.47 | 1.62 | 0.46 | 0.20 | 0.77 |
| | (0.28) | | | (.13) | | |

Table 8. Percentage time devoted to text/subtitles and images

| | Time in text/subtitles | | | Time in images | | |
|---|---|---|---|---|---|---|
| Format | *M(SD)* | Min. | Max. | *M(SD)* | Min. | Max. |
| Book (n=34) | 83.27% (8.6) | 58.4% | 96.04% | 16.73% (8.6) | 3.96% | 41.60% |
| Video (n=34) | 62.60% (16.13) | 29.05% | 90.44% | 37.40 (16.13) | 9.56% | 70.95% |

Figure 1. Time spent on text/subtitles

[i] Twenty-one seconds at the beginning of part 2 of the original video (including 5 subtitles and 26 words) were omitted because the scene was not included in the book format.

[ii] The pace of reading was calculated by comparing an excerpt of the story where there was no music without text.

[iii] The higher data loss in the present study is a combination of different related factors. First of all, short fixations were not merged, which leads to a higher data loss when deleting short fixations. Secondly, data was collected in remote mode. Studies have shown how the amount of data deteriorates when participants are unrestrained (e.g. Niehorster, Cornelissen, Holmqvist, Hooge, and Hessels, 2017), which may lead to more blinking and more frequent head movements. This increase movement and blinking can lead to cases of unreliable tracking or "flickery" contact with the eye-tracker (Was, Smith, & Johnson, 2013), which can result in fragmentary fixations that might be stored by the eye-tracker as multiple shorter fixations. The quality of the eye-tracker can also contribute to this unreliable contact with the eye-tracker, increasing the chances for the eye-tracker to register shorter fixations. In addition, percentages of data loss higher than 10% have been reported even in studies conducted with adults (e.g., Carroland Conklin, 2014).