# Journal of Educational Psychology

## Human and Automated Assessment of Oral Reading Fluency

Daniel Bolaños, Ron A. Cole, Wayne H. Ward, Gerald A. Tindal, Jan Hasbrouck, and Paula J. Schwanenflugel

# Human and Automated Assessment of Oral Reading Fluency

Daniel Bolaños, Ron A. Cole, and Wayne H. Ward
Boulder Language Technologies, Boulder, Colorado

Gerald A. Tindal
University of Oregon

Jan Hasbrouck
Gibson Hasbrouck & Associates, Wellesley, Massachusetts

Paula J. Schwanenflugel
The University of Georgia

This article describes a comprehensive approach to fully automated assessment of children's oral reading fluency (ORF), one of the most informative and frequently administered measures of children's reading ability. Speech recognition and machine learning techniques are described that model the 3 components of oral reading fluency: word accuracy, reading rate, and expressiveness. These techniques are integrated into a computer program that produces estimates of these components during a child's 1-min reading of a grade-level text. The ability of the program to produce accurate assessments was evaluated on a corpus of 783 one-min recordings of 313 students reading grade-leveled passages without assistance. Established standardized metrics of accuracy and rate (words correct per minute [WCPM]) and expressiveness (National Assessment of Educational Progress Expressiveness scale) were used to compare ORF estimates produced by expert human scorers and automatically generated ratings. Experimental results showed that the proposed techniques produced WCPM scores that were within 3–4 words of human scorers across students in different grade levels and schools. The results also showed that computer-generated ratings of expressive reading agreed with human raters better than the human raters agreed with each other. The results of the study indicate that computer-generated ORF assessments produce an accurate multidimensional estimate of children's oral reading ability that approaches agreement among human scorers. The implications of these results for future research and near term benefits to teachers and students are discussed.

Keywords: oral reading fluency, automated reading assessment, expressive reading, automatic speech recognition

Reading assessments provide school districts and teachers with critical and timely information for identifying students who need immediate help; for making decisions about reading instruction; for monitoring individual student's progress in response to instruc-tional interventions; for comparing different approaches to reading instruction; and for reporting annual outcomes in classrooms, schools, school districts, and states. One of the most common tests administered to primary school students is *oral reading fluency* (ORF). Over 25 years of scientifically based reading research has established that fluency is a critical component of reading and that effective reading programs should include instruction in fluency (Fuchs, Fuchs, Hosp, & Jenkins, 2001; Kuhn & Stahl, 2000; National Reading Panel, 2000). Although ORF does not measure comprehension directly, there is substantial evidence that estimates of ORF predict future reading performance and correlate strongly with comprehension (Fuchs et al., 2001; Shinn, 1998). According to Wayman, Wallace, Wiley, Tichá, and Espin (2007), ORF is a valid indicator of comprehension in early grades, though less so beyond Grade 4. Because ORF can be measured rather quickly (typically in 5–10 min) with good validity and reliability, it is widely used to screen individuals for reading problems and to measure reading progress over time.

In this article, we present a comprehensive approach to assessing ORF accurately and automatically through the use of speech recognition and machine learning techniques. The approach is comprehensive because all three measures of ORF—accuracy, rate (combined into a words correct per minute [WCPM] score), and expressiveness—can be measured automatically and in real time, whereas expressiveness is rarely scored in real-world educational

contexts. The ultimate goal of research leading to fully automatic and comprehensive assessment of ORF is to provide an accurate, accessible, and low-cost alternative to human-administered assessments. Successful outcomes of research in this area would substantially reduce the millions of hours that teachers spend each year assessing their students' reading abilities, which is mandated by federal law in the United States. In addition, computer-based assessments of ORF could generate detailed records of individual student's performance, including the digital recordings of each reading session that could be reviewed by teachers, parents, and students, and analyzed automatically for detailed information about the student's reading problems. Automatic administration of ORF will also enable collection of massive amounts of speech data that can be used to analyze and understand children's development of reading skills; these data can also be used to improve the performance of the speech recognition technologies.

We used a speech recognition system (Bolaños, 2012) specifically designed to process children's read speech to produce a word-level hypothesis of what the student read from a grade-level text during 1 min. From this hypothesis and the text passage, a WCPM score was computed reflecting the student's reading accuracy and rate. In order to assess prosodic reading, we developed a series of lexical and prosodic features that were extracted from the student's speech. These included analysis of the text syntax and its correlation with filled pauses and silence regions, syllable and word duration, pitch, and word co-occurrences, among other features described below. Machine learning classifiers were trained on these features, resulting in statistical models that were able to discriminate between different degrees of prosodic reading using the National Assessment of Educational Progress ORF Scale (NAEP; Daane, Campbell, Grigg, Goodman, & Oranje, 2005). A hierarchical classification scheme was used in order to assign 1-min reading sessions to levels in the NAEP scale.

The accuracy of these assessment methods was evaluated on approximately 13 hr of speech collected from the 313 first-through fourth-grade students who read grade-level text passages. WCPM scores as well as NAEP assessments generated by the system, FLuent Oral Reading Assessment (FLORA), were compared with those produced by at least two independent human judges.

The remainder of the article is organized as follows: The next section provides the scientific rationale for assessing ORF. We then describe the corpus of children's read speech that was collected for this study. We then describe the system and features used to assess WCPM (accuracy and rate) and expressive reading using lexical and prosodic features extracted from the speech. The last section presents the discussion and conclusions.

## Scientific Rationale for FLORA

### ORF

*ORF* is typically defined as a student's ability to read words in grade-level texts accurately and effortlessly, at a natural speech rate and with appropriate prosodic expression. A synthesis of scientifically based reading research by the National Reading Panel (2000) concluded that

> Reading fluency is one of several critical factors necessary for reading comprehension, but it is often neglected in the classroom. If children

read out loud with speed, accuracy and proper expression, they are more likely to comprehend and remember the material than if they read with difficulty and in an inefficient way.

**Accuracy and automaticity.**    Accurate reading speed is both a strong discriminator of reading ability (e.g., Jenkins, Fuchs, van den Broek, Espin, & Deno, 2003; Perfetti, 1985) and a strong predictor of later reading proficiency (Lesgold & Resnick, 1982; Scarborough, 1998; see review by Compton & Carlisle, 1994.) As Jenkins et al. (2003) put it: "Together with listening comprehension, word-reading skill accounts for nearly all of the reliable variance in reading ability, and individual differences in word recognition explain significant variance in reading ability, even after controlling for reading comprehension" (Curtis, 1980; Hoover & Gough, 1990).

ORF depends on the ability to recognize words in a text *quickly and automatically*. As defined by Fuchs et al. (2001), automaticity is "the oral translation of text with speed and accuracy." Automaticity theory (LaBerge & Samuels, 1974; Samuels, 1985; Wolf, 1999) and related verbal efficiency accounts of reading (Perfetti, 1985) hold that students who have learned to decode printed words automatically are able to devote more attention (cognitive resources) to comprehending what they are reading. Readers who have not achieved automaticity during word recognition must devote significant attention to recognizing words (at the expense of devoting this attention to making sense of the text), resulting in slower reading times and weaker comprehension. Support for automaticity and verbal efficiency theories of reading is provided by the strong association between the speed of reading words, either in word lists or in context, and measures of reading comprehension.

**Expressiveness.**    Although readers who have achieved fluency can read texts rapidly and accurately, they may not read expressively (i.e., they may not pause between sentences, at major phrase boundaries within sentences, or produce appropriate prosody when reading out loud). Expressive reading is the third critical component of *reading fluency*, typically defined as reading a text with the appropriate expression, intonation, and phrasing in order to preserve meaning (Miller & Schwanenflugel, 2008).

**Connection between ORF and comprehension.**    For over 25 years, researchers have documented the association between reading fluency and comprehension. Reviews of the research on ORF have demonstrated consistently moderate to strong correlations between ORF and comprehension (Marston, 1989; Shinn, 1998). Research results have demonstrated high concurrent validity between ORF and measures of word recognition and reading comprehension (Hosp & Fuchs, 2005; Jenkins et al., 2003), and between ORF and nationally normed standardized tests of reading comprehension (Roehrig, Petscher, Nettles, Hudson, & Torgesen, 2008; Schilling, Carlisle, Scott, & Zeng, 2007; Schwanenflugel et al., 2006). Measures of ORF in early grades have also been found to predict comprehension in later grades (Kim, Petscher, Schatschneider, & Foorman, 2010). Thus, the relation between ORF and reading comprehension has been well established by previous research, particularly for students in elementary school (Kim et al., 2010; Roberts, Good, & Corcoran, 2005; Roehrig et al., 2008).

## Previous Work Using Automatic Speech Recognition to Assess and Improve ORF

**Automatic assessment of reading accuracy and rate.** Over two decades of research has investigated the use of automatic speech recognition (ASR) to assess and improve reading. Seminal research conducted by Jack Mostow and his colleagues in Project Listen at Carnegie Mellon University has demonstrated the effectiveness of ASR for improving reading fluency and comprehension for both native and nonnative speakers of English (Mostow et al., 2003; Reeder, Shapiro, & Wakefield, 2007). Mostow et al. (2003) used an ASR system to measure a student's *interword latency*, defined as the elapsed time between certain words read aloud by the student that were scored as correctly read by the ASR system. Their model of interword latency produced a correlation of over .7 with independent WCPM measures of ORF using grade-level passages.

In the context of Project Tball (Technology Based Assessment of Language and Literacy) at the University of California, Los Angeles and University of Southern California, Black, Tepperman, Lee, and Narayanan (2008) investigated oral reading of 55 isolated words produced by kindergarten, first-, and second-grade children with the aim of detecting reading miscues automatically, such as sounding-out, hesitations, whispering, elongated onsets, and question intonations. Black et al. developed an ASR system that used specialized grammars to model word-level disfluencies using the subword-modeling approach developed by Hagen and Pellom (2005). Scores produced by the recognition system correlated highly (.91) with fluency judgments provided by human listeners.

A series of studies by Bryan Pellom and Andreas Hagen and their collaborators (Hagen, Pellom, & Cole, 2007) investigated ways to optimize an ASR system for children's read speech. The research resulted in a reduction in the word error rate from 17.4% to 7.6%. Hagen et al. (2007) developed a version of the ASR system that used subword-modeling rather than whole-word scoring to detect reading errors. In the study, several subword lexical units and approaches were evaluated for detection of reading disfluencies, and modest gains were reported. Bolaños (2008) reported that additional detection gains were achieved by using syllable graphs to represent hypotheses from the ASR system.

**Automatic assessment of expressive oral reading.** Although the National Reading Panel (2000) and research community define ORF in terms of word recognition accuracy, reading rate, and how expressively the student reads (see Kuhn, Schwanenflugel, & Meisinger, 2010, for a discussion of this topic), expressiveness is rarely measured in assessments of ORF. Only recently has the expressiveness aspect of the reading fluency construct found its way into automated assessments of fluency. Duong, Mostow, and Sitaram (2011) investigated two alternative methods of measuring prosody during children's oral reading. The first method, which was text dependent, consisted of generating a prosodic template model for each sentence in the text. The template was based on word-level features like pitch, intensity, latency, and duration extracted from fluent adult narrations. The second method investigated adult narrations to train a general duration model that could be used to generate expected prosodic contours of sentences for any text, so an adult reader was no longer required to generate sentence templates for each new text. Both methods were evaluated for their ability to predict student's scores on fluency and comprehension tests, and each produced promising results, with the second, automated method for generating prosodic sentence templates outperforming the system that compared children's read speech with adult narrations of each individual sentence in the text. However, neither of these methods could satisfactorily classify sentences using the NAEP expressiveness rubric relative to human judgments, which was probably due to the low human interrater reliability reported in this study.

## Development of the FLORA System

### Development of a Corpus for Assessing ORF

**Data collection setting.** Data were collected from 313 first-through fourth-grade students in four elementary schools (nine classrooms) in the Boulder Valley School District in Colorado. Data were collected from students in their classrooms at their schools. School 1 had 53.8% students receiving free or reduced lunches, and the lowest literacy achievement scores of the three schools on the Colorado state literacy test given to third-grade students; 53% third-grade students in School 1 scored proficient or above on the state reading assessment. School 2 had 51.7% students with free or reduced lunch (similar to School 1), but 79% of third-grade students tested as proficient or above on the state literacy test. School 2 was a bilingual school with nearly 100% English learners (ELs) who spoke Spanish as their first language. School 3 had 18.4% of students with free or reduced lunch, 85% of students were proficient or above in the state literacy test. School 3 also had relatively few ELs.

**Text passages.** Twenty text passages were available for reading at each grade level. The standardized text passages were downloaded from a website (Good & Kaminski, 2002) and are freely available for noncommercial use. The text passages were designed specifically to assess ORF and are about the same level of difficulty within each grade level. ORF norms have been collected for these text passages for tens of thousands of students at each grade level in fall, winter, and spring semesters, so that students can be assigned to percentiles based on national WCMP scores (Hasbrouck & Tindal, 2006).

**Data collection protocol.** The data were collected using the FLORA system (Bolaños, Cole, Ward, Borts, & Svirsky, 2011), which was configured to enroll each student, randomly select one passage from the set of 20 standardized passages for the student's grade level, and present the passage to the student for reading out loud. Because testing was conducted in May, near the end of the school year, classroom teachers had recently assessed their student's oral reading performance (using text passages different from those used in our study). About 20% of the time, teachers requested that specific students be presented with text passages either one or two levels below or one or two levels above the student's grade level. Thus, about 80% of students in each grade read passages at their grade level, whereas 20% of students read passages above or below their grade level, based on their teachers' recommendations. Depending on the number of students who needed to be tested on a given day, each student was presented with two or three text passages to read aloud.

During the testing procedure, the student was seated before a laptop and wore a set of headphones with an attached noise-cancelling microphone. The experimenter observed or helped the

student enroll in the session, which involved entering the student's gender, age, and grade level. FLORA then presented a text passage, started the 1-min recording at the instant the passage was displayed, recorded the student's speech, and relayed the speech to a server.

**Corpus summary.** The corpus comprised 783 recordings from 313 first- through fourth-grade students for a total of approximately 13 hr of speech data. Each recording was scored manually by two human judges. Words were scored as reading errors if the word was skipped over, or the judge decided that the word was misread. Insertions of words (intrusions) were not scored as reading errors, as insertions were not counted as errors in the national norms collected by Hasbrouck and Tindal (Hasbrouck & Tindal, 2006).

## Automatic Generation of WCPM Scores

The number of words that a student read correctly during 1 min was computed automatically by ReadToMe, the reading tracker built on top of our ASR system (Bolaños, 2012). The computation of the WCPM score was done as follows. (a) ReadToMe used the Bavieca speech recognition toolkit (Bolaños, 2012) to produce a word-level hypothesis representing what the student read. (b) ReadToMe aligned the hypothesis to the reference text (the words in the text passage the student read) and tagged each of the words in the reference text as correctly or incorrectly read or skipped over. (c) Finally, ReadToMe counted the number of words scored as correctly read during the 1-min reading; this number is the resulting WCPM score for the text passage.

## Automatic Assessment of Expressive Reading

In order to assess expressive reading automatically, we proposed a set of lexical and prosodic features that can be used to train a machine learning system to classify how expressively students read text passages aloud using the 4-point NAEP scale. The proposed features were designed to measure the speech behaviors associated with each of the four levels of fluency described in the NAEP rubric and were informed by research on acoustic-phonetic, lexical, and prosodic correlates of fluent and expressive reading described in the research literature (Kuhn et al., 2010). Features were extracted from multiple sources, including the recognition hypothesis, a pitch-extractor, and a syllabification tool. Features included the WCPM score itself, the speaking rate, sentence reading rate, number of word repetitions, location of the pitch accent, word and syllable durations, and filled and unfilled pauses and their correlation to punctuation marks in the text passage. A detailed description, motivation, and analysis of all the features proposed and used for the study can be found in Bolaños et al. (2013).

**Classification method.** In order to classify the 783 one-min recordings using the features proposed, we used a powerful classification technique called support vector machines (Vapnick, 1995). We experimented with difference classification strategies and found a strategy based on a decision directed acyclic graph (DAG) to be most successful (Platt, Cristianini, & Shawe-Taylor, 2000). The DAG approach makes sense conceptually because it maps directly to the NAEP scale; that is, it distinguishes disfluent reading (Levels 1 and 2 in the NAEP scale) from fluent reading

(Levels 3 and 4 in the NAEP scale) and then makes finer distinctions (1 vs. 2 and 3 vs. 4). To implement the DAG strategy, we trained three classifiers. The first classifier was trained on samples from all classes and separated samples from Classes 1 and 2 and 3 and 4. This classifier was placed at the root of the tree, whereas two other classifiers, trained on samples from Classes 1 and 2 and 3 and 4, respectively, were placed on the leaves to make the finer-grained decisions. A detailed description of the classification scheme can be found in Bolaños et al. (2013).

## Speech Recognition System

A total of 106 hr or read speech from three different children's speech corpora were used to train the recognition system. The recognizer was not trained on the corpus of read speech, described above, that was used to evaluate FLORA. We note that the system is *text independent*; that is, for new text passages, the system automatically generates the expected pronunciation(s) of each word in a text passage from a pronunciation dictionary.

The speech recognition system combines two main sources of information to produce a score for each word. These sources are (a) the score produced by matching the system's acoustic models for the expected sequence of phonemes in a word (based on a pronunciation dictionary) to the student's pronunciation of the word and (b) the probability of the word occurring in the text (the statistical language model, based on the co-occurrence of words in the text passage). These two sources of information are combined to produce the most likely hypothesis string given the speech input. Additionally, phone-level alignments from each of the 1-min recordings were generated automatically for feature extraction purposes. Two complementary speaker adaptation techniques were used in order to tailor the speaker-independent acoustic models to the speech characteristics and vocal tract length of each speaker.

## Comparison Between Automated and Human Assessments of ORF

### Human Scoring of Recorded Sessions

In order to evaluate the ability of FLORA to produce reliable WCPM scores, each of the 783 one-min recordings in the evaluation corpus was scored independently by two former elementary school teachers. Each teacher had more than a decade of experience administering reading assessments to elementary school children. The scorers were able to listen to, review, and modify their judgments within each recording until they were satisfied with their WCPM score. Thus, they were allowed to listen to the recording more than once.

Additionally, each of the 783 recordings was scored from 1 to 4 using the NAEP ORF scale by at least two independent scorers, who were former elementary school teachers with experience assessing reading proficiency. A set of 70 stories of the total 783 stories were scored by the five available teachers, whereas the other recordings were scored by just two of them, which were randomly assigned to each scorer. A training session was scheduled before the scoring process to review the NAEP scoring instructions and unify criteria. The judges first listened to passages rated by two experienced researchers whose area of expertise is expressive reading (Paula Schwanenflugel and Melanie Kuhn).

The teachers who scored the stories then rated these passages and compared their ratings with the experts. The teachers then rated several additional passages and reviewed their ratings based on the definitions of each of the NAEP levels. The training was concluded when the teachers' level of agreement approximated the agreement exhibited by the two experts.

For the actual scoring of the evaluation corpus, the judges listened to each 60-s story in 20-s intervals and provided a 1–4 rating for each interval. The NAEP ORF scale (Daane et al., 2005) comprises four levels from less to more fluent. Level 1 is characterized by word-by-word reading, Level 2 by reading using two-word phrases with some three- or four-word groupings, and Level 3 is characterized by a majority of three- or four-word phrase groups while preserving the syntax of the author. Readers at Level 4 produce larger, meaningful phrase groups with expressive interpretation. Finally, scorers attached a global NAEP score to the recording based on the NAEP scores assigned to each 20-s segment. The global score was based on a review of the scores and their best judgment rather than using a deterministic method like the mean or mode.

## Assessment of Reading Accuracy and Automaticity

Table 1 shows the means and standard deviations (between parentheses) for accuracy, words per minute (WPM), and WCPM scores for the human scorers and FLORA. Statistics are shown per reading level for students in the four schools. As noted above, although the evaluation data were collected from students from Grades 1 to 4, about 20% of the time, teachers requested that specific students be presented with text passages either one or two levels below or above the student's grade level, resulting in reading levels for text passages from Grades 1 to 6. In Table 1, accuracy is expressed in percentages and WPM, which measures fluency from the perspective of speed-ignoring accuracy, and the score is based on the average across the two human scorers for each recording. It can be seen that accuracy (percentage of words read correctly) is higher for higher grade levels, from 70.3% for first grade to 92.6% and 90.5% for fifth- and sixth-grade levels, respectively. WPM are displayed in Column 5 for each grade level; as expected, they are highly correlated, with WCPM measured by human scorers (Column 5); however, WCPM computed by FLORA (Column 7) are much closer to human WCPM scores (Column 6) than WPM.

A major result can be observed by comparing the WCPM scores from the human scorers and FLORA, which present very similar distributions (means and standard deviations). In addition, we observed very similar distributions of WCPM scores from humans and FLORA within each of the nine classrooms in which we conducted the study, even for classrooms in schools in which the majority of students spoken Spanish as their first language and were officially designated as English learners.

Column 8 shows the expected number of WCPM for each grade level according to Hasbrouck and Tindal (2006) reading norms. It can be seen in the table that students were assigned by teachers to reading levels at which they read around the 50th percentile. We believe that there is no credible evidence to link higher WCPM scores to improved comprehension, but there is substantial support for the need for readers to have an accuracy and rate (WCPM score) in the range of the 50th percentile to support both comprehension and motivation.

Another pattern of results is revealed by examining the numbers in Column 9, which shows the mean difference in WCPM scores for the two human scores for the recordings in each classroom, and the numbers in Column 10, which shows the mean difference between the averaged human scores and FLORA for each classroom. Note that differences in WCPM scores are expressed in absolute value. Viewing the numbers in Column 9 reveals the remarkable agreement between the two human scorers (1.2-WCPM difference across all schools) and the low variance. Across all recordings, the mean difference between FLORA and the averaged human scores was 3.6 words, whereas the mean difference between human scores was 1.2 words.

Figure 1a displays a scatter plot of the WCPM scores from the two human scorers for all recordings, whereas Figure 1b displays a scatterplot of the WCPM scores from FLORA with respect to the average human scores for all recordings. If agreement were perfect, all points would lie on the diagonal. These figures show the strong agreement between WCPM scores for human scorers on each recording, and the very good agreement between FLORA and the human scores, with relatively few outliers.

We were interested in determining whether FLORA might be a useful tool for providing WCPM scores that could be used as one valuable indicator, along with other measures, to identify students who are at risk for failing to learn to read. One way to do this is to compare human and FLORA WCPM scores with the national reading norms developed by Hasbrouck and Tindal (2006), which measured WCPM scores, for first- through sixth-grade students during each trimester of a school year. The interrater agreement in the task of mapping recorded stories to percentiles was 0.97 for the

Table 1

*Summary of Accuracy, WPM, and WCPM According to Human Scorers (H) and FLORA (F). Expected WCPM (E) Are Also Shown*

| Level | Stu. | Rec. | Acc. (%) | H-WPM | H-WCPM | F-WCPM | E-WCPM | H-diff | FH-diff |
|-------|------|------|----------|-------|--------|--------|--------|--------|---------|
| 1 | 68 | 171 | 70.3 (19.7) | 54.6 (25.5) | 41.9 (26.4) | 42.5 (25.9) | 53 | 1.2 (1.8) | 2.7 (2.7) |
| 2 | 97 | 242 | 84.6 (10.1) | 99.3 (31.9) | 85.7 (33.1) | 86.1 (31.8) | 89 | 1.2 (2.0) | 3.8 (4.4) |
| 3 | 52 | 128 | 87.3 (7.6) | 113.4 (28.1) | 100.1 (29.6) | 101.6 (28.0) | 107 | 1.2 (1.4) | 3.6 (2.8) |
| 4 | 59 | 147 | 87.4 (8.1) | 124.4 (26.6) | 109.9 (27.3) | 112.7 (27.3) | 123 | 1.3 (1.8) | 4.1 (3.1) |
| 5 | 30 | 76 | 92.6 (3.6) | 156.9 (26.4) | 145.6 (26.6) | 145.6 (24.5) | 139 | 1.1 (1.2) | 4.6 (4.5) |
| 6 | 7 | 19 | 90.5 (14.1) | 145.9 (46.1) | 137.3 (49.6) | 137.4 (49.1) | 150 | 1.5 (2.0) | 2.8 (2.6) |
| All | 313 | 783 | 83.3 (14.0) | 103.3 (42.2) | 90.1 (43.1) | 91.1 (42.6) | | 1.2 (1.8) | 3.6 (3.6) |

*Note.* WPM = words per minute; WCPM = words correct per minute; FLORA = FLuent Oral Reading Assessment; Stu. = number of students; Rec. = number of recordings; Acc. = accuracy; H-diff = difference between the human scorers; FH-diff = difference between the FLORA and human scorers.
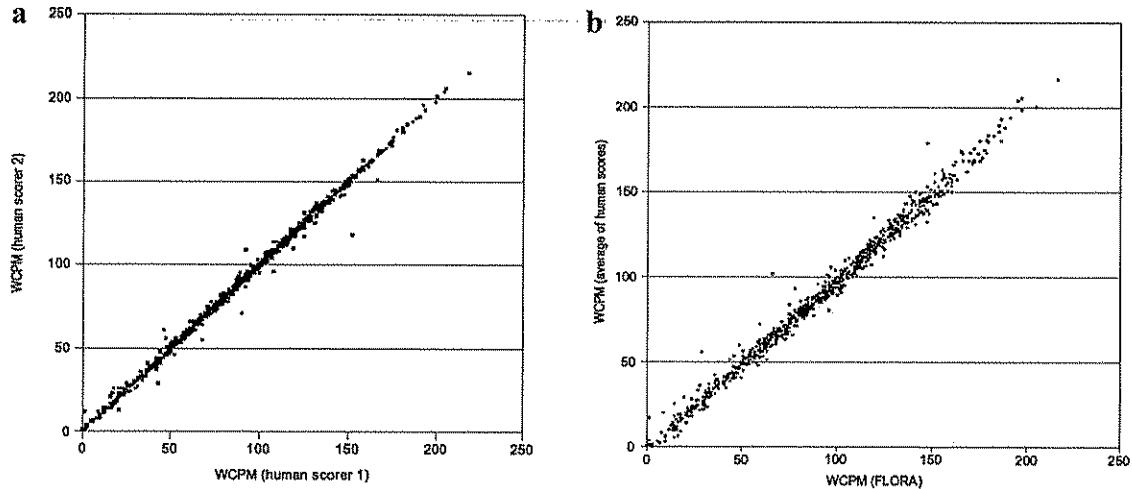
*Figure 1.* Correlation between WCPM scores produced by two independent human scorers (a) and between FLORA and the average of the two independent human scorers (b) for each of the 1-min recordings assessed. WCPM = words correct per minute; FLORA = FLuent Oral Reading Assessment.

human scorers and 0.89 between FLORA and each of the human scorers. The interrater agreement in the task of mapping recorded stories above and below the 50th percentile (which is used normally as a reference to identify at-risk students) was 0.98 for the human scorers and 0.92 between FLORA and each of the human scorers. Agreement was computed using the weighted kappa coefficient ($\kappa$) (Cohen, 1968), which is suitable for ordinal categories. In sum, the interhuman agreement and the FLORA to human agreement is very close, which means that FLORA performs well at identifying students who might require additional reading assessments and instruction.

## Assessment of Expressive Reading

In this section, we present results on assessing expressive oral reading using FLORA. First, we briefly analyze the classification accuracy for the lexical and prosodic features proposed in relation to human assessments. We then analyze agreement and correlation between human scores and FLORA's automatic scoring system using the NAEP scale.

**Classification accuracy.** In order to derive the most effective combination of features to assess expressive reading, we measured the classification accuracy (percentage of recordings that FLORA assigned the same label than the human labelers) of FLORA on the corpus described above. Each recording was labeled by FLORA according to the NAEP scale, and labels were compared with those from all the available human labelers. We note that there exists an upper bound to the classification accuracy that can be attained by the classifier. The reason is that whenever the human raters score the same recording differently, there is an unrecoverable classification error.

Results showed that both lexical and prosodic features contributed similarly to the classification accuracy for the NAEP-2 (disfluent vs. fluent) task (89.27% and 89.02%, respectively). This can be initially considered an unexpected result because lexical aspects like the number of words read correctly are expected to dominate the discrimination between fluent and nonfluent readers. However,

it is important to note that some of the prosodic features defined in this study are highly correlated with the lexical features. For example, it is obvious that the number of words correctly read in a 1-min reading session should correlate highly with the average duration of a silence region or the number of filled pauses made.

For both the NAEP-2 and NAEP-4 tasks, lexical and prosodic features provided complementary information that led to improved classification accuracy when combined. For the NAEP-4 tasks, lexical features seem to have a dominant role (73.24% and 69.73%, respectively). We attribute this to the WCPM score, which is taken as a lexical feature; this score by itself provides a 71.78% accuracy for the NAEP-4 task. As expected, the automatically computed WCPM, which comprises two of the three reading fluency cornerstones (accuracy and rate), plays a fundamental role. In particular, the combination resulted in accuracies of 90.72% and 75.87% for the NAEP-2 and NAEP-4 tasks, respectively. Finally, note that the distribution of recordings across the NAEP levels according to humans and machine was very similar.

## Interrater Agreement and Correlation

In this section, we present interrater agreement and correlation results for the best system from the previous section (multilabel training using all the features). Table 2 shows the interrater agreement for the tasks of classifying recordings into the broad NAEP categories (fluent vs. nonfluent; NAEP-2), or the four levels of expressiveness using the NAEP rubric (NAEP-4). For the NAEP-2 task, the interrater agreement was measured using Cohen's kappa coefficient ($\kappa$) (Cohen, 1960); p(a) is the probability of observed agreement, whereas p(e) is the probability of chance agreement.

For the NAEP-4 task, we measured the interrater agreement using the weighted kappa coefficient ($\kappa$) (Cohen, 1968), which is more suitable for ordinal categories given that it weights disagreements differently depending on the distance between the categories (we used linear weightings). As a complementary metric for this task, we computed the Spearman's rank correlation coefficient (Spearman, 1904). In a number of classification problems, like

Table 2

*Interrater Agreement and Correlation Coefficients on the NAEP Scale*

| Scorer | # recordings | NAEP-2 | | | NAEP-4 | |
|---|---|---|---|---|---|---|
| | | p(a) | p(e) | κ | κ | ρ |
| Human 1 | 571 | 0.87 | 0.50 | 0.73 | 0.66 | 0.80 |
| Human 2 | 391 | 0.90 | 0.50 | 0.80 | 0.69 | 0.81 |
| Human 3 | 698 | 0.87 | 0.50 | 0.74 | 0.68 | 0.81 |
| Human 4 | 799 | 0.86 | 0.50 | 0.71 | 0.69 | 0.81 |
| Human 5 | 367 | 0.86 | 0.50 | 0.71 | 0.68 | 0.80 |
| FLORA | 1,776 | 0.94 | 0.50 | 0.84 | 0.77 | 0.86 |

*Note.* NAEP = National Assessment of Educational Progress; FLORA = FLuent Oral Reading Assessment.

emotion classification, the data are annotated by a group of human raters who may exhibit consistent disagreements on similar classes or similar attributes. In such classification tasks, it is inappropriate to assume that there is only one correct label because different individuals may consistently provide different annotations (Steidl, Levit, Batliner, Nöth, & Niemann, 2005). Although the NAEP scale is based on clear descriptions of reading behaviors at each of four levels, children's reading behaviors can vary across these descriptions while reading, and individuals scoring the stories may differ consistently in how they interpret and weight children's oral reading behaviors. For this reason, we believe that examining correlations between human raters and between human raters and the machine classifiers is a meaningful and useful metric for this task.

Each row in Table 2 shows the agreement and correlation coefficients of each rater with respect to the other raters (excluding FLORA in the case of the human raters; note that not all the scorers scored the same number of recordings). In order to interpret the computed kappa values, we have used as a reference the interpretation of the kappa coefficient provided in Landis and Koch (1977), which attributes *good* agreement to kappa values within the interval (0.61–0.80) and *very good* agreement to higher kappa values (0.81–1.00). According to this interpretation, Table 2 reveals that (a) there is *good* interhuman agreement for both the NAEP-2 and NAEP-4 tasks, (b) there is *good* FLORA-to-human agreement for the NAEP-4 task, and (c) there is *very good* FLORA-to-human agreement for the NAEP-2 task. It can be observed that the kappa agreement between FLORA and the humans is higher than the agreement between each human scorer and the rest of the human scorers. This is true for both the NAEP-2 and NAEP-4 tasks. This difference in agreement is statistically significant, which indicates the ability of the proposed features and

classification scheme to provide a useful method to automatically assess expressive oral reading using the NAEP scale.

In terms of the Spearman's rank correlation coefficient (ρ), we obtained relatively strong interhuman correlation (.80–.81) and an even stronger machine-to-human correlation (.86) in the NAEP-4 task. This indicates that NAEP scores from every pair of scorers are closely related, which is consistent with the weighted kappa values obtained.

In Table 3, we display cross-tabs of agreement and disagreement between humans and between FLORA and humans (in percentages). In both cases, most of the data lie in the main diagonal, and we believe that there are no obvious biases between humans and FLORA.

## Connection Between Reading Accuracy, Reading Rate, and Expressive Reading

We conducted a set of analyses to gain insights into the relationship between the two main measures of ORF, WCPM, and expressiveness. These analyses are displayed in Figure 2a and 2b. In each panel of the figure, we sorted students according to their WCPM percentile using the Hasbrouck and Tindal (2006) norms. Thus, the leftmost bar of each panel represents students with WCPM scores below the 10th percentile, whereas the rightmost bar shows students in the 90th percentile. Figure 2a displays percentile assignments based on average human scorers rating, and Figure 2b displays percentile assignments based on FLORA WCPM estimates. The tones of gray within each bar indicate the percentage of students at each NAEP score; in Figure 2a, these numbers are based on the NAEP scores assigned by the human scorers, and in Figure 2b these numbers were assigned by FLORA.

Table 3

*Cross-Tabs of Agreement/Disagreement Between FLORA and Human-Generated NAEP Scores (in %)*

| | | FLORA | | | | | | Human | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | | | 1 | 2 | 3 | 4 |
| Human | 1 | 16.6 | 2.9 | 0.1 | 0 | Human | 1 | 14.9 | 4.2 | 0.1 | 0 |
| | 2 | 3.5 | 21.3 | 3.9 | 0.2 | | 2 | 4.4 | 19.6 | 5.7 | 0.2 |
| | 3 | 0 | 3.9 | 32.4 | 5.6 | | 3 | 0.1 | 7.2 | 27.7 | 5.2 |
| | 4 | 0 | 0 | 3 | 6.6 | | 4 | 0 | 0 | 4.7 | 6.1 |

*Note.* FLORA = FLuent Oral Reading Assessment; NAEP = National Assessment of Educational Progress.
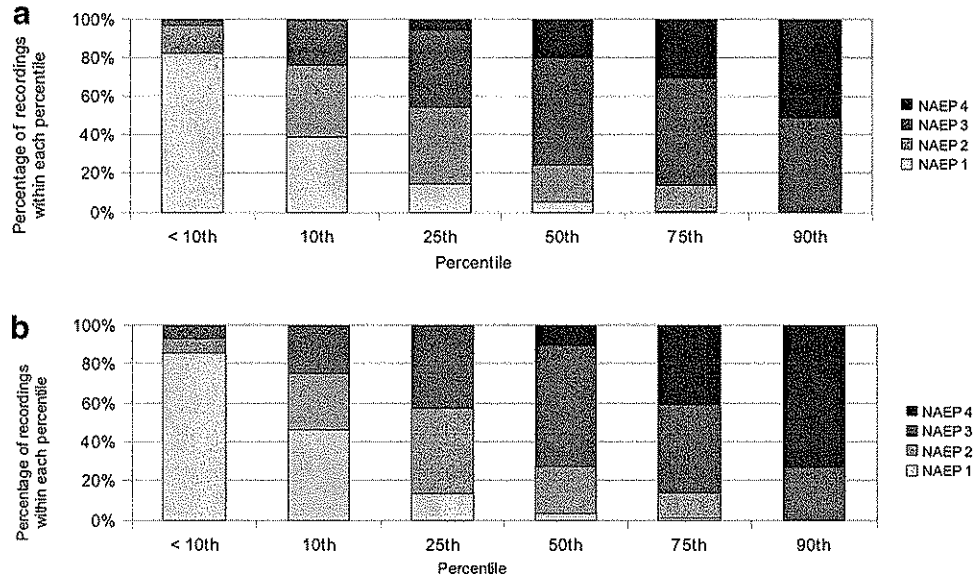
*Figure 2.* a: Distribution of recordings across the NAEP scale for each WCPM percentile according to human scorers. b: Distribution of recordings across the NAEP scale for each WCPM percentile according to FLORA. NAEP = National Assessment of Educational Progress; WCPM = words correct per minute.

It is clear from this figure that recordings in the highest percentiles (highest reading accuracy and rate) correspond to more expressive readers (higher levels in the NAEP scale). For example, all of the recordings for students in the 90th percentile based on WCPM were assigned to Levels 3 and 4 in the NAEP scale. Moreover, about 97.0% of the recordings below the 10th percentile were assigned to Levels 1 and 2 in the NAEP scale. Figures 2a and 2b reveal several interesting patterns: A significant percentage of recordings placed below the 50th percentile (which might be used to identify students in need for fluency support) were placed in the higher levels of the NAEP scale according to our expert human annotators (3.08%, 24.02%, and 45.19% for recordings below the 10th percentile, in the 10th percentile, and in the 25th percentile, respectively). This means that there are a number of speakers who, despite reading below the expected rate according to the percentiles published by Hasbrouck and Tindal (2006), read with appropriate/good expression and would be considered fluent readers according to the NAEP scale. Another interesting observation is that a significant percentage of recordings placed above the 50th percentile were assigned to the lower levels in the NAEP scale by our expert human annotators. Those recordings likely correspond to speakers who are reading for speed rather than for comprehension in order to get as many words read as possible within the 1-min session. In particular, 24.88% of the recordings in the 50th percentile were assigned to Levels 1 and 2 in the NAEP scale (nonfluent), whereas 13.92% of the recordings in the 75th percentile were assigned to those levels. We note that the instructions provided to students before recording stories emphasized the importance of reading the text naturally, rather than as fast as they could; these percentage might have been higher if we had not emphasized reading naturally in the instructions. These observations suggest that measuring both expressiveness and WCPM is likely to be both informative and beneficial to understanding

individual student's oral reading abilities. Finally, we note that Figure 2b, which is analogous to Figure 2a but was built using FLORA scores, presents very similar information.

## Discussion and Conclusions

We investigated the automatic assessment of ORF in children's speech according to two standard rubrics: WCPM (to measure accuracy and rate) and the NAEP Expressiveness scale. Compared with human scoring of WCPM and expressiveness on 783 one-min recordings of children reading grade-level text passages, results show that automatically generated WCPM scores differ by an average of 3.5 words with respect to the human-average score for each recorded story, whereas humans differ by an average of 1.5 words for each story.

For expressiveness, FLORA had an accuracy of 90.93% classifying recordings according to the binary NAEP scale ("fluent" vs. "nonfluent") and 76.05% on the more difficult 4-point NAEP scale. According to the classification of kappa strength proposed by Landis and Koch (1977), the kappa agreement for both NAEP-2 and NAEP-4 tasks between each human scorer and the rest of the human scorers was *good*, whereas the kappa agreement between the machine and the human scorers was *good* and *very good*, respectively. In addition, the kappa agreement between FLORA and each human scorer was always significantly higher than the kappa agreement between the human scorers. In terms of the Spearman's rank correlation coefficient ($\rho$), correlation between the machine and each human scorer was always significantly higher than the correlation between human scorers.

The results of the research reveal that speech recognition and machine learning systems can produce accurate assessments of WCPM and expressiveness that approach (WCPM) or exceed human performance. Without question, the results of the WCPM

scores reported above can be improved substantially in the near future using known ASR solutions, such as collecting more training data to model children's speech patterns. For example, Vergyri, Lamel, and Gauvain (2010) reported that accent-dependent acoustic modeling (which implies training/adapting on data from the target accent) produces a significant increase in recognition performance compared with accent-independent modeling. In a recent study that we conducted on 191 native Spanish children learning to read English text in Spanish schools (Bolaños, Elhazaz, Ward, & Cole, 2012), we determined experimentally that statistical models trained on speech from the target population were significantly more accurate than models trained on native English children. Results from that study showed a mean difference in WCPM scores of 5.49 and 4.96, respectively, between FLORA and each of the human scorers, whereas the mean difference between the human scorers was about 5.92 words.

Perhaps the major limitation of this study is the relatively small number of students (313) used in our research. To fully demonstrate the feasibility and validity of a fully automatic assessment of ORF, speech data during oral reading of leveled texts must be collected for a large and diverse population of students at different grade levels, representing students with different dialects and accents. The system must also be tested with data collected from many different classrooms or computer labs to model the acoustic environments and the realities of real-world use.

## Toward Valid Automatic Assessment of ORF

We believe there are great potential benefits of incorporating measures of expressiveness into assessments of ORF. One of the major criticisms of using WCPM to measure individual student's improvements in reading over time (i.e., in response to instruction) is that students strive to read texts as quickly as possible in order to increase their WCPM scores, which teachers often set as learning targets within a reading instruction program. When a student's ability is measured in terms of how quickly he or she can read the words in a text, teachers and students learn to focus on reading fast, rather than reading the text at a normal reading rate with intonation and phrasing that communicates the meaning of text, and thus reflects its comprehension by the student. Fast readers have shorter segment durations, muted stress marking, and reduced phrase-final bracketing than slow readers, so the normal comprehension benefits children might experience by reading with good prosody may not be derived by students who are trying to read fast (Benjamin & Schwanenflugel, 2010; Kuhn et al., 2010). In sum, the emphasis on speed that can result from using WCPM as the only measure of ORF may undermine the goal of helping students develop strategies for reading with deep understanding.

Incorporating measures of expressiveness into assessments of ORF could mitigate this problem. One can easily imagine a weighted measure of ORF that combines WCPM and expressiveness estimates, such that students receive the highest score when the words in a text are read at a natural speaking rate with prosody appropriate to the discourse structure of the text. In fact, some rating systems of reading expressiveness such as the Multidimensional Fluency Guide (Rasinski, Rikli, & Johnston, 2009) already do this.

One of the major benefits of the automated scoring of reading prosody by FLORA that neither the NAEP nor the other various teacher rating systems for evaluating reading fluency have is that these reading fluency scales have not (as yet) been grounded in research on reading prosody. We do not know whether the ratings obtained using these scales would be spectrographically valid, that is, that children rated as expressive on these scales would be the same ones who would appear expressive when their readings are viewed on a spectrogram. Because the features used in FLORA to classify expressive reading were derived directly from spectrographic measures derived from children's speech (Kuhn et al., 2010), FLORA can make this claim. Conversely, because the teacher NAEP ratings match the spectrographic distinctions made by FLORA, FLORA has also served to validate teacher impressions of reading prosody as determined by the NAEP. In sum, fully automatic assessment or ORF that combines its three components appears to be feasible with today's technologies. Additional research is needed to determine how to use these measures to provide the most useful feedback to teachers and students to assess students' reading abilities and inform instruction.

## References

Benjamin, R., & Schwanenflugel, P. J. (2010). Text complexity and oral reading prosody in young readers. *Reading Research Quarterly, 45,* 388–404.

Black, M., Tepperman, J., Lee, S., & Narayanan, S. (2008). *Estimation of children's reading ability by fusion of automatic pronunciation verification and fluency detection.* Proceedings of Interspeech, Brisbane, Australia.

Bolaños, D. (2012, December). *The Bavieca open-source speech recognition toolkit.* In Proceedings of IEEE Workshop on Spoken Language Technology (SLT), Miami, FL.

Bolaños, D., Cole, R. A., Ward, W., Borts, E., & Svirsky, E. (2011). FLORA: Fluent oral reading assessment of children's speech. *ACM Transactions on Speech and Language Processing, 7,* 1–19. doi: 10.1145/1998384.1998390

Bolaños, D., Cole, R. A., Ward, W., Tindal, G., Schwanenflugel, P., & Kuhn, M. (2013). Automatic assessment of expressive oral reading. *Speech Communication, 55,* 221–236. doi:10.1016/j.specom.2012.08 .002

Bolaños, D., Elhazaz, P., Ward, W., & Cole, R. (2012). Automatic assessment of oral reading fluency for native Spanish ELL children. In *Proceedings of WOCCI 2012: Workshop on Child, Computer and Interaction: Satellite Event of INTERSPEECH.*

Bolaños, D., Ward, W. H., Wise, B., & Vuuren, S. V. (2008, September). *Pronunciation error detection techniques for children's speech.* Paper presented at INTERSPEECH 2008, Brisbane, Australia.

Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement, 20,* 37–46. doi:10.1177/ 001316446002000104

Cohen, J. (1968). Weighted kappa: Nominal scale agreement with provision for scaled disagreement or partial credit. *Psychological Bulletin, 70,* 213–220. doi:10.1037/h0026256

Compton, D. L., & Carlisle, J. F. (1994). Speed of word recognition as a distinguishing characteristic of reading disabilities. *Educational Psychology Review, 6,* 115–140. doi:10.1007/BF02208970

Curtis, M. E. (1980). Development of components of reading skill. *Journal of Educational Psychology, 72,* 656–669. doi:10.1037/0022-0663.72.5 .656

Daane, M. C., Campbell, J. R., Grigg, W. S., Goodman, M. J., & Oranje, A. (2005). *Fourth-grade students reading aloud: NAEP 2002 special study of oral reading* (NCES 2006-469). Washington, DC: U.S. Department of Education. Institute of Education Sciences, National Center for Education Statistics.

Duong, M., Mostow, J., & Sitaram, S. (2011). Two methods for assessing oral reading prosody. *ACM Transactions on Speech and Language Processing, 7*.

Fuchs, L. S., Fuchs, D., Hosp, M. K., & Jenkins, J. R. (2001). Oral reading fluency as an indicator of reading competence: A theoretical, empirical, and historical analysis. *Scientific Studies of Reading, 5*, 239–256. doi: 10.1207/S1532799XSSR0503_3

Good, R. H., & Kaminski, R. A. (2002). *Dynamic indicators of basic early literacy skills*. Eugene, OR: Institute for the Development of Educational Achievement.

Hagen, A., & Pellom, B. (2005, April). *A multi-layered lexical-tree based recognition of subword speech units*. Paper presented at the Second Language and Technology Conference: Human Language Technologies as a Challenge for Computer Science and Linguistics, Poznan, Poland.

Hagen, A., Pellom, B., & Cole, R. (2007). Highly accurate children's speech recognition for interactive reading tutors using subword units. *Speech Communication, 49*, 861–873. doi:10.1016/j.specom.2007.05.004

Hasbrouck, J., & Tindal, G. (2006). Oral reading fluency norms: A valuable assessment tool for reading teachers. *The Reading Teacher, 59*, 636–644. doi:10.1598/RT.59.7.3

Hoover, W. A., & Gough, P. B. (1990). The simple view of reading. *Reading and Writing, 2*, 127–160. doi:10.1007/BF00401799

Hosp, M. K., & Fuchs, L. S. (2005). Using CBM as an indicator of decoding, word reading, and comprehension: Do the relations change with grade? *School Psychology Review, 34*, 9–26.

Jenkins, J. R., Fuchs, L. S., van den Broek, P., Espin, C. L., & Deno, S. L. (2003). Sources of individual differences in reading comprehension and reading fluency. *Journal of Educational Psychology, 95*, 719–729. doi: 10.1037/0022-0663.95.4.719

Kim, Y.-S., Petscher, Y., Schatschneider, C., & Foorman, B. (2010). Does growth rate in oral reading fluency matter in predicting reading comprehension achievement? *Journal of Educational Psychology, 102*, 652–667. doi:10.1037/a0019643

Kuhn, M. R., Schwanenflugel, P. J., & Meisinger, E. B. (2010). Aligning theory and assessment of reading fluency: Automaticity, prosody, and definitions of fluency. *Reading Research Quarterly, 45*, 230–251. doi: 10.1598/RRQ.45.2.4

Kuhn, M. R., & Stahl, S. (2000). *Fluency: A review of developmental and remedial practices*. Ann Arbor, MI: Center for the Improvement of Early Reading Achievement.

LaBerge, D., & Samuels, S. (1974). Toward a theory of automatic information processing in reading. *Cognitive Psychology, 6*, 293–323. doi: 10.1016/0010-0285(74)90015-2

Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical data. *Biometrics, 33*, 159–174. doi:10.2307/2529310

Lesgold, A. M., & Resnick, L. B. (1982). How reading disabilities develop: Perspectives from longitudinal study. In J. P. Das, R. Mulcahy, & A. Wall (Eds.), *Theory and research in learning disability*. New York, NY: Plenum Press.

Marston, D. (1989). Curriculum-based measurement: What is it and why do it? In M. Shinn (Ed.), *Curriculum-based measurement: Assessing special children* (pp. 18–78). New York, NY: Guilford Press.

Miller, J., & Schwanenflugel, P. J. (2008). A longitudinal study of the development of reading prosody as a dimension of oral reading fluency in early elementary school children. *Reading Research Quarterly, 43*, 336–354. doi:10.1598/RRQ.43.4.2

Mostow, J., Aist, G., Burkhead, P., Corbett, A., Cuneo, A., Eitelman, S., & Tobin, B. (2003). Evaluation of an automated reading tutor that listens: Comparison to human tutoring and classroom instruction. *Journal of Educational Computing Research, 29*, 61–117. doi:10.2190/06AX-QW99-EQ5G-RDCF

National Reading Panel, National Institute of Child Health and Human Development. (2000). *Teaching children to read: An evidence-based assessment of the scientific research literature on reading and its implications for reading instruction*. Washington, DC: National Institute of Child Health and Human Development, National Institutes of Health.

Perfetti, C. (1985). *Reading ability*. Oxford, England: Oxford University Press.

Platt, J. C., Cristianini, N., & Shawe-Taylor, J. (2000). Large margin DAGs for multiclass classification. *Advances in Neural Information Processing Systems, 12*, 547–553.

Rasinski, T., Rikli, A., & Johnston, S. (2009). Reading fluency: More than automaticity? More than a concern for the primary grades? *Literacy Research and Instruction, 48*, 350–361. doi:10.1080/19388070802468715

Reeder, K., Shapiro, J., & Wakefield, J. (2007). The effectiveness of speech recognition technology in promoting reading proficiency and attitudes for Canadian immigrant children In *Proceedings of the 9th European Conference on Reading*. Berlin, Germany: IDEC.

Roberts, G., Good, R., & Corcoran, S. (2005). Story retell: A fluency-based indicator of reading comprehension. *School Psychology Quarterly, 20*, 304–317. doi:10.1521/scpq.2005.20.3.304

Roehrig, A. D., Petscher, Y., Nettles, S. M., Hudson, R. F., & Torgesen, J. K. (2008). Accuracy of the DIBELS oral reading fluency measure for predicting third grade reading comprehension outcomes. *Journal of School Psychology, 46*, 343–366. doi:10.1016/j.jsp.2007.06.006

Samuels, J. (1985). *Automaticity and repeated reading*. Lexington, MA: Lexington Books.

Scarborough, H. S. (1998). Early identification of children at risk for reading difficulties: Phonological awareness and some other promising predictors. In B. K. Shapiro, P. J. Accardo, & A. J. Capute (Eds.), *Specific reading disability: A view of the spectrum* (pp. 75–199). Timonium, MD: York Press.

Schilling, S. G., Carlisle, J. F., Scott, S. E., & Zeng, J. (2007). Are fluency measures accurate predictors of reading achievement? *Elementary School Journal, 107*, 429–448. doi:10.1086/518622

Schwanenflugel, P. J., Meisinger, E. B., Wisenbaker, J. M., Kuhn, M. R., Strauss, G. P., & Morris, R. D. (2006). Becoming a fluent and automatic reader in the early elementary school years. *Reading Research Quarterly, 41*, 496–522.

Shinn, M. (1998). *Advanced applications of curriculum based measurement*. New York, NY: Guilford Press.

Spearman, C. (1904). The proof and measurement of association between two things. *American Journal of Psychology, 15*, 72–101. doi:10.2307/1412159

Steidl, S., Levit, M., Batliner, A., Nöth, E., & Niemann, H. (2005). "Of all things the measure is man": Automatic classification of emotions and interlabeler consistency. In *Proceedings ICASSP* (Vol. 1, pp. 317–320). doi:10.1109/ICASSP.2005.1415114

Vapnik, V. (1995). *The nature of statistical learning theory*. New York, NY: John Wiley & Sons.

Vergyri, D., Lamel, L., & Gauvain, J. L. (2010, September). *Automatic speech recognition of multiple accented English data*. Paper presented at INTERSPEECH 2010, Makuhari, Japan.

Wayman, M. M., Wallace, T., Wiley, H. I., Tichá, R., & Espin, C. A. (2007). Literature synthesis on curriculum-based measurement in reading. *The Journal of Special Education, 41*, 85–120. doi:10.1177/00224669070410020401

Wolf, M. (1999). What time may tell: Towards a new conceptualization of developmental dyslexia. *Annals of Dyslexia, 49*, 1–28. doi:10.1007/s11881-999-0017-x