



# A Universal Screening Tool for Dyslexia by a Web-Game and Machine Learning

Maria Rauschenberger<sup>1,2,3\*</sup>, Ricardo Baeza-Yates<sup>3,4</sup> and Luz Rello<sup>5</sup>

<sup>1</sup>Technology, University of Applied Science Emden/Leer, Emden, Germany, <sup>2</sup>Max-Planck-Institute for Software Systems, Saarbrücken, Germany, <sup>3</sup>Universitat Pompeu Fabra, Barcelona, Spain, <sup>4</sup>Institute for Experiential AI, Northeastern University, Boston, MA, United States, <sup>5</sup>Department of Information Systems and Technology, IE Business School, IE University, Madrid, Spain

Children with dyslexia have difficulties learning how to read and write. They are often diagnosed after they fail school even if dyslexia is not related to general intelligence. Early screening of dyslexia can prevent the negative side effects of late detection and enables early intervention. In this context, we present an approach for universal screening of dyslexia using machine learning models with data gathered from a web-based language-independent game. We designed the game content taking into consideration the analysis of mistakes of people with dyslexia in different languages and other parameters related to dyslexia like auditory perception as well as visual perception. We did a user study with 313 children (116 with dyslexia) and train predictive machine learning models with the collected data. Our method yields an accuracy of 0.74 for German and 0.69 for Spanish as well as a F1-score of 0.75 for German and 0.75 for Spanish, using Random Forests and Extra Trees, respectively. We also present the game content design, potential new auditory input, and knowledge about the design approach for future research to explore Universal screening of dyslexia. universal screening with language-independent content can be used for the screening of pre-readers who do not have any language skills, facilitating a potential early intervention.

## OPEN ACCESS

### Edited by:

Faustina Hwang,  
University of Reading, United Kingdom

### Reviewed by:

Kai Kunze,  
Keio University, Japan  
Gil Aguilar,  
Intuit, Eagle, Idaho, United States

### \*Correspondence:

Maria Rauschenberger  
maria.rauschenberger@hs-emden-leer.de

**Keywords:** dyslexia, screening tool, game, machine learning, German, Spanish, study setup, online experiment

## 1 INTRODUCTION

Dyslexia is a *specific learning disorder* which affects 5–15% of the global population (American Psychiatric Association, 2013; World Health Organization, 2010, 2019). A person with dyslexia has difficulties with reading and writing that are independent from intelligence, mother tongue, social status, or education level. Hence, people with dyslexia understand the meanings of words, but do not always know how to spell or pronounce them correctly. However, children with dyslexia do not show any obvious difficulties in other areas. This is why *dyslexia* is considered to be a *hidden* disorder. This often results in bad grades in school and frustration for the children and parents over many years. Around 40–60% of children with dyslexia show symptoms of psychological disorders (Schulte-Körne, 2010) such as negative thoughts, sadness, sorrow, or anxiety. A study showed that even if the child is diagnosed by the age of eight, they achieve lower school performance (Esser et al., 2002). Also, according to the same study, the unemployment rate for adults with dyslexia is higher. Moreover, these are common indicators for detecting a person with dyslexia.

Generally, dyslexia manifestations can be observed when children reach a certain age and literary knowledge. Current approaches to screen (pre-)readers require expensive personnel,

### Specialty section:

This article was submitted to  
Human-Media Interaction,  
a section of the journal  
Frontiers in Computer Science

**Received:** 12 November 2020

**Accepted:** 02 November 2021

**Published:** 03 January 2022

### Citation:

Rauschenberger M, Baeza-Yates R  
and Rello L (2022) A Universal  
Screening Tool for Dyslexia by a Web-  
Game and Machine Learning.  
Front. Comput. Sci. 3:628634.  
doi: 10.3389/fcomp.2021.628634

such as a professional therapist or special hardware such as fMRI scans (Paulesu et al., 2014). Previous research have studied signs of dyslexia that are not related to reading and writing such as visual perception, short-term memory, executive functions or auditory perception (Goswami et al., 2016). These signs could be used to screen potential dyslexia in pre-readers and our work shows a possible approach for doing this by using machine learning with data coming from a language-independent content integrated in a web-based game [The collected user data is online available (Rauschenberger et al., 2021b)]. Our game has the potential of being easily accessible, making parents aware the potential risk of dyslexia to further look for more help, e.g., a medical doctor or therapist.

The game and the user study is designed with the *human-centered design* framework (ISO/TC 159/SC 4 Ergonomics of human-system interaction, 2010) to collect the data set. This is relevant since collecting personal data is challenging because of privacy and trust issues (Baeza-Yates, 2018; Faraway and Augustin, 2018; Weigand et al., 2021). As a result, the final data sets are small and *small data* makes the prediction with machine learning models more difficult. That is, there is the risk of over-fitting or having a data set too small to be divided into meaningful test, training and validation sets. Hence, we followed our own recommendations coming from experience analyzing small data (Rauschenberger and Baeza-Yates, 2020a; Rauschenberger and Baeza-Yates, 2020b; Weigand et al., 2021).

We use standard machine learning classifiers like Random Forest with and without class weights, Extra Trees and Gradient Boosting from the *Scikit-learn* library for the prediction of dyslexia. Our models yields an accuracy of 0.74 and F1-score of 0.75 in German using a Random Forest and an accuracy of 0.69 and F1-score of 0.75 in Spanish using Extra Trees (Rauschenberger et al., 2020).

Historically, the rates of spelling mistakes and reading errors have been the most common way to detect persons with dyslexia, using the popular paper and pencil assessments in different languages (Cuetos et al., 2002, 2007; Fawcett and Nicolson, 2004; Grund et al., 2004). Therefore, we compare our game measures and found in our pilot study ( $n = 178$ ) four significant game measurements for Spanish, German, and English as well as eight significant game measurements for Spanish (Rauschenberger et al., 2018b), e.g., total clicks or time to first click.

Early, accurate prediction of dyslexia remains a challenge (Bandhyopadhyay et al., 2018) because dyslexia is known for causing reading and writing problems but no obvious deficits in other areas. Therefore, we need to design language-independent content fit to differentiate between children with and without dyslexia.

Another challenge is finding language-independent content that can show measurable differences between children with and without dyslexia that are comparable to differences in reading and writing mistakes. Designing language-independent content is probably the greatest challenge [also according to a report from the *National Center on Improving*

*Literacy* (Petscher et al., 2019)] because the new indicators, though related to the reading and writing difficulties, are probably not the main causes.

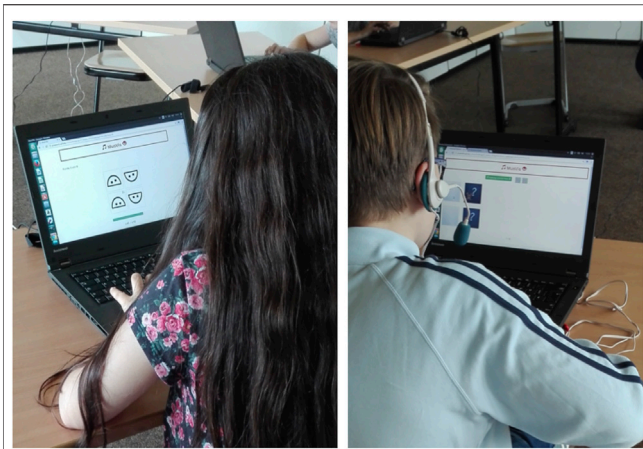
Therefore, we also share here additional **Supplement Material** of the content design, promising potential new auditory content, and knowledge about the design approach for others to use. To gather the data of this study, we had participants already diagnosed with dyslexia, instead of using pre-readers (younger children), since that would have required a long-term study. At this point, a long-term study with pre-readers would be very time-consuming, since the effort to find participants is high, participants are less likely to be diagnosed, and much time passes before results are available. An online study with readers has the advantage of reducing the effort and time required to design content, conduct various experiments for optimization, and increase the number of participants. Nevertheless, the language-independent content can be used to screen pre-readers who do not yet have any language skills. Additionally, we present the design decisions for the content creation for the auditory content and the new potential acoustic parameters that can be used in future applications. Our results show that the approach is feasible and that a higher prediction accuracy is obtained for German than for Spanish participants.

The rest of the paper is organized as follows: **Section 2** covers the related work while **Section 3** explains the rationale behind the game design. In **Section 4** we cover the methodology and in **Sections 5, 6** the predictive models and their results. We discuss the results in **Section 7**, finishing with conclusions and future work in **Section 8**.

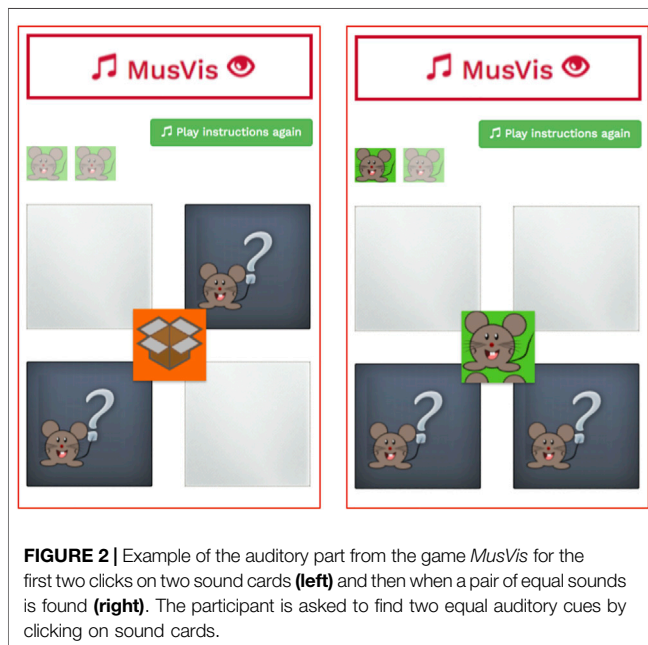
## 2 RELATED WORK

Over the last decades, dyslexia has been studied from different fields, but no scientific agreement of the causal origin has been achieved (Borleffs et al., 2019). There are two main theories at this point (De Zubicaray and Schiller, 2018). One considers visual perception (Vidyasagar and Pammer, 2010) to be a key attribute for the cause of dyslexia depending on the information processing and memory, while the other considers it to be auditory perception (Goswami, 2011).

Various applications and games to support, detect and treat dyslexia have been developed (Rauschenberger et al., 2019b). *Gamification* has been used to design various use cases, applications as well as frameworks (Hamari et al., 2014; Ritzhaupt et al., 2014; Mora et al., 2015; Seaborn and Fels, 2015; Thomas et al., 2021). Gamification designs the *game play* of games with game elements to engage and motivate users (Rouse, 2004; Rauschenberger et al., 2019c). Games are developed to screen readers (Rello et al., 2020, 2018) using linguistic content and to screen pre-readers (Gaggi et al., 2017; Geurts et al., 2015; Rauschenberger et al., 2019a, 2018a) focusing on the gameful experience. Apart from our own work (Rauschenberger et al., 2020) only *Lexa* (Poole et al., 2018) published an accuracy (89.2%) using features related to phonological processing. However, they did not include game elements, and features are collected with costly and long tests.



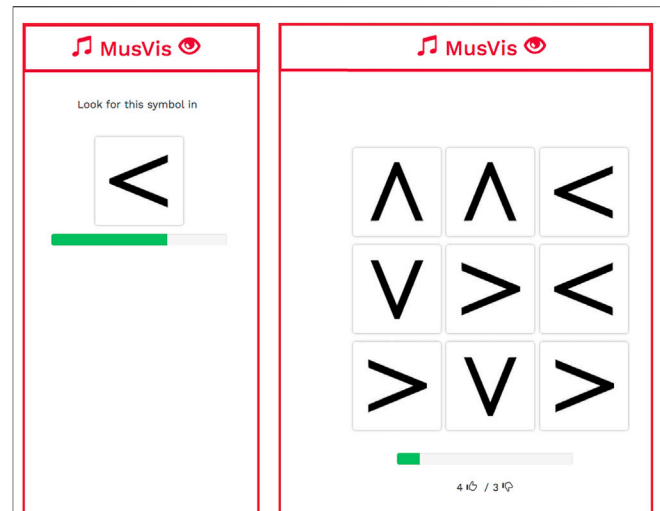
**FIGURE 1** | Participants playing the visual part (**left**) and the musical part (**right**) of MusVis. Photos included with the adults' permission.



**FIGURE 2** | Example of the auditory part from the game *MusVis* for the first two clicks on two sound cards (**left**) and then when a pair of equal sounds is found (**right**). The participant is asked to find two equal auditory cues by clicking on sound cards.

In addition, the classification is carried out on a small sample ( $n = 56$ ), without any validation and no discussion about over-fitting.

To the best of our knowledge, others have not published the details of the design decisions, or iterations of content design. Here we advance previous approaches by taking precautions on over-fitting, by not focusing on linguistic knowledge, and by using the same game content for every language, publishing our raw data (Rauschenberger et al., 2021a; Rauschenberger et al., 2021b; Rauschenberger et al., 2021c; Rauschenberger et al., 2021d). This will reduce the effort and time to design different content for different languages but more importantly, the content could be used and tested for pre-readers in different applications from different research labs.



**FIGURE 3** | Example of the visual part of the game *MusVis* with the priming of the target cue *symbol* (**left**) and the nine-squared design including the distractors for each *symbol* (**right**).

### 3 GAME DESIGN

The aim of our web-game called *MusVis* (**Figure 1**) is to measure the reaction of children with and without dyslexia while playing, in order to find differences on their behavior. A video of *MusVis* is available at <http://bit.ly/MusVisContent>. We designed our game with the assumption that non-linguistic content like rhythm or frequency (Poole et al., 2018) can represent the difficulties that a child with dyslexia has with writing and reading (Yuskaitis et al., 2015; Goswami et al., 2016), and dyslexia can be measured through the interaction of a person (Rello et al., 2020, 2018) like total number of clicks or play duration. We measure the reactions of children with and without dyslexia while playing in order to find differences in the groups' behavior. The auditory (**Figure 2**) and visual (**Figure 3**) content refers mainly to one single acoustic or visual indicator, e.g., frequency or horizontal similarity. Participants need to find the visual or auditory cue that has been shown to them before.

The game is implemented as a web application using JavaScript, jQuery, CSS, and HTML5 for the front-end, and a PHP server plus a MySQL database for the back-end. One reason for this is simplicity for remote online studies. Another reason is the advantage of adapting the application for different devices in future research studies.

We designed the language-independent game content taking into account the knowledge of previous literature selecting the most challenging content for people with dyslexia that was also easy to design in a web-game, namely auditory and visual cues.

The auditory part is shown in **Figure 2** while the visual part is shown in **Figure 3**. The game play is different due to the unequal perception of auditory and visual cues but both parts targets general skills, e.g., short-term memory (Johnson, 1980; Overy, 2000; Goswami et al., 2016), the phonological similarity effect

**TABLE 1** | Description of the auditory attributes which show promising relations to the prediction of dyslexia.

Key	Name	Description
CS	Complex vs. simple	Children with dyslexia (DG) recall significantly fewer items correctly in a lab study for long memory spans Goswami et al. (2016). The rhythmic complexity did not have an effect on the difference between DG and children without dyslexia (CG) Huss et al. (2011).
Pi	Pitch	Pitch perception is essential for prosodic performance Huss et al. (2011), is correlated to language development, and can be used as a predictor for language Yuskaitis et al. (2015).
SD	Sound duration	Acoustic parameter differences in short tones (<350 ms) are difficult to distinguish for a person with language difficulties Overy, (2000).
RT	Rise time	Both groups showed significant differences when comparing <i>rise time</i> Goswami et al. (2016). Rise time and prosodic development are strongly connected and were shown to be most sensitive to dyslexia Huss et al. (2011).
Rh	Rhythm	DG show deficits in recalling the patterns of auditory cues Overy, (2000). However, rhythm modulations show no effect on the children performance Huss et al. (2011).
STM	Short-term memory	DG show weaknesses in short-term memory tasks Overy, (2000) when more items are presented Goswami et al. (2016). Also, deficits can be frequently observed for the short-term auditory memory span Johnson. (1980).
PSE	Phonological similarity effect	DG have difficulties with similar sounds and the <i>phonological neighborhood</i> when long memory spans are used Goswami et al. (2016).
CAPS	Correlated acoustic parameters speech	Since the <i>phonological grammar</i> of music is similar to the prosodic structure of language, music (i.e., a combination of acoustical parameters) can be used to imitate these features Yuskaitis et al. (2015). DG are " <i>reliably impaired in prosodic tasks</i> " Goswami et al. (2016).

(Goswami et al., 2016), or the correlation of acoustic parameters in speech (Yuskaitis et al., 2015; Goswami et al., 2016).

As is well-known, children have more difficulty paying attention over a longer period of time. Therefore, the two parts have four stages which are counter-balanced with *Latin Squares* (Field and Hole, 2003). Each stage has two rounds, which sums up to 16 rounds in total for the whole game. Each stage first has a round with four cards and then with six cards, needing less than 10 min to play. We aim to address participants' motivation for both game parts with the design of the following game mechanics frequently used in learning environments (Rauschenberger et al., 2019c): rewards (points), feedback (instant feedback) or challenges (time limit), plus the game components (story for the game design).

The content design, user interface, game play, interaction and implementation for the auditory and visual parts of the game are described in the following sections. First, we describe the selection of content and follow with the description of the game *MusVis*, which already integrates the changes suggested after an usability test (Rauschenberger et al., 2017b).

### 3.1 Selection of Content

The selection of the content for the game is crucial, because the content links the key features extracted from previous literature connected to dyslexia into a game format. For this we need to design the game with the proper indicators (content) and game constraints in order to collect solid dependent measures that reveal differences between the participant groups. Our language-independent content to measure differences between children with and without dyslexia that represent reading and writing difficulties people with dyslexia have is shown in **Table 1**. Furthermore, this new content needs to be integrated into a game context, designed to be used as online experiment and pre-tested to avoid unintentional influences.

Previously studied language-independent indicators have been used in lab settings, which means these indicators have been tested in controlled environments. That is not the case for online

experiments. Consequently, external factors must be controlled and influences made transparent for the analysis. For example, we asked all participants to use the *Google Chrome* browser since browsers behave differently.

We decided to use an approach similar to the *Memory*<sup>1</sup> game for auditory content because of the easy and well-known gameplay for young children.

We describe the selection of content for auditory, to give an example how to inventory and select possible content (Rauschenberger et al., 2021a). The design iterations and files are available at *GitHub* (Rauschenberger et al., 2017a) and *Researchgate* (Rauschenberger et al., 2021b)<sup>2</sup>. Our goal was to reduce ideas, gameplay and acoustic parameters due to the following main requirements.

- The acoustic parameters integrated indicators strongly and significant connected to dyslexia.
- The acoustic parameters need to be easily deployed in a web-game.
- The acoustic parameters need to be easily deployed in the gameplay.
- The game duration fits pre-readers attention span.

We present the different iterations starting with the first iteration where we collected ideas (possible acoustic parameters connected to dyslexia) from literature (Rauschenberger et al., 2021d). We selected with a semi structured literature review the literature relevant to find indicators related to auditory difficulties (Rauschenberger et al., 2021a). We found the first core paper like (Overy, 2000) and looked into similar wording and publication for the second

<sup>1</sup>Pairs of identical cards (face down) must be identified by flipping them over (Wikipedia, 2019)

<sup>2</sup>The musical content used in the final game *MusVis* is available at <https://github.com/Rauschii/DysMusicMusicalElements>.

**TABLE 2** | Description of auditory parameters.

Participant features	Description
1 Implementation Priority	Our priority to implement this option.
2 Main-Round ID	ID for rounds with the same acoustic parameter, e.g., frequency, rhythm.
3 Sub-Round ID	Additional ID to distinguish different main-rounds with different settings, e.g., different <i>difficulty level</i> like <i>easy</i> vs. <i>difficult</i> .
4 Difficulty Level	It indicates depending mainly on the amount of cards how difficult this content is.
5 Instructions	Short description of what the participants should do in English and German
6 Input Description	Short description of the parameters of the auditory elements for this round.
7 Auditory Parameters Part	It indicates which acoustic parameters are considered.
8 Feedback Loop Example	It shows example feedbacks from the domain experts and researchers for the round.
9 Interaction	It is a description of how the child should interact with the game and content.
10 Reason	It is a short description why the researcher thinks this round will work for the goal.
11 Citation Key	It is an example of the citation key to point to the literature we use as baseline. The final connections are presented in <b>Section 3.2</b> .

iteration. Next, we explored acoustic parameters with a strong connection to dyslexia in other lab studies and redefine the acoustic parameter collection (see iteration two in the **Supplementary Material**). For example, we added new acoustic parameters ideas such as presenting one frequency at different times on one side of the ear and then ask “*Where did you hear the frequency first?*”. Also, we included the first ideas of a gameplay and how to make game rounds. How we came up with the game round was a very unsupervised creative approach. In the third and last iteration included the main requirements mentioned above.

The main parameters to describe our auditory features picked are detailed in **Table 2** and elaborate here on important decisions. Iteration three has promising indicators not tested in a game environment, yet (Rauschenberger et al., 2021d). Example indicators are the lab sound (Huss et al., 2011), volume level, or timing. Also, we describe promising game rounds such as “Find the same sound” behind a card but the card is making a sound from a certain direction and needs to be found. Another game round could be “Which sound came first?” that schedules the same sound on different timing for the left and right ear.

The consideration of *difficulty-level* is important as we want to address pre-readers that can be easily overwhelmed by information. Hence, *MusVis* implements only single acoustic indicators and only up to three different choices (one choice equals 2 cards). We decided against a game round with only one choice (two cards) as this would increase finding the correct answer by chance (50% to be correct) and a change of gameplay.

We included a short description (*instructions*) in different languages already in the phase of *collection of content* to ensure the feasibility and consistency between user studies in different languages. We recommend to make first one draft of the game and description in one language and then the translation into other languages. Reason is the iterative creative process until a first draft is reached that would make also a lot of changes to the text. But due to the differences of languages designer should not wait too long as German has longer names compared to English and therefore titles and descriptions need more space in the game design. The space needed for different language needs to be taken into account for pre-testing.

The short description of the input and short description of acoustic parameters is to ensure a simple overview to other rounds when deciding between different sounds. It might be usefully to separate the information from these two columns in the future to reduce redundant information.

The feedback loop is used between researchers and an expert in the creation of MP3 files to verify the artifact with the goal we have: Finding differences in the interaction behavior of groups with and without dyslexia when playing our game with this content. We shortly explain here the dyslexia auditory theory that was also mentioned in *Section 2*. Researchers argue that dyslexia might be mainly based on phonological and perception differences (Goswami, 2011). Moreover, previous research has related speech perception difficulty to auditory processing, phonological awareness, and literacy skills (Tallal, 2004; Rolka and Silverman, 2015; De Zubicaray and Schiller, 2018). Phonological deficits of dyslexia have also been linked to basic auditory processing (Hämäläinen et al., 2013). However, there are musicians with dyslexia who scored better on auditory perception tests than the general population (Männel et al., 2017). At the same time, these participants score worse on tests of auditory working memory, i.e., the ability to keep a sound in mind for seconds. This observation is in line with the results on perceptions for short duration sounds (Huss et al., 2011) and the findings on the *prosodic similarity effects* of participants with dyslexia (Goswami et al., 2016). Still, it is challenge to design auditory cues connection to the auditory perception that can be used in our gameplay and not measuring the musical knowledge or hearing range. Hence, we included knowledge from an domain expert and an example feedback is presented in *feedback loop*. An example of screenshots with the parameters for different stages is available on *Researchgate* (Rauschenberger et al., 2021c).

We hope other research working on language-independent screening find our collection of auditory indicator, insight about the design decision and how we inventoried our content useful for their own research and present next our selection for the game *MusVis*.

### 3.2 Auditory Game Design

The auditory part is inspired in the traditional game *Memory* in which pairs of identical cards (face down) must be identified by flipping them over (Wikipedia, 2019). We chose this game play

**TABLE 3** | Mapping of the evidence from literature to distinguish a person with dyslexia, the attributes and general assumptions, and the stages of the game *MusVis*.

Attributes	Auditory				General			
	CS	Pi	SD	RT	Rh	STM	PSE	CAPS
Literature								
Goswami et al. (2016)	✓			✓		✓	✓	✓
Huss et al. (2011)	✓	✓		✓	✓			
Johnson. (1980)						✓		
Overy, (2000)			✓		✓	✓		
Yuskaitis et al. (2015)		✓						✓
Stage								
Frequency	✓	✓	✓			✓	✓	✓
Length			✓			✓	✓	✓
Rise time	✓		✓	✓		✓	✓	✓
Rhythm	✓		✓		✓	✓	✓	✓

because it is a well-known children game and could be easily transformed to use auditory cues. To create the auditory cues, we used acoustic parameters; for example, to imitate the *prosodic* structure of language which is similar to the *phonological grammar* of music (Port, 2003).

Musicians with dyslexia score better on auditory perception tests than the general population, but not on auditory working memory tests (Männel et al., 2016). Auditory working memory helps a person to keep a sound in mind. We combined, for example, the deficits of children with dyslexia in auditory working memory with the results on the short duration of sounds (Huss et al., 2011) while taking the precaution of not measuring hearing ability (Fastl and Zwicker, 2007). Each stage is assigned to one acoustic parameter like frequency or rhythm which is designed with the knowledge of the analysis from previous literature (Rauschenberger et al., 2018b, 2021a,d).

Therefore, we used the acoustic parameters *frequency*, *length*, *rise time* and *rhythm* as auditory cues. Each auditory cue was assigned to a game stage (Table 3), which we mapped to the attributes and literature references (Table 1) that provide evidence for distinguishing a person with dyslexia.

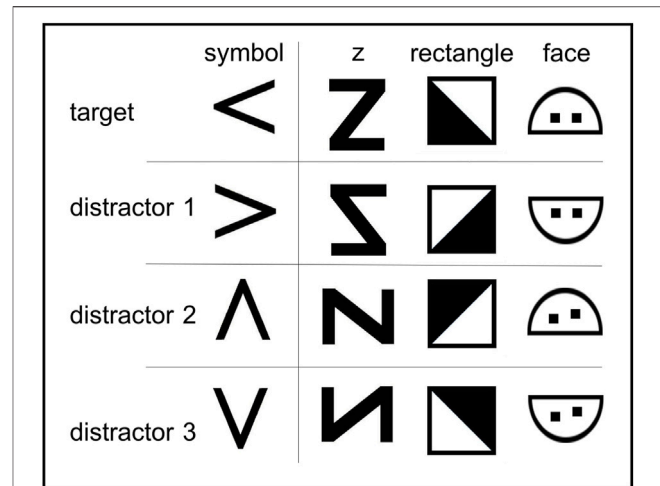
For example, our *rhythm* stage uses the following characteristics: *complex vs. simple* (Huss et al., 2011; Goswami et al., 2016), *sound duration*, *rhythm* (Huss et al., 2011), *short-term memory* (Johnson, 1980; Goswami et al., 2016), *phonological similarity effect* (Goswami et al., 2016), and *correlated acoustic parameters speech* (Yuskaitis et al., 2015; Goswami et al., 2016).

Each acoustic stage has three auditory cues (we use MP3 for sound files). Each stage is assigned to one acoustic parameter of sound, which is designed with knowledge of the analysis from previous literature (e.g., frequency or rhythm).

The auditory cues are generated with a simple sinus tone using the free software *Audacity*<sup>3</sup>. The exact parameters of each auditory cue are already published (Rauschenberger et al., 2018b) and the auditory cues are available at *GitHub* (Rauschenberger et al., 2017a)<sup>4</sup>. Each stage has two rounds,

<sup>3</sup>*Audacity* is available at <http://audacity.es/>, Last access: May 2019

<sup>4</sup><https://github.com/Rauschii/DysMusicMusicalElements>



**FIGURE 4** | Overview of the designed visual cues. The figure shows the target cue (**top**) and distractor cues (**below**) for the four different stages (*z*, *symbol*, *rectangle*, *face*) of the visual part of the game *MusVis*.

with first two and then three auditory cues that must be assigned by choosing the same sound (see Figure 2). The arrangement of sounds (which auditory cue matches which card) is random for each round.

### 3.3 Visual Game Design

The visual game play uses a Whac-A-Mole interaction similar to the first round of *Dyctective* (Rello et al., 2020). But instead of using letter recognition as does *Dyctective*, we used language-independent visual cues. An example for letter recognition would be finding the graphical representation of the letter /e/. We adapted the interaction design and content for this purpose (Figure 3). For the visual game, we designed cues that have the potential of making more cues with similar features and represent horizontal and vertical symmetries that are known to be difficult for a person with dyslexia in different languages (Vidyasagar and Pammer, 2010; Rello et al., 2016a; Rauschenberger et al., 2016).

To create the visual cues, we designed different visual representations similar to visual features of annotated error words from people with dyslexia (Vidyasagar and Pammer, 2010; Rello et al., 2016a; Rauschenberger et al., 2016) and designed the game as a simple search task, which does not require language acquisition.

In the beginning, participants are shown the target visual cues (see Figure 3, left) for 3 seconds. They are asked to remember this visual cue. After that, the participants are presented with a setting where the target visual cue and distractors are displayed (see Figure 3, right). The participants try to click on the target visual cue as often as possible within a span of 15 s. The arrangement of the target and distractor cues randomly changes after every click.

The visual part has four stages, which are counter-balanced with *Latin Squares* (Field and Hole, 2003). Each stage is assigned to one visual type (*symbol*, *z*, *rectangle*, *face*) and four visual cues for each stage are presented. One visual cue is the target, which the participants need to find and click (see Figure 4, top). The

other three visual cues are *distractors* for the participants. Each stage has two rounds with first a 4-squared and then a 9-squared design (see **Figure 3**, right). The target and all three distractors are displayed in the 4-squared design. In the 9-squared design, the target is displayed twice as well as distractors two and three. Only distractor one is displayed three times.

## 4 USER STUDY METHODOLOGY

We use the human-centered design framework to design our study and to collect the data for the prediction of dyslexia. We conducted a within-subject design study ( $n = 313$ ) which means that all participants played all game rounds (Field and Hole, 2003) with the same language-independent content. Only the game instructions were translated into each native language.

Spanish participants diagnosed with dyslexia were mainly recruited from public social media calls by non-profit organizations. We recruited German participants diagnosed with dyslexia mainly over support groups on social media. Also, some English speakers contacted us through this call as our location is international. The control groups for Spanish and German were recruited mostly with the collaboration of four schools, two in each country.

### 4.1 Online Data Collection

Collecting data is costly in terms of time consumption and privacy issues, especially if the data is related to education and health. Therefore, we must make the best of the limited resource (Baeza-Yates, 2018; Faraway and Augustin, 2018; Rauschenberger and Baeza-Yates, 2020a; Weigand et al., 2021). In our case, we need a certain age range to make sure a person with dyslexia is already diagnosed and has not been fully treated yet. Since our collected data is considered *small data* (Baeza-Yates, 2018; Faraway and Augustin, 2018), we need to analyze them accordingly, i.e., avoid over-fitting using cross-validation instead of training, test and validation sets as well as using classifiers configured to avoid over-fitting.

### 4.2 Procedure and Ethics Statement

First, the parents were informed about the purpose of the voluntary study. Next, only after the parents gave the consent, children were allowed to participate in this user study from home or from school, with the first author of this work present or always available through digital communication.

The data collection for this user study has been approved by the German Ministry of Education, Science and Culture in Schleswig-Holstein (*Ministerium für Bildung, Wissenschaft und Kultur*) and Lower Saxony State Education Authority (*Niedersächsische Landesschulbehörde*). In Spain governmental approval was not needed in addition to the school approval.

If the study was conducted in a school or learning center, the parents or the legal guardian consent was obtained in advance and the user study was supervised by a teacher or therapist. After the online consent form was approved, we collected demographic data which was completed by the participant's supervisor (e.g.,

**TABLE 4** | Overview of the participants per data set.

Data set	N	Dyslexia (DG)			
		n	age	female	male
DE	149	59	10.22	21	38
ES	153	49	9.47	26	23
ALL	313	116	9.77	50	66

Data set	N	Control (CG)			
		n	age	female	male
DE	149	90	9.58	42	48
ES	153	104	9.99	58	46
ALL	313	197	9.76	103	94

parent/teacher), including the age of the participant, the dyslexia diagnosis (yes/no/maybe) and the native language. We ask the participant's supervisor to only say *YES* for a participant if the child had an official diagnosis, for example from an authorized specialist or a medical doctor.

After that participants played both parts of the game. At the end, two feedback questions are asked and the participant's supervisor could leave contact details to be informed about the results of the study. Personal information of the participant's supervisor such as name or email is not published and is stored separately from the participants data, if given. On the other hand, the name of the child is not collected and all data is stored on a password secured web server.

### 4.3 Participants

The data includes only participants that completed all 16 rounds of the web game using a computer or a tablet. Dropouts happened mostly because participants used a different browser (e.g., *Internet Explorer* instead of *Google Chrome*) or a different device (tablet instead of a computer).

For the predictive models, we took 313 participants into account, including the 178 participants from the pilot study (Rauschenberger et al., 2018b). To have precise data, we took out participants that reported in the background questionnaire that they suspected of having dyslexia but did not have a diagnosis ( $n = 48$ ).

The remaining participants were classified as diagnosed with dyslexia (DG) or not showing any signs of dyslexia (control group, CG), as reported in the background questionnaire.

We separated our data into three data sets: one for the Spanish participants (ES,  $n = 153$ ), a second for the German participants (DE,  $n = 149$ ), and one for all languages (ALL,  $n = 313$ ) in which we included participants that spoke English ( $n = 11$ ). Participants ranged in age from 7 to 12 years old. The users in the data sets are described in **Table 4**.

Participants played the game either in English, German or Spanish depending on their native language. We had some bilingual participants ( $n = 48$ ) in the Spanish data set (Spanish and Catalan) since the media call was done from the non-profit organization *ChangeDyslexia*<sup>6</sup>. For these cases, we used the

<sup>6</sup><https://changedyslexia.org/>

**TABLE 5** | Description of participant features.

Participant features	Description
1 Age	It ranges from 7 to 12 years old.
2 Gender	It is a binary feature, either with a <i>female</i> or <i>male</i> value.
3 Language	It is either <i>Spanish</i> , <i>German</i> or <i>English</i> .
4 Native Language	It indicates if the language used for the instructions is the first language of the participants, being <i>Yes</i> , <i>No</i> or <i>Maybe</i> .
5 Instrument	It indicates if a participant plays a musical instrument, being <i>No</i> , <i>Yes, less than 6 months</i> or <i>Yes, over 6 months</i> .
6 Memory	It indicates how well the participant knows the visual <i>Memory</i> game, being <i>Participant gave no answer</i> , <i>Participant does not know the game</i> , <i>Played once</i> , <i>Played a few times</i> or <i>Played a lot</i> .
7 Rating Auditory Part	It indicates the self-reported answer with a 6-level <i>Likert scale</i> Field and Hole, (2003) to the statement: "the auditory part was easy for the participants." The values are <i>Answer unknown</i> , <i>Strongly disagree</i> , <i>Disagree</i> , <i>Undecided</i> , <i>Agree</i> or <i>Strongly Agree</i> .
8 Rating Visual Part	It indicates the self-reported answer of the statement: "the visual part was easy for the participants." (same <i>Likert scale</i> from feature 7).
9 Device	It is the device the participants used and is a binary feature with the value <i>Computer</i> or <i>Tablet</i> .

**TABLE 6** | On the left are features 10 to 105 for the auditory part and on the right are features 106 to 201 for the visual part of the game MusVis.

Auditory features	Visual features
<b>10–17</b> Time to click	<b>106–113</b> Time to click
<b>18–25</b> Total clicks	<b>114–121</b> Total clicks
<b>26–33</b> Duration per round	<b>122–129</b> Correct answers
<b>34–41</b> Duration interaction	<b>130–137</b> Wrong answers
<b>42–49</b> Average click time	<b>138–145</b> Accuracy
<b>50–57</b> Logic	<b>146–153</b> Efficiency
<b>58–65</b> 2nd click interval	<b>154–161</b> 2nd click interval
<b>66–73</b> 3rd click interval	<b>162–169</b> 3rd click interval
<b>74–81</b> 4th click interval	<b>170–177</b> 4th click interval
<b>82–89</b> 5th click interval	<b>178–185</b> 5th click interval
<b>90–97</b> 6th click interval	<b>186–193</b> 6th click interval
<b>98–105</b> Instructions	<b>194–201</b> Time last click

language they reported to be more comfortable with, which was used for the instructions of the game. We do not use the native language, but rather the language the game was played in as the criterion to split the data sets for three reasons. First, the definition of a native language or mother tongue can be made easily when a participant speaks only one language. But this is not the case for bilingual participants because they might not be able to choose, and then we cannot distinguish the mother tongue or native language clearly (Kecskes and Papp, 2000). Second, this question is a self-reported question and every participant's supervisor might define it differently for each child. Finally, some bilingual speakers spoke similar Latin languages (Spanish and Catalan). We consider these participants in the ES data set, as the instructions of the game were in Spanish.

#### 4.4 Dependent Variables and Features

The participant features are detailed in **Table 5** while the dependent variables collected through the game are listed in **Table 6**. These variables were used for the statistical comparison of the pilot study and for the selection of the features for the predictive models. Feature three was set with the language selected for the instructions. Features 1, 2, 4 to 8 were answered with the online questions by the participants' supervisor. Feature 9 was collected from the browser during the study experiment.

We used the following dependent variables for the statistical comparison:

##### Auditory game part

- *Duration round* (milliseconds) starts when round is initialized.
- *Duration interaction* (milliseconds) starts after the player clicks the first time on a card in each round.
- *Average click time* (milliseconds) is the duration of a round divided by the total number of clicks.
- *Time interval* (milliseconds) is the time needed for the second, third, fourth, fifth and sixth clicks.
- *Logic* we define it as *True* when in a round the first three clicked cards are different, otherwise, it is *False*.
- *Instructions* is the number of times the game instructions were listened by the player.

##### Visual game part

*Number of hits* is the number of correct answers. *Number of misses* is the number of incorrect answers. *Efficiency* is the number of hits multiplied by the total number of clicks. *Accuracy* is the number of hits divided by the total number of clicks.

##### Both parts

- *Time to the first click* (milliseconds) is the duration between the round start and the first user click.
- *Total number of clicks* is the number of clicks during a round.

We would like to further elaborate on the game measurement *Logic*, which is based on the direct experience of the user study. Some children may not have *really listened* to the sounds and played *logically*. As each round is designed such that the first two clicks never match, if the participant chooses for the third click a different card, s/he is increasing the chances of finding a match independent of the total amount of cards.

The descriptions of the participant features are in **Table 5**. The features for the data sets ALL, ES, and DE are the same. Each data



set has 201 features per participant, where features 10 to 105 are the variables from the auditory part and features 106 to 201 are the variables from the visual part.

## 5 PREDICTIVE MODELS

In this section we present the machine learning techniques used for the data sets ALL ( $n = 313$ ), ES ( $n = 153$ ), and DE ( $n = 149$ ). First, we explain the choice of predictive models and then the feature selection.

### 5.1 Model Selection

We used Random Forest (RF), Random Forest with class weights (RFW), Extra Trees (ETC), Gradient Boosting (GB), and the Dummy Classifier (Baseline), which are described in the Scikit-learn version 0.21.2 (Scikit-learn Developers, 2019). We address the risk of over-fitting on our small data sets with 10-fold cross-validation and the default parameters suggested in the Scikit-learn library to avoid training a model by optimizing the parameters specifically for our data (Scikit-learn Developers, 2019). As we have small data, we are not optimizing the input parameters of classifiers until we can hold out a test data set as proposed by scikit-learn 0.21.2 documentation to evaluate the changes (Scikit-learn, 2019) and to avoid biases (Varma and Simon, 2006). To explore the best prediction conditions, we used the feature selection as described in the next section.

### 5.2 Informative Features

We address the danger of selecting the correct features (Jain and Zongker, 1997) by taking into account the knowledge of previous literature about the differences of children with an without dyslexia. For example, since there are two theories of the cause of dyslexia [visual vs. auditory (De Zubicaray and Schiller, 2018)], we use subsets of visual and auditory features to explore the influence on the classifiers.

We rank the most informative features with *Extra Trees*. The results show a flat distribution for all three data sets and a step at the information score of 0.008: ALL ( $n = 33$  features), ES ( $n = 41$  features), and DE ( $n = 38$  features). The comparison of the most informative features reveals that the data sets have only a few features in common, e.g., four features for Spanish and German (Logic, sixth click interval, total clicks, duration interaction) or only 16 features in ALL compared to Spanish and German. Visual and auditory features are equally represented in the ranking of the most informative features; for example, ALL has 16 auditory features and 14 visual features.

The biggest step in the informative ranking for all three data sets is between the fifth and sixth informative features, e.g., for ALL the step is between the visual part (cue Z, 4 cards) *Efficiency* with the informative score of 0.0128 and the auditory part (cue *Rhythm*, 6 cards), *Time fifth click* with a score of 0.0104. The only dependent variables with the same tendency are *Number of misses* and *Total clicks* from the visual game part, but the features from the different rounds for the different data sets are mainly not

under the 33 informative features (ALL 2/16, ES 3/16 and DE 6/16).

## 6 RESULTS

We followed the same steps of the pilot study to compare the statistical findings before giving the machine learning results.

### 6.1 Statistical Validation

The pilot study collected data from 178 participants (which were later included into our current data set,  $n = 313$ ) to find significant differences on the game measurements (Rauschenberger et al., 2018b). Therefore, we apply first the *Shapiro-Wilk Test* and then the *Wilcoxon Test* since all game measures are not normally distributed. We use the Bonferroni correction ( $p < 0.002$ ) to avoid type I errors. We present the results of the statistical analysis for the validation data ( $n = 313$ ) separated by language and for all languages (see **Table 7**). Additionally, we compare the statistical analysis results from the pilot-study ( $n = 178$ ) with the new data set ( $n = 313$ ).

The ES data set ( $n = 153$ ) has seven dependent variables with significant differences between groups: *fourth click interval*, *duration round*, *average click time*, *total number of clicks*, *time to the first click*, *number of hits*, and *efficiency*. The ES data set ( $n = 153$ ) confirmed the results of the pilot study ( $n = 178$ ). All other game measurements decreased the significance by slightly increasing the  $p$ -value (visual efficiency from  $4e - 5$  to  $1e - 4$ ). The data set ES has seven significant variables that distinguish a person with or without dyslexia.

For the data set ALL ( $n = 313$ ) we consider only dependent variables with the same tendency as for the pilot study ( $n = 178$ ). We categorize the tendency (e.g., *playing faster or having more clicks*) by the group (dyslexia compared to control group) *mean* of the dependent variables within the same language. ALL ( $n = 313$ ) has two visual game measurements (*number of misses* and *total clicks*) with the same tendency while the pilot study had five for the visual game (*total clicks*, *time to the first click*, *hits*, *accuracy*, and *efficiency*).

The DE data set ( $n = 149$ ) confirmed the results of the pilot study ( $n = 57$ ) with no significant dependent variables. The *means* of the dependent measurements for DE are all very close (e.g., the *time to the first click* is 2.58s for the control group and 2.50s for the dyslexia group).

We can confirm that misses did not reveal significant differences for German or Spanish, even though the tendency is now the same for both languages. On the other hand, the total number of clicks is still significant.

To sum up, we confirmed one significant dependent variable in ALL ( $n = 313$ ), seven significant dependent variables for ES ( $n = 153$ ), and no significant dependent variables for DE ( $n = 149$ ).

### 6.2 Predictive Results

We processed our data sets with different classifiers and different subsets of features, following the description from the previous section. We follow our criteria for analyzing small (tiny) data to avoid wrong results as wrong results have a negative impact on a

**TABLE 7** | Overview of dependent variables for visual (top) and auditory (below) features of *MusVis*.

Part	Data set	Variable	Control		Dyslexia		Mann-Whitney U		
			Mean	sd	Mean	d	W	p-value	Effect size
Visual	ALL	Total clicks	6.8	2.7	7.2	3.2	670194	<b>2e-04</b>	0.14
		Misses	1.2	2	1.3	2.7	713627	0.14	0.05
	ES	Total clicks	6.8	2.7	7.7	3	132207	<b>3e-08</b>	<b>0.31</b>
		First click	2.63s	1.69s	2.26s	1.22s	141938	<b>1e-04</b>	0.27
		Hits	5.8	3	6.5	2.9	136904	<b>2e-06</b>	0.25
		Misses	1	1.7	1.2	2.7	157086	0.12	0.07
		Accuracy	0.82	0.27	0.85	0.26	153012	0.03	0.10
		Efficiency	3.1s	2.6s	2.75	2.4s	142162	<b>1e-04</b>	0.14
	DE	Total clicks	6.7	2.6	6.8	3.3	169439	0.47	0.03
		First click	2.50s	1.32s	2.58s	1.56s	168932	0.43	0.06
		Hits	5.4	2.6	5.3	2.8	164224	0.16	0.05
		Misses	1.3	2.1	1.5	2.8	166140	0.24	0.09
		Accuracy	0.81	0.27	0.78	0.29	165688	0.22	0.08
		Efficiency	3.2s	2.4s	3.5s	2.9s	167288	0.33	0.10
	Auditory	ES	Total clicks	11.3	6	10.9	5.5	157282	0.15
4th click			2.0s	1.3s	1.7s	1.0s	131228	<b>1e-08</b>	0.29
6th click			1.7s	1.1s	1.6s	0.9s	152772	0.04	0.15
Duration			32.6s	69.9s	24.7s	18.2s	142726	<b>2e-04</b>	0.19
Average			3.0s	2.7s	2.6s	0.9s	121966	<b>5e-13</b>	0.29
DE		Total clicks	11.1	5.5	11.5	6.6	166340	0.27	0.07
		4th click	1.9s	1.0s	2.0s	1.0s	167184	0.32	0.01
		6th click	1.8s	0.8s	1.9s	1.3s	163076	0.12	0.12
		Duration	27.1s	18.6s	29.4s	22.9s	163994	0.15	0.11
		Average	2.7s	0.8s	2.8s	1.0s	166194	0.26	0.11

Significant results are in bold.

**TABLE 8** | Best results of the different classifiers, features and data sets. Results are ordered by the best F1-score and accuracy.

Model	Data	Feat.	Recall	Precis.	F1	Acc.
<b>RF</b>	<b>DE</b>	5	0.77	0.78	<b>0.75</b>	<b>0.74</b>
RFW	DE	5	0.75	0.75	0.74	0.73
Baseline	DE		0.60	0.37	0.46	0.50
<b>ETC</b>	<b>ES</b>	20	0.76	0.76	<b>0.75</b>	<b>0.69</b>
RF	ES	5	0.74	0.73	0.72	0.65
Baseline	ES		0.68	0.46	0.55	0.50
<b>GB</b>	<b>ALL</b>	20	0.66	0.65	<b>0.65</b>	<b>0.61</b>
GB	ALL	5	0.64	0.64	0.63	0.59
Baseline	ALL		0.63	0.40	0.49	0.50

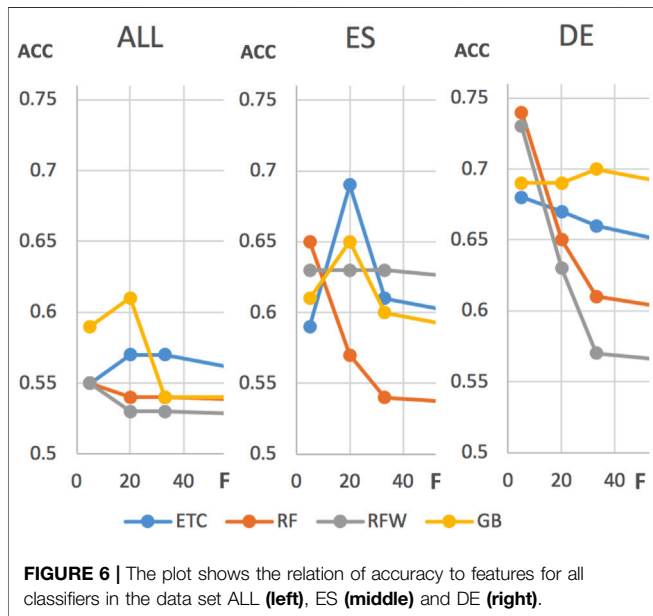
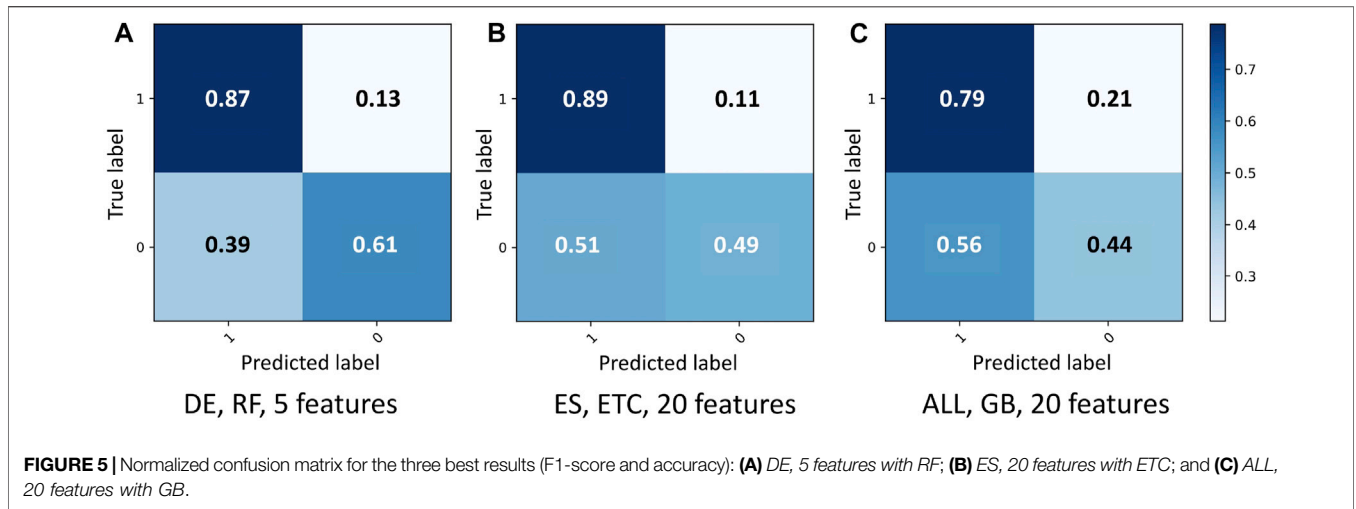
person such as missing out a person with dyslexia (Rauschenberger and Baeza-Yates, 2020a; Rauschenberger and Baeza-Yates, 2020b).

We computed the *balanced accuracy* for our binary classification problem to deal with imbalanced data sets; for example, the ALL data set has dyslexia 37% vs. control 63%. The Dummy Classifier is computed for our imbalanced data with the most frequent label and reported with the balanced accuracy (Scikit-learn Developers, 2019). We do not apply over- or under-sampling to address our imbalanced data because the variances among people with dyslexia are broad, for example, difficulty level or the individual causes for perception differences.

As described in the previous section the ranking of the informative features is different for the three data sets. Hence, we explore the influence of different subsets of features, namely: 1) all represented features (201 features); 2) the 5 most informative features; 3) the 33 most informative features, as this was the next natural informative subset; 4) 20 random features selected from (3); and 5) 27 features that have the same tendency and which have been answered by the participants' supervisors, because they are mainly not under the most informative feature subsets (although *total clicks* is significant in the statistical comparison).

We report the two best F1-scores and *balanced accuracy* scores for each data set as well as the baseline, as can be seen in **Table 8**. We outperform our baseline for all data sets. The best F1-score, *0.75*, is achieved for both languages, the DE and ES data sets. DE uses 5 features with RF and ES uses ETC with 20 features. The second best F1-score, *0.74*, is achieved with the DE data set using 5 features and RFW. The best accuracy, *0.74*, is achieved with RF while the second best of *0.73* is achieved with RFW, both in the DE data set using just 5 features.

For ES, the best F1-score is also *0.75* with ETC and the selection of 20 features. The second best F1-score for ES is *0.72* with RF and a selection of 5 features. The F1-score is reduced by 0.1 when combining the two data sets (DE and ES), since the best F1-score for ALL is *0.65* using GB and 20 features. The second best F1-score for ALL is *0.63* with GB and 5 features. For ES, the best accuracy is *0.69* with ETC and the



selection of 20 features. The second best accuracy for ES is 0.65 with RF and a selection of 5 features. The accuracy is reduced by nearly 0.1 when combining the two data sets (DE and ES), since the best accuracy for ALL is 0.61 using GB and 20 features. The second best accuracy for ALL is 0.59 with GB and 5 features. This shows that there are differences across languages.

The normalized confusion matrix (see Figure 5) does not show over-fitting for the best results for DE, ES and ALL. The fact that the best results are with few features imply that the rest are highly correlated or noisy.

The reduction of features improves the accuracy for DE but not consistently for ES and ALL, as can be seen for the different classifiers and data sets in Figure 6. For example, reducing the features for DE improves the accuracy for ET, RF, and RFW, but not for GB. For ES, the accuracy improves only for RF and stagnates for RFW when reducing the number

of features, otherwise the accuracy inverts for ETC and GB. For the data set ALL, RFW and RF improve but ETC and GB decrease.

## 7 DISCUSSION

Most children with dyslexia show a varying severity of deficits in more than one area (Black et al., 2016), which makes dyslexia more a spectrum than a binary disorder. Additionally, we rely on current diagnostic tools (e.g., DRT (Grund et al., 2004; Steinbrink and Lachmann, 2014)) to select our participant groups, which do not yet represent the diversity of people with dyslexia. We accept that our participants have a high variance because of the measurement of our current diagnostic tools and the spectrum that dyslexia have.

### 7.1 Group Comparison

The measurement data taken from the game *MusVis* show that Spanish participants with dyslexia behave differently than their control group. Differences can be reported for the auditory game part for: *fourth click interval, duration, and average click time*. For the visual part, the following measurements can be reported as indicators: *total clicks, time to the first click, hits, and efficiency*.

We can show with our results over all languages that the effect for each measurement is confirmed even if we cannot draw strong conclusions about our sample size on the comparison of German vs. Spanish speaking participants. Spanish had eight significant indicators in the pilot study and we expected to reproduce the same number of significant indicators with more German participants.

In general, all participants found the game easy to understand, and only children at the age of 12 complained about missing challenges. The amount of positive feedback and engagement of all age groups let us conclude that the game mechanics and components applied are also positive for perceiving *MusVis* as a game and not as a test.

Dyslexia is known to be present across different languages and cultures (American Psychiatric Association, 2013). The assumption that the tendencies for the indicators are similar over all languages cannot (yet) be proven for all indicators in our study (e.g., German participants with dyslexia start to click faster than the Spanish participants compared to their language control group in the auditory part). We can exclude external factors such as different applications or study setups as possible influences on this opposite tendency. According to the results, we may have to assume that not all indicators for dyslexia are language-independent and that some have cultural dependencies, or we have *omitted variable bias*. To confirm this assumption, we will need to obtain larger numbers of participants for both language groups (Spanish and German) or investigate further measurements (indicators).

The variables *time to first click (visual and auditory)* and *total number of clicks (visual and auditory)* provide dependencies of the game content and game design. Otherwise, we could not explain the trend difference between the auditory and visual parts for *total number of clicks* (i.e., *total clicks* for visual is significantly different than for auditory). Additionally, the analysis of the auditory game part presents one limitation: participants could select a correct pair by chance, e.g., participants could click through the game board without listening to the sounds.

Children with dyslexia are detected by their slower reading or spelling **error rate** (Schulte-Körne et al., 1996; Coleman et al., 2008). Therefore, we designed our game with content that is known to be difficult for children with dyslexia to measure the errors and duration. Nevertheless, from previous literature we knew that children with dyslexia do not make more mistakes in games than the control group (Rello et al., 2020). We can confirm that *misses* did not reveal significant differences for German or Spanish either. It might be possible that we cannot compare errors in reading and writing with errors in this type of game. Then, we cannot explain (yet) why the Spanish control group made more mistakes than the Spanish group with dyslexia. It might also be possible that participants with dyslexia show generally different behavior that is separated from the content but depends on the *game play*.

Spanish children without dyslexia take significantly more time to find all pairs and finish the auditory game part. Children without dyslexia take more time before they *click the first time* (visual) for all languages. This might be due to the time they need to process the given auditory information (Tallal, 2004) or recall the auditory and visual information from short-term memory (Goswami et al., 2016). However, participants with dyslexia from the German group are nearly as fast as the control group in finding all pairs (auditory) which might be due to cultural differences (e.g., more musical training).

The auditory and visual cues are designed on purpose to be more difficult to process for people with dyslexia than without. Therefore, children with dyslexia are expected to need more time (duration), which might be due to a less distinctive encoding of prosody (Goswami et al., 2016) and is in line with the indicator of slower reading. Considering that children with dyslexia need more time to process information, we observe this behavior as well for our indicators. For example, participants with dyslexia

from the Spanish group take more time on the *fourth click interval* and also on the *average click time* compared to the control group. Both results are significant and have medium effect sizes of 0.29, so we can estimate what the effects would be in the whole population (Field and Hole, 2003).

A person with dyslexia has difficulties with reading and writing independent of the mother tongue, which also appear when learning a second language (Helland and Kaasa, 2005; Nijakowska, 2010). The analysis of errors from children with dyslexia show similar error categories for Spanish, English (Rello et al., 2016a), and German (Rauschenberger et al., 2016), revealing similarities of perception between the languages.

Our results from the pilot study (Rauschenberger et al., 2018b) suggest that we can measure a significant difference on four indicators for the visual game with the same tendency between Spanish, German, and English. With all our data ( $n = 313$ ), we can confirm just one significant dependent variable with the same tendency for Spanish and German.

Still this means that people with dyslexia might perceive our visual game content similarly, independent of the mother tongue. Further research needs to be done to confirm the results, but this validation study provides strong evidence that it will be possible to screen dyslexia with our content, approach, and game design using the same language-independent content for different languages.

## 7.2 Screening Differences

Our approach aims to screen dyslexia with indicators that do not require linguistic knowledge. These indicators are probably not as strong or visible as the reading and spelling mistakes of children with dyslexia. Therefore, we consider our results (highest accuracy of 0.74 and highest F1-scores of 0.75) for German with Random Forest as a promising way to predict dyslexia using language-independent auditory and visual content for pre-readers.

Having an early indication of dyslexia before spelling or reading errors appear can have a positive impact on the child's development, as we can intervene earlier in her/his education. Therefore, we aim to optimize the recall and F1-score by finding as many participants with dyslexia as possible.

We have set ourselves this goal because early detection in a person with dyslexia has a greater positive effect on the person with dyslexia than a misjudgment in a person without dyslexia. However, to avoid over-fitting we did not modify the default value for the threshold (typically 0.5), something that we plan to study in the near future as we need to increase recall for the dyslexia class keeping a reasonable number of false positives.

If a person with dyslexia is not discovered (early), they are prone to face additional issues such as anxiety, sadness and decreased attention (Schulte-Körne, 2010). Also, a person with dyslexia needs around 2 years to compensate for their reading and spelling difficulties. Early treatment among children at risk of dyslexia as well as children without dyslexia can serve, both, as a preventive measure and as early stimulation of literacy skills.

Our results support the hypothesis that dyslexia cannot be reduced to one cause, but is rather a combination of characteristics (De Zubizaray and Schiller, 2018). The equal distribution of auditory and visual features in the informative

features ranking supports the hypothesis of dyslexia being related to auditory and visual perception in different people. We might be able to measure stronger effects when we design visual and auditory cues that have more attributes related to dyslexia, including some that favor the latter.

The ALL data set reached *only* an accuracy of 0.61, which might be due to the following reasons. First, the informative features for each data set are different from each other, which indicates different informativeness in German and Spanish. Combining the data sets into ALL probably adds noise for the prediction, which results in a lower accuracy. The noise might be that features are not as informative anymore because they cancel each other out as they are highly correlated.

In addition, reducing the features only to the features with the same tendency as used for the statistical analysis did not reveal any improvement, which supports the hypothesis that features in ALL cancel each other out.

The results of our current game measures with 313 participants confirm differences in the behavior of Spanish *vs.* German participants (i.e., 1) seven significant dependent variables in Spanish *vs.* none in German and 2) only two dependent variables with the same tendency over all languages).

These results might be explained by bilingualism. It is argued that a person who speaks more than one language has more knowledge of their first language than a monolingual person (Kecskes and Papp, 2000), and it is unclear whether this also has an influence on “how people perceive differences as well”. Additionally, dyslexia detection differences are reported for transparent (like Spanish) *vs.* deep (like English) orthographies [quoted after (Rello et al., 2019)]. In a transparent orthography mainly a single grapheme (letter) corresponds to a single phoneme (sound) and dyslexia is reported to be more distinct in deep orthographies.

If so, this might explain the difference we have in the significance for the statistical analysis as well as the tendency of values, and the need for separate models to predict dyslexia for our German *vs.* Spanish data set (Spanish has bilingual participants).

Overall, having fewer features improves the accuracy, but this is less so when we run experiments for ALL or ES. There, the influence of the different informative features for ES and DE seem to cancel each other out. The high correlation between features would explain why, for example, taking into account 27 features (GB) performs no better than using 20 features (GB) for the ALL data set. The fact that the accuracy does not increase when more features are used supports the argument that features are highly correlated.

As described before, small data can help to understand the data and results better. In our case, we see that ALL does not perform as well as ES or DE. This is probably due to the facts described above (e.g., bilingualism, features canceling each other, English-speaking participants). The prediction for dyslexia is therefore possible with the data taken from the same game, but needs different models for the prediction in different

languages as was proposed by (Bandhyopadhyay et al., 2018), something that made sense in retrospect.

## 8 CONCLUSIONS AND FUTURE WORK

We processed our game data with Extra Trees, Random Forest without and with class weights, and Gradient Boost to predict dyslexia using a data set of 313 participants. We reached the best accuracy of 74% for the German case using RF while the best accuracy for Spanish was 69% using ETC.

Our approach can optimize resources for detecting and treating dyslexia, however, it would need at the beginning more personnel to screen many more children at a young age to enlarge our training data. As children with dyslexia need around 2 years to compensate their difficulties, our approach could help to decrease school failure, late treatment and most importantly, to reduce suffering for children and parents.

The main advantage of our language-independent content approach is that has the potential to screen pre-readers in the near future. Indeed, we aim to collect more data with younger children to improve our results, use different input related to more characteristics of dyslexia, and other game design.

Future work includes improving our machine learning models and do further feature analysis. More explainable models should also be considered.

## DATA AVAILABILITY STATEMENT

The datasets presented in this study can be found in online repositories. The names of the repository/repositories and accession number(s) can be found in the article/**Supplementary Material**.

## ETHICS STATEMENT

The studies involving human participants were reviewed and approved by the institutions responsible. The data collection for this user study has been approved by the German Ministry of Education, Science and Culture in Schleswig-Holstein (Ministerium für Bildung, Wissenschaft und Kultur) and Lower Saxony State Education Authority (Niedersächsische Landesschulbehörde). In Spain, governmental approval was not needed in addition to the school approval. Written informed consent to participate in this study was provided by the participants' legal guardian/next of kin.

## AUTHOR CONTRIBUTIONS

All authors listed have made a substantial, direct, and intellectual contribution to the work and approved it for publication.

## ACKNOWLEDGMENTS

This is an extension of the previously published conference paper from W4A '20: Proceedings of the 17th International Web for All Conference in Tapei. This paper and content were partially funded by the *fem:talent Scholarship* from the *Applied University of Emden/Leer* as well as by the *Deutschen Lesepreis 2017* from the *Stiftung Lesen* and the *Commerzbank-Stiftung*. First, we would like to thank all teachers, students, and parents for their participation and time as well as all supporters to distribute my participation calls! Special thanks goes to one class and one teacher which cannot be named due to the anonymous regulations. We deeply thank for their support L. Albó, Barcelona; *ChangeDyslexia*, Barcelona; M. Jesús Blanque and R. Noé López, school *Hijas de San José*, Zaragoza; A. Carrasco, E. Méndez and S. Tena, innovation team of school *Leonardo da Vinci*, Madrid; in Spain, and L. Niemeier, *Fröbel Bildung und Erziehung gemeinnützige GmbH*, Berlin; E. Prinz-Burghardt, *Lerntherapeutische Praxis*, Duderstadt; L. Klaus, *Peter-Ustinov-Schule*, Eckernförde; H. Marquardt, *Gorch-Fock-*

*Schule*, Eckernförde; M. Batke and J. Thomaschewski, *Hochschule Emden/Leer*, Emden; N. Tegeler, *Montessori Bildungshaus Hannover gGmbH*, Hannover; Y. Schulz, *Grundschule Heidgraben*, Heidgraben; T. Westphal, *Leif-Eriksson-Gemeinschaftsschule*, Kiel; F. Goerke, *Grundschule Luetjensee*, Luetjensee; B. Wilke, *Schule am Draiberg*, Papenburg; P. Stümpel, *AncoraMentis*, Rheine; A. Wendt, *Grundschule Seth*, Seth; K. Usemann, *OGGS Meyerstraße*, Wuppertal; in Germany. We also thank all parents and children for playing *MusVis*. Finally, thanks to H. Witzel for his advice during the development of the visual part and to M. Blanca, and M. Herrera for the translation of the Spanish version.

## SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fcomp.2021.628634/full#supplementary-material>

## REFERENCES

- American Psychiatric Association (2013). *Diagnostic and Statistical Manual of Mental Disorders*. London, England: American Psychiatric Association. doi:10.1176/appi.books.9780890425596.744053
- Baeza-Yates, R. (2018). Big, Small or Right Data: Which Is the Proper Focus? Available at: <https://www.kdnuggets.com/2018/10/big-small-right-data.html>.
- Bandyopadhyay, A., Dey, D., and Pal, R. K. (2018). *Prediction of Dyslexia Using Machine Learning — A Research Travelogue*. Springer Singapore. doi:10.1007/978-981-10-6890-4
- Black, D. W., and Grant, J. E. American Psychiatric Association (2016). *DSM-5 Guidebook: The Essential Companion to the Diagnostic and Statistical Manual of Mental Disorders*. fifth edition. American Psychiatric Association.
- Borleffs, E., Maassen, B. A. M., Lyytinen, H., and Zwarts, F. (2019). Cracking the Code: The Impact of Orthographic Transparency and Morphological-Syllabic Complexity on Reading and Developmental Dyslexia. *Front. Psychol.* 9, 1–19. doi:10.3389/fpsyg.2018.02534
- Coleman, C., Gregg, N., McLain, L., and Bellair, L. W. (2008). A Comparison of Spelling Performance across Young Adults with and without Dyslexia. *Assess. Eff. Intervention* 34, 94–105. doi:10.1177/1534508408318808
- Cuetos, F., Ramos, J. L., and Ruano, E. (2002). *PROESC. Evaluación de los procesos de escritura (Writing processes assessment)*. Madrid: TEA.
- Cuetos, F., Rodríguez, B., Ruano, E., and Arribas, D. (2007). PROLEC-R: Batería de Evaluación de los Procesos Lectores, Revisada (Battery of reading processes assessment—Revised).
- De Zubicaray, G., and Schiller, N. O. (2018). *The Oxford Handbook of Neurolinguistics*. New York, NY: Oxford University Press.
- Ergonomics of human-system interaction (2010). Part 210: Human-Centred Design for Interactive Systems. In *Ergonomics of human-system interaction*. Brussels: International Organization for Standardization, 132.
- Esser, G., Wyszkon, A., and Schmidt, M. H. (2002). Was wird aus Achtjährigen mit einer Lese- und Rechtschreibstörung. *Z. für Klinische Psychol. Psychotherapie* 31, 235–242. doi:10.1026/0084-5345.31.4.235
- Faraway, J. J., and Augustin, N. H. (2018). When Small Data Beats Big Data. *Stat. Probab. Lett.* 136, 142–145. doi:10.1016/j.spl.2018.02.031
- Fastl, H., and Zwicker, E. (2007). *Psychoacoustics*. third edn. Berlin, Heidelberg: Springer Berlin Heidelberg.
- Fawcett, A., and Nicolson, R. (2004). *The Dyslexia Screening Test: Junior (DST-J)*. London, UK: Harcourt Assessment.
- Field, A., and Hole, G. (2003). *How to Design and Report Experiments*. London: SAGE Publications.
- Gaggi, O., Palazzi, C. E., Ciman, M., Galiazzo, G., Franceschini, S., Ruffino, M., et al. (2017). Serious Games for Early Identification of Developmental Dyslexia. *Comput. Entertain.* 15, 1–24. doi:10.1145/2629558
- Geurts, L., Vanden Abeele, V., Celis, V., Husson, J., Van den Audenaeren, L., Loyez, L., et al. (2015). “DIESEL-X: A Game-Based Tool for Early Risk Detection of Dyslexia in Preschoolers,” in *Describing and Studying Domain-specific Serious Games* (Switzerland: Springer), 93–114. doi:10.1007/978-3-319-20276-1\_7
- Goswami, U. (2011). A Temporal Sampling Framework for Developmental Dyslexia. *Trends Cogn. Sci.* 15, 3–10. doi:10.1016/j.tics.2010.10.001
- Goswami, U., Barnes, L., Mead, N., Power, A. J., and Leong, V. (2016). Prosodic Similarity Effects in Short-Term Memory in Developmental Dyslexia. *Dyslexia* 22, 287–304. doi:10.1002/dys.1535
- Grund, M., Naumann, C. L., and Haug, G. (2004). *Diagnostischer Rechtschreibtest Für 5. Klassen: DRT 5 (Diagnostic Spelling Test for Fifth Grade: DRT 5)*. aktual edn. Göttingen, Germany: Deutsche Schultests Göttingen: Beltz Test, 2.
- Hämäläinen, J. A., Salminen, H. K., and Leppänen, P. H. T. (2013). Basic Auditory Processing Deficits in Dyslexia. *J. Learn. Disabil.* 46, 413–427. doi:10.1177/0022219411436213
- Hamari, J., Koivisto, J., and Sarsa, H. (2014). Does Gamification Work? -- A Literature Review of Empirical Studies on Gamification. In 2014 47th Hawaii International Conference on System Sciences. IEEE, 3025–3034. doi:10.1109/HICSS.2014.377
- Helland, T., and Kaasa, R. (2005). Dyslexia in English as a Second Language. *Dyslexia* 11, 41–60. doi:10.1002/dys.286
- Huss, M., Verney, J. P., Fosker, T., Mead, N., and Goswami, U. (2011). Music, Rhythm, Rise Time Perception and Developmental Dyslexia: Perception of Musical Meter Predicts reading and Phonology. *Cortex* 47, 674–689. doi:10.1016/j.cortex.2010.07.010
- Jain, A., and Zongker, D. (1997). Feature Selection: Evaluation, Application, and Small Sample Performance. *IEEE Trans. Pattern Anal. Machine Intell.* 19, 153–158. doi:10.1109/34.574797
- Johnson, D. J. (1980). Persistent Auditory Disorders in Young Dyslexic Adults. *Bull. Orton Soc.* 30, 268–276. doi:10.1007/BF02653723
- Kecskes, I., and Papp, T. n. (2000). *Foreign Language and Mother Tongue*. 1 edn. New York: Psychology Press. doi:10.4324/9781410606464
- Männel, C., Schaadt, G., Illner, F. K., van der Meer, E., and Friederici, A. D. (2017). Phonological Abilities in Literacy-Impaired Children: Brain Potentials Reveal

- Deficient Phoneme Discrimination, but Intact Prosodic Processing. *Dev. Cogn. Neurosci.* 23, 14–25. doi:10.1016/j.dcn.2016.11.007
- Mora, A., Riera, D., Gonzalez, C., and Arnedo-Moreno, J. (2015). A Literature Review of Gamification Design Frameworks. In 7th International Conference on Games and Virtual Worlds for Serious Applications. doi:10.1109/VSGAMES.2015.7295760
- Nijkowska, J. (2010). *Dyslexia in the Foreign Language Classroom*. Bristol, United Kingdom: Multilingual Matters.
- Overy, K. (2000). Dyslexia, Temporal Processing and Music: The Potential of Music as an Early Learning Aid for Dyslexic Children. *Psychol. Music* 28, 218–229. doi:10.1177/0305735600282010
- Paulesu, E., Danelli, L., and Berlinger, M. (2014). Reading the Dyslexic Brain: Multiple Dysfunctional Routes Revealed by a New Meta-Analysis of PET and fMRI Activation Studies. *Front. Hum. Neurosci.* 8, 830. doi:10.3389/fnhum.2014.00830
- Petscher, Y., Fien, H., Stanley, C., Gearin, B., Fletcher, J. M., Johnson, E., et al. (2019). *Screening for Dyslexia*. Washington, DC: U.S. Department of Education, Office of Elementary and Secondary Education, Office of Special Education Programs, National Center on Improving Literacy. Retrieved from: [improvingliteracy.org](http://improvingliteracy.org).
- Poole, A., Zulkernine, F., and Aylward, C. (2017). Lexa: A Tool for Detecting Dyslexia through Auditory Processing. In 2017 IEEE Symposium Series on Computational Intelligence, SSCI 2017 - Proceedings 2018-January, 1–5. doi:10.1109/SSCI.2017.8285191
- Port, R. F. (2003). Meter and Speech. *J. Phonetics* 31, 599–611. doi:10.1016/j.wocn.2003.08.001
- Rauschenberger, M., and Baeza-Yates, R. (2021a). How to Handle Health-Related Small Imbalanced Data in Machine Learning? *i-com* 19, 215–226. doi:10.1515/icom-2020-0018
- Rauschenberger, M., and Baeza-Yates, R. (2020b). “Recommendations to Handle Health-Related Small Imbalanced Data in Machine Learning,” in *Mensch und Computer 2020 - Workshopband (Human and Computer 2020 - Workshop proceedings)*. Editor B. Hansen (Bonn: Gesellschaft für Informatik e.V.), 1–7. Christian AND Nurnberger, Andreas AND Preim. doi:10.18420/muc2020-ws111-333
- Rauschenberger, M., Baeza-Yates, R., and Rello, L. (2021a). Data: Documentation Semi-structure Literature Review 0.1. doi:10.13140/RG.2.2.19378.94401
- Rauschenberger, M., Baeza-Yates, R., and Rello, L. (2021b). MusVis User Dataset - as Csv (Resource, Content). doi:10.13140/RG.2.2.20633.95846
- Rauschenberger, M., Baeza-Yates, R., and Rello, L. (2021c). Protocol of the Generated Audio Files Selected for A Universal Screening Tool for Dyslexia by a Web-Game and Machine Learning. doi:10.13140/RG.2.2.27348.12162
- Rauschenberger, M., Baeza-Yates, R., and Rello, L. (2021d). Protocol of the Iterations to Select the Content for A Universal Screening Tool for Dyslexia by a Web-Game and Machine Learning. doi:10.13140/RG.2.2.17281.79209
- Rauschenberger, M., Baeza-Yates, R., and Rello, L. (2020). *Screening Risk of Dyslexia through a Web-Game Using Language-independent Content and Machine Learning*. Taipei: ACM Press, 1–12. doi:10.1145/3371300.3383342
- Rauschenberger, M., Baeza-Yates, R., and Rello, L. (2019b). “Technologies for Dyslexia,” in *Web Accessibility Book*. Editors Y. Yesilada and S. Harper. 2 edn (London: Springer-Verlag London), 1, 603–627. doi:10.1007/978-1-4471-7440-010.1007/978-1-4471-7440-0\_31
- Rauschenberger, M., Lins, C., Rousselle, N., Fudickar, S., and Hain, A. (2019a). A Tablet Puzzle to Target Dyslexia Screening in Pre-readers. In Proceedings of the 5th EAI International Conference on Smart Objects and Technologies for Social Good - GOODTECHS. Valencia, 155–159.
- Rauschenberger, M., Rello, L., and Baeza-Yates, R. (2018a). Barcelona: ACM Press, 306–312. doi:10.1145/3236112.3236156
- Rauschenberger, M., Rello, L., Baeza-Yates, R., and Bigham, J. P. (2018b). Towards Language Independent Detection of Dyslexia with a Web-Based Game. In W4A '18: The Internet of Accessible Things. Lyon, France: ACM, 4–6. doi:10.1145/3192714.3192816
- Rauschenberger, M., Rello, L., Baeza-Yates, R., Gomez, E., and Bigham, J. P. (2017a). Supplement: DysMusicMusicalElements: Towards the Prediction of Dyslexia by a Web-Based Game with Musical Elements. doi:10.5281/zenodo.809783
- Rauschenberger, M., Rello, L., Baeza-Yates, R., Gomez, E., and Bigham, J. P. (2017b). Towards the Prediction of Dyslexia by a Web-Based Game with Musical Elements. In The Web for All conference Addressing information barriers - W4A'17. Western Australia: PerthACM Press, 4–7. doi:10.1145/3058555.3058565
- Rauschenberger, M., Rello, L., Fuchsel, S., and Thomaschewski, J. (2016). “A Language Resource of German Errors Written by Children with Dyslexia,” in Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016) (Paris, France: European Language Resources Association).
- Rauschenberger, M., Willems, A., Teinieden, M., and Thomaschewski, J. (2019c). Towards the Use of Gamification Frameworks in Learning Environments. *J. Interactive Learn. Res.* 30.
- Rello, L., Baeza-Yates, R., Ali, A., Bigham, J. P., and Serra, M. (2019). Predicting Risk of Dyslexia with an Online Gamified Test. arXiv preprint arXiv:1906.03168 V.1, 1–13.
- Rello, L., Baeza-Yates, R., and Llisterri, J. (2016a). A Resource of Errors Written in Spanish by People with Dyslexia and its Linguistic, Phonetic and Visual Analysis. *Lang. Resour. Eval.* 51, 379–408. doi:10.1007/s10579-015-9329-0
- Rello, L., Baeza-Yates, R., Ali, A., Bigham, J. P., and Serra, M. (2020). Predicting risk of dyslexia with an online gamified test. *PLoS ONE* 15 (12), e0241687. doi:10.1371/journal.pone.0241687
- Rello, L., Romero, E., Rauschenberger, M., Ali, A., Williams, K., Bigham, J. P., et al. (2018). Screening Dyslexia for English Using HCI Measures and Machine Learning. In Proceedings of the 2018 International Conference on Digital Health - DH '18. New York, New York, USA: ACM Press, 80–84. doi:10.1145/3194658.3194675
- Ritzhaupt, A. D., Poling, N. D., Frey, C. A., and Johnson, M. C. (2014). A Synthesis on Digital Games in Education: What the Research Literature Says from 2000 to 2010. *Jl. Interactive Learn. Res.* 25, 263–282.
- Rolka, E. J., and Silverman, M. J. (2015). A Systematic Review of Music and Dyslexia. *The Arts in Psychotherapy* 46, 24–32. doi:10.1016/j.aip.2015.09.002
- Rouse, R. (2004). *Game Design: Theory and Practice in Theory and Practice*. Second Edition Second Edition (Plano, TX: Wordware Publishing, Inc.).
- Schulte-Körne, G., Deimel, W., Müller, K., Gutenbrunner, C., and Remschmidt, H. (1996). Familial Aggregation of Spelling Disability. *J. Child. Psychol. Psychiat* 37, 817–822. doi:10.1111/j.1469-7610.1996.tb01477.x
- Schulte-Körne, G. (2010). The Prevention, Diagnosis, and Treatment of Dyslexia. *Deutsches Ärzteblatt Int.* 107, 718–727. doi:10.3238/arztebl.2010.0718
- Scikit-learn (2019). 3.1. Cross-Validation: Evaluating Estimator Performance. Available at: <https://scikit-learn.org/stable/modules/crossvalidation.html>.
- Scikit-learn Developers (2019). Scikit-learn Documentation. Available at: <https://scikit-learn.org/stable/documentation.html>.
- Seaborn, K., and Fels, D. I. (2015). Gamification in Theory and Action: A Survey. *Int. J. Human-Computer Stud.* 74, 14–31. doi:10.1016/j.ijhcs.2014.09.006
- Steinbrink, C., and Lachmann, T. (2014). *Lese-Rechtschreibstörung (Dyslexia)*. Springer Berlin Heidelberg. doi:10.1007/978-3-642-41842-6Lese-Rechtschreibstörung
- Tallal, P. (2004). Improving Language and Literacy Is a Matter of Time. *Nat. Rev. Neurosci.* 5, 721–728. doi:10.1038/nrn1499
- Thomas, A., Bader, F., Thomaschewski, J., and Rauschenberger, M. (2021). Integrating Gamification: The Human-Centered Gamification Process. In Proceedings of the 17th International Conference on Web Information Systems and Technologies. Online, 2021, 430–435. doi:10.1038/nrn1499
- Varma, S., and Simon, R. (2006). Bias in Error Estimation when Using Cross-Validation for Model Selection. *BMC Bioinformatics* 7, 91. doi:10.1186/1471-2105-7-91
- Vidyasagar, T. R., and Pammer, K. (2010). Dyslexia: a Deficit in Visuo-Spatial Attention, Not in Phonological Processing. *Trends Cogn. Sci.* 14, 57–63. doi:10.1016/j.tics.2009.12.003
- Weigand, A. C., Lange, D., and Rauschenberger, M. (2021). How Can Small Data Sets Be Clustered in Mensch und Computer 2021 {Workshopband} {Workshop on User-Centered Artificial Intelligence (UCAI '21), 1. Bonn, Germany: Association for Computing Machinery. doi:10.18420/muc2021-mci-ws02-284
- Wikipedia (2019). Memory (Spiel) (Memory Game).

- World Health Organization (2019). *International Classification of Diseases 11th Revision*. World Health Organization.
- World Health Organization (2010). *International Statistical Classification of Diseases and Related Health Problems 10th Revision*. Geneva, Switzerland: World Health Organization.
- Yuskaitis, C. J., Parviz, M., Loui, P., Wan, C. Y., and Pearl, P. L. (2015). Neural Mechanisms Underlying Musical Pitch Perception and Clinical Applications Including Developmental Dyslexia. *Curr. Neurol. Neurosci. Rep.* 15, 51. doi:10.1007/s11910-015-0574-9

**Conflict of Interest:** The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

**Publisher's Note:** All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

*Copyright © 2022 Rauschenberger, Baeza-Yates and Rello. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.*