

Florida Assessments for Instruction in Reading, Aligned to the Language Arts Florida Standards

FAIR – FS

Grades 3 through 12

Technical Manual

Barbara R. Foorman, Ph.D.

Yaacov Petscher, Ph.D.

Chris Schatschneider, Ph.D.

Florida Center for Reading Research
Florida State University



Acknowledgements

The items, dynamic flow, computer-adaptive task flow and placement algorithms, creation of the development application, and psychometric work for this component skills battery (called the Florida Center for Reading Research Reading Assessment; FRA) were funded by grants from the Institute of Education Sciences (IES) to Florida State University:

Institute of Education Sciences, USDOE (\$4,447,990), entitled “Assessing Reading for Understanding: A Theory-Based, Developmental Approach,” subcontract to the Educational Testing Service for five years (R305F100005), 7/1/10-6/30/15 (Foorman, PI on subcontract; Petscher and Schatschneider, Co-Is).

Institute of Education Sciences, USDOE (R305A100301; \$1,499,741), entitled “Measuring Reading Progress in Struggling Adolescents,” awarded for four years, 3/1/10-2/28/14. (Foorman, PI; Petscher and Schatschneider, Co-Is).

The Florida State University licensed the FRA to the Florida Department of Education at no cost in perpetuity in 2012. We would like to acknowledge the following individuals for their leadership in executing the work funded by the above two IES grants: Dr. Adrea Truckenmiller, Karl Hook, and Nathan Day. We also would like to thank the numerous school districts administrators and teachers who participated in the research funded by these two grants.

For the purpose of this technical manual, we refer to the portions of the FRA licensed to the Florida Department of Education as the Florida Assessments for Instruction in Reading, Aligned to the Language Arts Florida Standards (FAIR-FS). Components of the FAIR-FS in grades 3-12 which are owned by the Florida Department of Education (e.g., ORT) are described in a separate technical manual.

Table of Contents

Acknowledgements	2
Introduction	5
<i>Mastering the Alphabetic Principle.....</i>	<i>5</i>
<i>Comprehending Written Language (better known as Reading Comprehension)</i>	<i>6</i>
<i>Summary of FAIR-FS Constructs and Tasks</i>	<i>9</i>
<i>Description of the Tasks in the FAIR-FS.....</i>	<i>9</i>
Description of Method	11
<i>Item Response Theory.....</i>	<i>11</i>
<i>Guidelines for Retaining Items.....</i>	<i>13</i>
<i>Linking Design & Item Response Analytic Framework.....</i>	<i>14</i>
<i>Norming Studies.....</i>	<i>15</i>
<i>Score Definitions</i>	<i>16</i>
Reliability	18
<i>Marginal Reliability.....</i>	<i>18</i>
<i>Standard Error of Measurement.....</i>	<i>20</i>
<i>Test-Retest Reliability.....</i>	<i>23</i>
Validity.....	25
<i>Predictive Validity.....</i>	<i>25</i>
<i>Differential Accuracy of Prediction.....</i>	<i>28</i>
<i>Construct Validity.....</i>	<i>36</i>
References	39

Introduction

The first question to ask when designing an assessment of reading and language skills is what predicts success in comprehending written language, that is, success in word reading and in reading comprehension? We are fortunate to have several consensus documents that review decades of literature about what predicts reading success (NRC, 1998; NICHD, 2000; NIFL, 2008; Rand, 2002; Rayner, Foorman, Perfetti, Pesetsky, & Seidenberg, 2001).

Mastering the Alphabetic Principle

What matters the most to success in reading words in an alphabetic orthography such as English is mastering the alphabetic principle, the insight that speech can be segmented into discrete units (i.e., phonemes) that map onto orthographic (i.e., graphemic) units (Ehri, Nunes, Willows, et al., 2001; Rayner et al., 2001). Oral language is acquired largely in a natural manner within a hearing/speaking community; however, written language is not acquired naturally because the graphemes and their relation to phonological units in speech are invented and must be taught by literate members of the community. The various writing systems (i.e., orthographies) of the world vary in the transparency of the sound-symbol relation. Among alphabetic orthographies, the Finnish orthography is highly transparent: phonemes in speech relate to graphemes in print (i.e., spelling) in a highly consistent one-to-one manner and graphemes in print relate to phonemes in speech (i.e., decoding) in a highly consistent one-to-one manner. Thus, learning to spell and read Finnish is relatively easy. English, however, is a more opaque orthography. Phonemes often relate to graphemes in an inconsistent manner and graphemes relate to phonemes in yet a different inconsistent manner. For example, if we hear the “long sound of *a*” we can think of words with many different vowel spellings, such as *crate*, *brain*, *hay*, *they*, *maybe*, *eight*, *great*, *vein*. If we see the orthographic unit *-ough*, we may struggle with the various pronunciations of *cough*, *tough*, *though*, *bough*. The good news is that 69% of monosyllabic English words—those Anglo-Saxon words most used in beginning reading instruction—are consistent in their letter to pronunciation mapping (Ziegler, Stone, & Jacobs, 1997). Most of the rest can be learned with grapheme-phoneme correspondence rules (i.e., phonics), with only a small percentage of words being so irregular in their letter-sound relations that they should be taught as sight words (Ehri, Nunes, Stahl, & Willows, 2001; Foorman & Connor, 2011).

Redacted K-2 detail no longer part of FAIR 3-12

Redacted K-2 detail no longer part of FAIR 3-12

In grades 3-12, alphabetic skills are measured with a word recognition task. In this computer-adaptive task, three words are presented on the computer monitor and students must select the word that best matches the word pronounced by the computer. About 10% of target words are nonsense words so that phonological decoding skills are tapped. When the target is a real word, distractors tap orthographic knowledge. For example, a distractor for “prerogative” might be *perogative*. By tapping orthographic knowledge in this task, the quality of a student’s lexical representation for a printed word is assessed. The more complete and accurate the lexical representation of a word is, the more efficient the student’s word recognition and reading comprehension (Perfetti & Stafura, 2014).

Comprehending Written Language (better known as Reading Comprehension)

Knowledge of word meanings. Mastering the alphabetic principle is a necessary but not sufficient condition for understanding written text. We may be able to pronounce printed words, but if we don’t know their meaning our comprehension of the text is likely to be impeded. Hence, our knowledge of word meanings is crucial to comprehending what we read. Grasping the meaning of a word is more than knowing its definition in a particular passage. Knowing the meaning of a word means knowing its full lexical entry in a dictionary: pronunciation, spelling, multiple meanings in a variety of contexts, synonyms, antonyms, idiomatic use, related words, etymology, and morphological structure. For example, a dictionary entry for the word *exacerbate* says that it is a verb meaning: 1) to increase the severity, bitterness, or violence of (disease, ill feeling, etc.); aggravate or 2) to embitter the feelings of (a person); irritate; exasperate (e.g., foolish words that only exacerbated the quarrel). It comes from the Latin word *exacerbātus* (the past participle of *exacerbāre*: to *exasperate*, *provoke*), equivalent to *ex + acerbatus* (*acerbate*). Synonyms are: *intensify*, *inflame*, *worsen*, *embitter*. Antonyms are: *relieve*, *sooth*, *alleviate*, *assuage*. Idiomatic equivalents are: add fuel to the flame, fan the flames, feed the fire, or pour oil on the fire. The more a reader knows about the meaning of a word like *exacerbate*, the greater the lexical quality the reader has and the more likely the reader will be able to recognize the word quickly in text, with full comprehension of its meaning (Perfetti & Stafura, 2014).

In the grades 3-12 FAIR-FS, knowledge of word meanings is measured by a Vocabulary Knowledge Task that taps morphological awareness. In the Vocabulary Knowledge Task, the student reads a sentence that has a missing word. The student selects among three words the one that best completes the sentence. The distractors and target vary in their morphological structure (i.e., prefixes or suffixes consisting of inflectional morphemes or derivational morphemes). It is relatively easy to read derived words that are pronounced similarly to their base (e.g., *reason, reasonable*). Words that contain a phonological shift (e.g., *vine, vineyard*) or an orthographic shift (e.g., *pity, piteous*) are harder to read, and words that contain both a phonological and an orthographic shift (e.g., *theory, theoretical*) are the hardest of all (Carlisle & Stone, 2005). The Vocabulary Knowledge Task in the FAIR-FS explained 2%-9% unique variance beyond prior reading comprehension, text reading efficiency, and spelling in predicting spring reading comprehension (Foorman, Petscher, & Bishop, 2012) and, by doing so, addresses aspects of language critical to understanding written language, language often called *academic language* because it is found in books and at school but not in informal conversations at home or outside school. Part of academic language is *inferential language* or *decontextualized language*, which allows speakers or writers to go beyond the present context and to predict, hypothesize, compare and contrast, and reason about events (e.g., an upcoming *referendum*) or abstract concepts (e.g., *photosynthesis, gravity*). Examples of words that signal such inferential or decontextualized language are *describe, analyze, hypothesize*.

Syntactic awareness. In addition to understanding word meanings, another important aspect of academic language is syntactic awareness. Syntax or grammar refers to the rules that govern how words are ordered to make meaningful sentences. Children typically acquire these rules in their native language prior to formal schooling. However, learning to apply these rules to reading and writing is a goal of formal schooling and takes years of instruction and practice. In the grades 3-12 FAIR-FS, there is a diagnostic task called Syntactic Knowledge Task (SKT). In this task the student listens to a sentence that is missing a word and selects the best word from a dropdown menu to complete the sentence. The words are verbs, pronouns, or connectives. Connectives are words that represent causal (e.g., *because*), temporal (e.g., *when*), logical (e.g., *if-then*), additive (e.g., *in addition*), or adversative (e.g., *although*) relations and are important linguistic devices for linking ideas and information within and across sentences. They link back to information already read through pronoun reference (anaphora) or repetition of nouns and verbs and provide clues to future meaning (e.g., *therefore, nonetheless*). Knowledge of the meaning and use of connectives is an important aid to comprehension (Cain & Nash, 2011; Crosson & Lesaux, 2013).

Reading comprehension. If a student can read and understand the meanings of printed words and sentences, then comprehending text should not be difficult, given the emphasis above on achieving the alphabetic principle, lexical quality, and syntactic awareness. Individual differences in readers' background knowledge, motivation, and memory and attention will create variability in word recognition skills, vocabulary knowledge, and syntactic awareness and this variability, in turn, will create variability in reading comprehension. Furthermore, genre differences—informational or literary text—may interact with reader skills to affect reading comprehension. For example, some students may have

better inferential language skills so critical to comprehending informational text; other students may have better narrative language skills of discerning story structure and character motivation and, therefore, be good comprehenders of literary text. Because reading comprehension is affected by the interactions of variables related to reader and text characteristics (RAND, 2002), tests of reading comprehension typically consist of informational *and* literary passages and provide as much relevant background information within the passage as possible.

States' reading comprehension tests typically have questions written to their state standards. One challenge for these tests are the trade-offs between coverage of the standards, time, and reliability. Typically, one should strive for about 15 items per standard. If a state has 14 standards per grade, then 210 questions would be needed to reliably cover the standards. If 7-9 questions are written for each passage, then students would need to read 23-30 passages, which would take them about 10 days. Most states prioritize testing the superordinate standards in order to reduce the testing time to 7 passages or so over two days. A limitation of many standards-based tests is their sole focus on grade-level proficiency. Students are given only grade-level passages; therefore, students who read below grade level tend to guess and students who read above grade level are not challenged. In both cases, no information about their actual reading ability is obtained. Furthermore, when the grade level of passages is determined by readability formulae or by qualitative ratings, the precision is not at a particular grade but rather within grade bands of two to three grades (e.g., upper elementary, middle school, high school; Foorman, 2009; Nelson, Perfetti, Liben, & Liben, 2012).

The FAIR-FS Reading Comprehension task in grades 3-12 avoids the problems with precision and efficiency noted above by being a computer-adaptive test. Students are placed into their first reading comprehension passage based on their ability on the computer-adaptive Word Recognition and Vocabulary Knowledge Tasks—which take 2-3 minutes each. The student reads the passage and answers the 7-9 multiple choice questions. Subsequent passage placement is based on relations among student ability, standard error, and discrimination parameters from a 2-parameter logistic item response theory (IRT) model. Students continue to receive passages until a precise estimate of reading comprehension is achieved (i.e., reliability $>.80$). In the FAIR-FS, students receive 1-3 passages in about 10-30 minutes. Given that the two Screening tasks and one Diagnostic task take, on average, 11 minutes, the entire 3-12 battery easily fits into a 45-minute class period. During the 2013-2014 implementation study in Pinellas County, reliability on the Reading Comprehension task was above .80 for 93 percent of students and above .90 for 54 percent of students.

Individual tasks in the FAIR-FS yield two score types—percentile ranks and ability scores. The ability score is used to measure growth and can be displayed against grade-level percentile ranks to communicate the important point that students are improving across the year even though they are performing far below or above grade-level peers.

Summary of FAIR-FS Constructs and Tasks

The FAIR-FS consists of computer-adaptive reading comprehension and oral language screening tasks that provide measures to track growth over time, as well as a Probability of Literacy Success (PLS) linked to grade-level performance (i.e., the 40th percentile) on the reading comprehension subtest of the Stanford Achievement Test (SAT-10) in the 2014-2015 school year and will predict to the Florida Standards Assessment once those data are available. Thus, the FAIR-FS provides universal screening and diagnostic tasks in a precise and efficient computer-adaptive framework with psychometrics and norms derived from large samples of Florida 3-12 students representative of Florida demographics. By including Vocabulary Knowledge and Syntax Knowledge Tasks, the FAIR-FS has excellent construct coverage of oral language, which has been shown to account for the vast majority (i.e., 72%-96%, with a median of 87%) of individual differences in reading comprehension in grades 4-10 (Foorman, Koon, Petscher, Mitchell, & Truckenmiller, 2014).

Description of the Tasks in the FAIR-FS

In grades 3 through 12, the FAIR-FS consists of four computer-adaptive tasks that each provide unique information regarding a student's literacy skills. Each of the tasks below, except for Reading Comprehension, have four stop rules that determine when administration of each task is complete¹.

1. A reliable estimate of the student's abilities is reached (i.e., standard error is less than 0.316).
2. The student has responded to 30 items.
3. The student responds correctly to all of the first 8 items.
4. The student responds incorrectly to all of the first 8 items.

At subsequent administrations of the tasks within the same school year, the student's prior score on that task determines the initial set of items administered to the student at that administration period.

The tasks in the FAIR-FS can be used as a highly efficient diagnostic tool due to the utilization of computer adaptive functionality. Computer administration allows for large groups of students to be assessed at once with a high degree of standardization. Adaptability in the items allows for a highly reliable score to be reached sooner and decreases the amount of time needed for each task. Although educators are most concerned with students' abilities in reading comprehension, it is a complex skill that takes significant amounts of time to assess (due to close reading of extended text) and poor performance does not necessarily signal which component skills of reading to target for instruction. The FAIR-FS efficiently assesses multiple research-based component skills of reading comprehension to help

¹ The stop rules for reading comprehension are a maximum of three passages or a reliable estimate of the student's ability (i.e., standard error < .316).

teachers diagnose skill weaknesses and target instruction. During the implementation study, more than 98% of students reached a highly reliable score (marginal reliability above .80) by taking an average of only 20 items on the WRT, 9 items on the VKT, and 18 items on the SKT. The increase in efficiency allows for more tasks to be administered to achieve a more complete diagnostic profile for a student. For example, in the implementation study 84% of students in grades 3 through 12 completed all four of the computer-adaptive tasks within one class period (i.e., 45 minutes).

Word Recognition Task (WRT). In the Word Recognition Task, the student listens to a word pronounced by the computer. The computer monitor displays a drop-down menu with the correctly spelled word and two distractors that are spelled incorrectly. The student may replay the audio for the word up to three times. The student has unlimited time to respond to each item. The item bank contains 274 available items and includes real words and some non-words.

Vocabulary Knowledge Task (VKT). Each item in the Vocabulary Knowledge Task consists of one sentence with a word missing. The missing word is replaced with a choice of three morphologically related words. The student selects the word that best completes the sentence. There are 374 items available. The student has unlimited time to respond to each item.

Reading Comprehension (RC). The Reading Comprehension task consists of passages that are between 200 and 1300 words in length. Each passage has between 7 and 9 multiple choice questions. Each question has one correct response and three distractors. All questions associated with the passage are displayed at the same time and the passage is also available on the computer monitor. Each question has an individual item difficulty and discrimination value. Each set of 7 to 9 questions has an average item difficulty, which is used to determine which set of questions (and associated passage) is administered to the student next. The Reading Comprehension task ends when a reliable score has been reached (i.e., the standard error is less than 0.316) or the student has responded to three sets of questions. The initial set of questions administered to a student is determined by a formula that includes the student's score on the WRT and the VKT. The computer will automatically log out students after 15 minutes of inactivity; otherwise, students have an unlimited amount of time to read the passage and respond to questions. There are a total of 139 sets of questions associated with passages available in the grades 3-12 FAIR-FS.

Syntactic Knowledge Task (SKT). In the Syntactic Knowledge Task, the student listens to a sentence or sentences read by the computer that is missing one word. The computer monitor also displays the sentence(s) for the student to read along. The missing word(s) in the sentence(s) is replaced by a dropdown box with the correct word or phrase and two distractors. There are a total of 240 items available. Some items require a student to select the correct connective word, the correct pronoun reference, or the correct verb that creates appropriate subject-verb agreement.

Description of Method

Item tryout and validation work with the above tasks occurred from 2010-2015 through the funding provided by two IES grants (see Acknowledgements). Once item writers had written items for each task, tasks were piloted with students in grades 3-12. Results from Item Response Theory (IRT) analyses were evaluated and in several cases items were deleted or more difficult items were written and further field trials were conducted. A large-scale linking study was conducted during the Spring of 2013 with approximately 45,000 students in grades 3 through grade 12 in two districts in Florida. Outcome data consisted of well-known standardized measures of reading comprehension (Gates-MacGinitie and the SAT-10). Item response and differential item function analyses were conducted. Parameters derived from these analyses are used in the look-up tables in the computer-adaptive system.

Item Response Theory

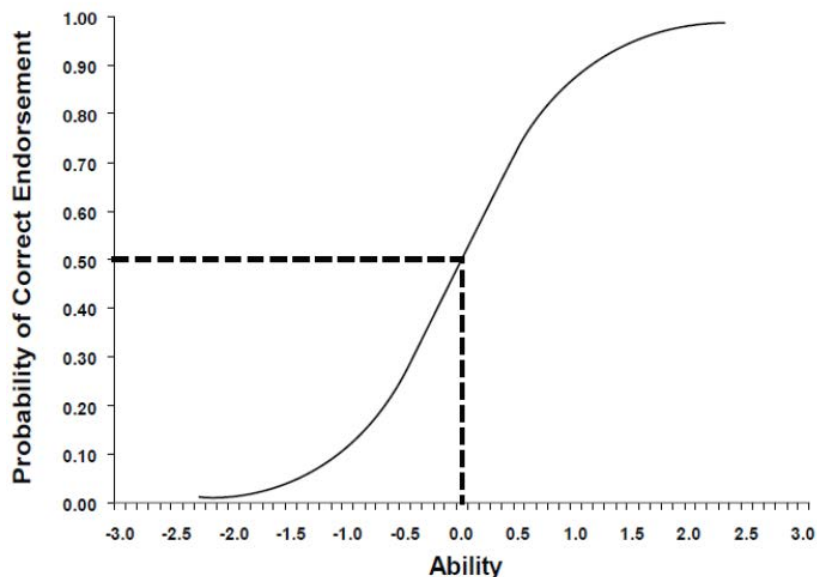
Data for the grades 3-12 FAIR-FS were analyzed using Item Response Theory (IRT). Traditional testing and analysis of items involves estimating the difficulty of the item (based on the percentage of respondents correctly answering the item) as well as discrimination (how well individual items relate to overall test performance). This falls into the realm of measurement known as classical test theory (CTT). While such practices are commonplace in assessment development, IRT holds several advantages over CTT. When using CTT, the difficulty of an item depends on the group of individuals on which the data were collected. This means that if a sample has more students that perform at an above-average level, the easier the items will appear; but if the sample has more below-average performers, the items will appear to be more difficult. Similarly, the more that students differ in their ability, the more likely the discrimination of the items will be high; the more that the students are similar in their ability, the lower the discrimination will be. One could correctly infer that scores from a CTT approach are entirely dependent on the makeup of the sample on which the items are tested.

The benefits of IRT are such that: 1) the difficulty, discrimination, and pseudo-guessing parameters are not dependent on the group(s) from which they were initially estimated; 2) scores describing students' ability are not related to the difficulty of the test; 3) shorter tests can be created that are more reliable than a longer test; and, 4) item statistics and the ability of students are reported on the same scale.

Item Difficulty. The difficulty of an item has traditionally been described for many tests as a “p-value”, which corresponds to the percent of respondents correctly answering an item. Values from this perspective range from 0% to 100% with high values indicating easier items and low values indicating hard items. Item difficulty in an IRT model does not represent proportion correct, but is rather represented as estimates along a continuum of -3.0 to +3.0. Figure 1 demonstrates a sample item characteristic curve which describes item properties from IRT. Along the x-axis is the ability of the individual, denoted by theta. As previously mentioned, the ability of students and item statistics are reported on the same scale. Thus, the x-axis is a simultaneous representation of student ability and item difficulty. Negative values along the x-axis will indicate that items are easier, while positive values describe harder items. Pertaining to students, negative values describe individuals who perform below average, while positive values identify students who perform above average. A value of zero for both students and items reflects average level of either ability or difficulty.

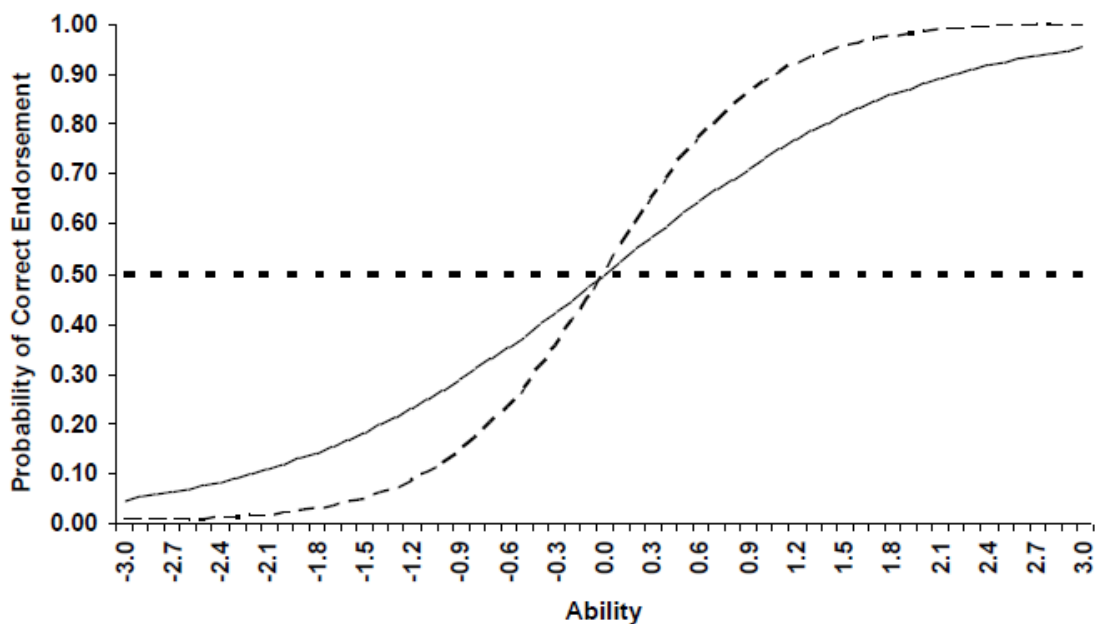
Along the y-axis is the probability of a correct response, which varies across the level of difficulty. Item difficulty is defined as the value on the x-axis at which the probability of correctly endorsing the item is 0.50. As demonstrated for the sample item in Figure 1, the difficulty of this item would be 0.0. Item characteristic curves are graphical representations generated for each item that allow the user to see how the probability of getting the item correct changes for different levels of the x-axis. Students with an ability of -3.0 would have an approximate 0.01 chance of getting the item correct, while students with an ability of 3.0 would have a nearly 99% chance of getting an item correct.

Figure 1: Sample Item Characteristic Curve



Item Discrimination. Item Discrimination is related to the relationship between how a student responds to an item and their subsequent performance on the rest of a test. In IRT it describes the extent to which an item can differentiate the probability of correctly endorsing an item across the range of ability (i.e., -3.0 to +3.0). Figure 2 provides an example of how discrimination operates in the IRT framework. For all three items presented in Figure 2, the difficulty has been held constant at 0.0, while the discriminations are variable. The dashed line (Item 1) shows an item with strong discrimination, the solid line (Item 2) represents an item with acceptable discrimination, and the dotted line (Item 3) is indicative of an item that does not discriminate. It is observed that for Item 3, regardless of the level of ability for a student, the probability of getting the item right is the same. Both high ability students and low ability students have the same chance of doing well on this item. Item 1 demonstrates that as the x-axis increases, the probability of getting the item correct changes as well. Notice that small changes between -1.0 and +1.0 on the x-axis result in large changes on the y-axis. This indicates that the item discriminates well among students, and that individuals with higher ability have a greater probability of getting the item correct. Item 2 shows that while an increase in ability produces an increase in the probability of a correct response, the increase is not as large as is observed for Item 1, and is thus a poorer discriminating item.

Figure 2: Sample Item Characteristic Curves with Varied Discriminations



Guidelines for Retaining Items

Several criteria were used to evaluate item validity. The first process was to identify items which demonstrated strong floor or ceiling effects in response rates $\geq 95\%$. Such items are not useful in

creating an item bank as there is little variability in whether students are successful on the item. In addition to evaluating the descriptive response rate, we estimated item-total correlations. Items with negative values are indicative of poor functioning such that it suggests individuals who correctly answer the question tend to have lower total scores. Similarly, items with low item-total correlations indicate the lack of a relation between item and total test performance. Items with correlations $<.15$ were flagged for removal. Following the descriptive analysis of item performance, difficulty and discrimination values from the IRT analyses were used to further identify items which were poorly functioning. Items were flagged for item revision if the item discrimination was negative or the item difficulty was greater than $+4.0$ or less than -4.0 .

Secondary criteria were used in evaluating the retained items, which was comprised of a differential item function (DIF) analysis. DIF refers to instances where individuals from different groups with the same level of underlying ability significantly differ in their probability to correctly endorse an item. Unchecked, items included in a test which demonstrate DIF will produce biased test results. For the FAIR-FS assessments, DIF testing was conducted comparing: Black-White students, Latino-White students, Black-Latino students, students eligible for Free or Reduced Priced Lunch (FRL) with students not receiving FRL, and English Language Learner to non-English Language Learner students.

DIF testing was conducted with a multiple indicator multiple cause (MIMIC) analysis in Mplus (Muthén & Muthén, 2008); moreover, a series of four standardized and expected score effect size measures were generated using VisualDF software (Meade, 2010) to quantify various technical aspects of score differentiation between the gender groups. First, the signed item difference in the sample (SIDS) index was created, which describes the average unstandardized difference in expected scores between the groups. The second effect size calculated was the unsigned item difference in the sample (UIDS). This index can be utilized as supplementary to the SIDS. When the absolute value of the SIDS and UIDS values are equivalent, the differential functioning between groups is equivalent; however, when the absolute value of the UIDS is larger than SIDS, it provides evidence that the item characteristic curves for expected score differences cross, indicating that differences in the expected scores between groups change across the level of the latent ability score. The D-max index is reported as the maximum SIDS value in the sample, and may be interpreted as the greatest difference for any individual in the sample in the expected response. Lastly, an expected score standardized difference (ESSD) was generated, and was computed similar to a Cohen's (1988) d statistic. As such, it is interpreted as a measure of standard deviation difference between the groups for the expected score response with values of $.2$ regarded as small, $.5$ as medium, and $.8$ as large.

Linking Design & Item Response Analytic Framework

A common-item, non-equivalent groups design was used for collecting data in our pilot, calibration, and validation studies. A strength of this approach is that it allows for linking multiple test forms via common items. For each task, a minimum of twenty-percent of the total items within a form were identified as vertical linking items to create a vertical scale. These items served a dual purpose of not only linking forms across grades to each other, but also linking forms within grades to each other.

FAIR-FS | Description of Method

Because the tasks in the FAIR-FS were each designed for vertical equating and scaling we considered two primary frameworks for estimating the item parameters: 1) a multiple-group IRT of all test forms or 2) test characteristic curve equating. We chose the latter approach using Stocking and Lord (1983) to place the items on a common scale. All item analyses were conducted using Mplus software (Muthen & Muthen, 2008) with a 2pl independent items model. Because the samples used for data collection did not strictly adhere to the state distribution of demographics (i.e., percent limited English proficiency, Black, White, Latino, and eligible for free/reduced lunch), sample weights according to student demographics were used to inform the item and student parameter scores.

Norming Studies

Students from several districts throughout Florida participated in the common-item, non-equivalent groups linking study to estimate and evaluate the item parameters and student ability score distributions for each of the computer adaptive tasks (CAT) in the FAIR-FS. A total of 44,780 students in grades 3-12 across six districts in Florida participated in the calibration and validation studies which consisted of students taking the FAIR-FS tasks appropriate to levels of performance. Table 1 provides a breakdown of the sample sizes used by grade level for each of the FAIR-FS adaptive assessments.

Table 1. Sample Size by Grade for FAIR-FS Tasks

Grade	Vocabulary Knowledge	Word Recognition	Syntactic Knowledge	Reading Comprehension
3	502	651	962	2,723
4	570	586	857	2,679
5	519	697	981	2,721
6	606	652	865	3,835
7	599	612	617	3,683
8	597	613	616	3,814
9	813	1,054	1,053	3,964
10	574	1,109	869	3,787
Total	4,780	5,974	6,820	27,206

Score Definitions

Several different kinds of scores are provided in order to facilitate a diverse set of educational decisions. In this section, we describe the types of scores provided for each measure, define each score, and indicate its primary utility within the decision making framework of the FAIR-FS. An ability score and a percentile rank are provided for each task (WRT, VKT, RC, and SKT) at each time point. One probability of literacy success score is provided at each assessment period.

Probability of Literacy Success (PLS). The Probability of Literacy Success score indicates the likelihood that a student will reach end of year expectations in literacy. For the purposes of the FAIR-FS in the 2014-2015 school year, reaching expectations is defined as performing at or above the 40th percentile on the Stanford Achievement Test, Tenth Edition (SAT-10). The PLS is used to determine which students are at-risk for meeting grade level expectations by the end of the school year. In addition to providing a precise probability of reaching grade level outcomes, the PLS is color-coded:

- red = the student is at high risk and needs supplemental and/or intensive instruction targeted to the student's skill weaknesses
- yellow = the student may be at-risk and educators may consider differentiating instruction for the student and/or providing supplemental instruction
- green = the student is likely not at-risk and will continue to benefit from strong universal instruction

In the grades 3-12 FAIR-FS, the components that are included in the PLS are an aggregate of the individual student's VKT, WRT, and RC scores.

Percentile Ranks. Percentile ranks can vary from 1 to 99, and they divide the distribution of scores from a large standardization sample (in this case a representative sample of students from Florida) into 100 groups that contain approximately the same number of observations in each group. Thus, a sixth grade student who scored at the 60th percentile would have obtained a score better than about 60% of the students in the standardization sample. The median percentile rank on all the tests of the grades 3-12 FAIR-FS is 50, which means that half the students in the standardization sample obtained a score above that point, and half scored below it. The percentile rank is an ordinal variable meaning that it cannot be added, subtracted, used to create a mean score, or in any other way mathematically manipulated. The median is always used to describe the midpoint of a distribution of percentile ranks. Since this score compares a student's performance to other students within a grade

level, it is meaningful in determining the skill strengths and skill weaknesses for a student as compared to other students' performance.

Ability Scores. Each computer-adaptive task has an associated ability score. The ability score provides an estimate of a student's development in a particular skill. This score is sensitive to changes in a student's ability as skill levels increase or decrease. Ability scores in the grades 3-12 FAIR-FS span the development of each of four important skills: Word Recognition, Vocabulary Knowledge, Reading Comprehension, and Syntactic Knowledge. The range of the developmental scale for each task is 200 to 1000, with a mean of 500 and standard deviation of 100. This score has an equal interval scale that can be added, subtracted, and used to create a mean score. Therefore, this is the score that should be used to determine the degree of growth in a skill for individual students.

Reliability

Marginal Reliability

Reliability describes how consistent test scores will be across multiple administrations over time, as well as how well one form of the test relates to another. Because the FAIR-FS uses Item Response Theory (IRT) as its method of validation, reliability takes on a different meaning than from a Classical Test Theory (CTT) perspective. The biggest difference between the two approaches is the assumption made about the measurement error related to the test scores. CTT treats the error variance as being the same for all scores, whereas the IRT view is that the level of error is dependent on the ability of the individual. As such, reliability in IRT becomes more about the level of precision of measurement across ability, and it may sometimes be difficult to summarize the precision of scores in IRT with a single number. Although it is often more useful to graphically represent the standard error across ability levels to gauge the range of abilities for which the test is more or less informative, it is possible to estimate a generic estimate of reliability known as marginal reliability (Sireci, Thissen, & Wainer, 1991) with:

$$\bar{\rho} = \frac{\sigma_{\theta}^2 - \overline{\sigma_{e^*}^2}}{\sigma_{\theta}^2}$$

where σ_{θ}^2 is the variance of ability score for the normative sample and $\overline{\sigma_{e^*}^2}$ is the mean-squared error. Marginal reliability coefficients for the three FAIR-FS Screening tasks are reported in Table 2 by grade and assessment period.

Table 2. Marginal Reliability for FAIR-FS Screening Tasks of Vocabulary Knowledge, Word Recognition, and Reading Comprehension at the Fall, Winter, and Spring Administrations.

Grade	Vocabulary Knowledge			Word Recognition			Reading Comprehension		
	Fall	Winter	Spring	Fall	Winter	Spring	Fall	Winter	Spring
3	.84	.86	.87	.73	.85	.89	.85	.86	.83
4	.81	.83	.86	.86	.84	.88	.76	.85	.89
5	.87	.87	.88	.87	.84	.90	.80	.83	.90
6	.85	.85	.86	.86	.85	.91	.84	.87	.91
7	.85	.85	.86	.86	.86	.91	.78	.83	.91
8	.83	.84	.84	.87	.83	.92	.81	.85	.92
9	.85	.82	.86	.88	.80	.91	.67	.78	.91
10	.85	.81	.84	.88	.78	.90	.76	.82	.92
All Grades	.91	.89	.90	.92	.88	.93	.86	.88	.93

Note. Reliability coefficients for the Fall and Winter Reading Comprehension scores are reflective of fixed item administrations. Spring reliability coefficients for Reading Comprehension are reflective of performance on the CAT version. Marginal reliability coefficients for Vocabulary and Word Recognition are reflective of CAT versions of the assessments.

Across all grades and assessment periods, the marginal reliability was quite high ranging from .86 for fall reading comprehension to .93 for spring word recognition and reading comprehension. Values of .80 are typically viewed as acceptable for research purposes while estimates at .90 or greater are acceptable for clinical decision making (Nunnally & Berstein, 1994). Marginal reliability coefficients for the diagnostic Syntactic Knowledge Task are reported in Table 3. Similar to the other tasks, marginal reliability coefficients were quite high across all grades ranging from .92 to .93.

Table 3. Syntactic Knowledge Marginal Reliability Coefficients

Grade	Syntax		
	Fall	Winter	Spring
3	.85	.87	.89
4	.88	.87	.88
5	.87	.88	.90
6	.88	.89	.91
7	.88	.89	.91
8	.91	.88	.92
9	.91	.87	.90
10	.91	.87	.90
All Grades	.93	.92	.93

Note. Reliability coefficients for all assessment periods are reflective of the CAT version of the assessment

Standard Error of Measurement

A standard error of measurement (SEM) is an estimate that captures the amount of variance that might be observed in an individual student's performance if they were tested repeatedly. That is, on any particular day of testing, an examinee's score may fluctuate and only through repeated testing is it possible to get closer to one's true ability. Because it is not reasonable to test a student enough to capture his/her true ability, we can construct an interval by which we can observe the extent to which the score may fluctuate. The SEM is calculated with:

$$SEM = \sigma_x \sqrt{1 - \rho^2}$$

where σ_x is the standard deviation associated with the mean for assessment x , and ρ^2 is the marginal reliability for the assessment. Means and SEM are reported in Tables 4-7 for the 3 Screening tasks, respectively.

Table 4. Means and Standard Error of Measurement for Vocabulary Knowledge Scores.

Grade	N	Fall		Winter		Spring	
		Mean	SEM	Mean	SEM	Mean	SEM
3	466	380.28	29.30	393.07	27.98	413.82	25.91
4	486	431.77	28.42	439.80	28.63	453.59	26.85
5	423	469.14	29.17	473.85	28.12	482.07	26.89
6	639	492.40	29.23	498.09	29.17	505.10	27.05
7	632	521.95	29.24	518.13	29.34	529.92	26.97
8	681	550.11	29.60	540.88	30.88	551.98	29.40
9	1014	555.66	29.40	560.26	32.00	562.86	28.62
10	887	571.88	30.28	575.32	36.19	574.38	30.44

Table 5. Means and Standard Error of Measurement for Word Recognition Scores.

Grade	N	Fall		Winter		Spring	
		Mean	SEM	Mean	SEM	Mean	SEM
3	470	341.36	29.72	351.25	29.79	377.59	24.21
4	491	407.69	31.06	405.81	30.43	427.49	29.73
5	426	437.77	30.92	440.94	30.42	466.91	27.06
6	646	465.32	31.28	458.53	31.06	490.20	26.41
7	634	498.42	32.22	482.32	31.74	518.74	27.85
8	690	531.50	32.88	515.55	36.63	555.32	27.06
9	1017	543.01	33.21	543.53	43.68	567.72	29.29
10	916	574.34	33.96	558.00	47.27	591.01	32.76

Table 6. Means and Standard Error of Measurement for Reading Comprehension Scores.

FAIR-FS | Reliability

Grade	N	Spring	
		Mean	SEM
3	325	386.03	28.69
4	322	440.07	32.96
5	302	497.25	36.49
6	431	499.96	37.63
7	426	524.45	39.67
8	461	571.71	48.61
9	703	583.06	39.26
10	626	589.72	44.65

Note. Data is only provided for Spring due to the CAT version only being administered in the Spring.

Means and standard error of measurement for the diagnostic Syntactic Knowledge Task are reported in Table 7.

Table 7. Means and Standard Error of Measurement for Syntactic Knowledge Scores.

Grade	N	Fall		Winter		Spring	
		Mean	SEM	Mean	SEM	Mean	SEM
3	377	328.84	30.80	358.06	30.58	402.12	25.29
4	376	403.74	30.06	417.15	30.80	452.63	24.85
5	340	430.52	30.12	452.58	30.82	483.09	25.29
6	383	456.01	31.18	473.15	31.59	505.59	25.04
7	396	510.01	30.40	504.94	31.41	529.24	25.49
8	380	523.01	30.16	533.04	34.28	554.57	25.73
9	457	554.38	32.05	551.09	36.27	571.61	27.52
10	443	554.98	31.07	549.89	38.55	562.49	28.15

Test-Retest Reliability

The extent to which a sample of students performs consistently on the same assessment across multiple occasions is an indication of test-retest reliability. Reliability was estimated for students participating in the field testing of the FAIR-FS by correlating their ability scores across three assessments. Retest correlations for vocabulary and word recognition (Table 8) were the strongest between winter and spring while the fall-winter correlations were strongest for reading comprehension. Correlations between the fall and spring were the lowest, which is expected as a weaker correlation from the beginning of the year to the end suggests that students were differentially changing over time (i.e., lower ability students may have grown more over time compared to higher ability students). Retest correlations for the diagnostic Syntactic Knowledge Task are reported in Table 9. Similar to the Vocabulary Knowledge and Word Recognition Tasks, the strongest correlations between time-points were the winter-spring associations.

Table 8. FAIR-FS Screening Test-Retest Correlations for Vocabulary Knowledge, Word Recognition, and Reading Comprehension.

Grade	Vocabulary Knowledge			Word Recognition			Reading Comprehension		
	Fall-Winter	Winter-Spring	Fall-Spring	Fall-Winter	Winter-Spring	Fall-Spring	Fall-Winter	Winter-Spring	Fall-Spring
3	.59	.61	.44	.46	.51	.31	.74	.66	.66
4	.58	.62	.51	.59	.62	.45	.83	.77	.71
5	.75	.74	.65	.63	.73	.64	.83	.77	.73
6	.60	.72	.51	.59	.65	.66	.85	.80	.77
7	.66	.69	.54	.65	.69	.73	.80	.79	.73
8	.63	.67	.63	.66	.72	.74	.81	.79	.71
9	.65	.64	.65	.65	.68	.76	.77	.72	.65
10	.62	.70	.64	.69	.70	.80	.75	.74	.66

Table 9. Test-Retest Correlations for Syntactic Knowledge Task.

Syntax			
Grade	Fall-Winter	Winter-Spring	Fall-Spring
3	.49	.55	.48
4	.62	.70	.56
5	.68	.75	.68
6	.63	.69	.65
7	.68	.74	.69
8	.66	.76	.70
9	.70	.73	.80
10	.67	.70	.72

Validity

Predictive Validity

The predictive validity of the Screening tasks to the SAT-10 Reading Comprehension test for grades 3-12 was addressed through a series of linear and logistic regressions. The linear regressions were run two ways. First, a correlation analysis was used to evaluate the strength of relations between each of the Screening tasks ability scores with the SAT-10. Second, a multiple regression was run to estimate the total amount of variance that the linear combination of the predictors explained in SAT-10 reading comprehension performance. Results from the linear regression analyses are reported in Table 10.

Table 10. Bivariate Correlations between FAIR-FS Screening Tasks and SAT-10. Percent Variance Explained in SAT-10 by FAIR-FS Vocabulary, Word Recognition, and Reading Comprehension.

Grade	Vocabulary Knowledge	Word Recognition	Reading Comprehension	Total R^2
3	.56	.43	.74	.62
4	.45	.39	.71	.56
5	.57	.41	.74	.59
6	.53	.46	.71	.53
7	.43	.43	.66	.45
8	.46	.47	.67	.48
9	.51	.55	.60	.47
10	.47	.51	.57	.39

For the logistic regressions, students' performance on the SAT-10 Reading Comprehension test was coded as '1' for performance at or above the 40th percentile, and '0' for scores below this target. This dichotomous variable was then regressed on a combination of vocabulary knowledge, word recognition, and reading comprehension scores at each grade level. Further, we evaluated the classification accuracy of scores from the FAIR-FS as it pertains to risk status on the SAT-10. By dichotomizing the combination of screening task scores as '1' for not at-risk for reading difficulties and '0' for at-risk for reading difficulties, students could be classified based on their dichotomized performances on both. As such, students could be identified as not at-risk on the combination of screening tasks and demonstrating

grade level performance on the SAT-10 (i.e., specificity or true-negatives), at-risk on the combination of screening task scores and below grade level performance on the SAT-10 (i.e., sensitivity or true-positives), not at-risk based on the combination of screening task scores and not at grade level on the SAT-10 (i.e., false negative error), or at-risk on the combination of screening task scores and at grade level on the SAT-10 (i.e., false positive error). Classification of students in these categories allows for the evaluation of cut-points on the combination of screening tasks (i.e., PLS) to determine which PLS cut-point maximizes predictive power

The concept of risk can be viewed in many ways, including the concept as a “percent chance” which is a number between 0 and 100, with 0 meaning there is no chance that a student will develop a problem, and 100 being there is no chance the student will not develop a problem. When attempting to identify children who are “at-risk” for poor performance on some type of future measure of reading achievement, this is typically a yes/no decision based upon a “cut-point” along a continuum of risk. Oftentimes this future measure of achievement is a state’s high-stakes assessment, which typically provides a standard score that describes the performance of each student. Grade-level cut-points are chosen that determine whether a student has passed or failed the state-wide assessment.

Decisions concerning appropriate cut-points for screening measures are made based on the level of correct classification that is desired from the screening assessments. While a variety of statistics may be used to guide such choices (e.g., sensitivity, specificity, positive and negative predictive power; see Schatschneider, Petscher & Williams, 2008), negative predictive power was utilized to develop the FAIR-FS cut-points. Negative predictive power is the percentage of students who are identified as “not at-risk” on the screening assessments that end up not passing based the outcome assessment. Predictive power is not considered to be a property of the screening assessments since it is known to fluctuate given the proportion of individuals who are at-risk on the selected outcome (Streiner, 2003).

The cut-point selected for the grades 3-12 FAIR-2009 (used in the State of Florida from 2009-2014, Florida Department of Education, 2009) was negative predictive power of 0.85, meaning that at least 85% of students identified as “not at-risk” on the FAIR-2009 (i.e., $FSP \geq 0.85$) would achieve at least a Level 3 on the Florida Comprehensive Assessment Test (FCAT) reading assessment at the end of the year. Greater emphasis was placed on negative predictive power than positive predictive power because the consequences of being identified as “at-risk” when the student is not actually at-risk are so much less than identifying students as “not at-risk” when they are actually at-risk for below grade-level performance on the FCAT. Prior research (Foorman & Petscher, 2010a; Foorman & Petscher, 2010b; Petscher & Foorman, 2011) demonstrated the technical adequacy of using .85 as an appropriate cut-point for risk on the FAIR 2009. As part of a continuing evaluation of the classification accuracy of FAIR 2009 scores, Petscher and Foorman (2011) found that an alternative cut-point (i.e., .70) could be used to maintain high negative predictive power and also minimize identification errors. As it pertains to the FAIR-FS, we tested the extent to which using a .85 cut-point for a student being identified as not at-risk yielded a negative predictive power value of at least 85%. Similarly, we also tested (a) how high negative

predictive power would be estimated when using a cut-point of .70, and (b) whether identification errors could be reduced. A summary of the classification results for FAIR-FS are reported in Table 11.

Table 11. Classification Accuracy of the Probability of Literacy Success (PLS) in Grades 3-12 using .85 and .70 Cut-Points.

Cut-Point	Grade	SE	SP	PPP	NPP	OCC	Base Rate
.85	3	.95	.54	.59	.94	.71	.41
	4	.95	.58	.52	.96	.70	.32
	5	.94	.60	.56	.95	.72	.35
	6	.96	.39	.61	.91	.68	.50
	7	.98	.46	.55	.97	.67	.40
	8	.94	.46	.54	.92	.64	.39
	9	.93	.50	.38	.96	.61	.25
	10	.87	.52	.42	.91	.62	.28
.70	3	.85	.69	.66	.87	.76	.41
	4	.77	.74	.59	.88	.75	.32
	5	.83	.76	.65	.89	.78	.35
	6	.92	.56	.68	.87	.86	.50
	7	.91	.60	.61	.91	.73	.40
	8	.85	.67	.62	.88	.74	.39
	9	.76	.69	.45	.90	.71	.25
	10	.64	.74	.49	.84	.71	.28

Note. SE= Sensitivity, SP = Specificity, PPP = Positive Predictive Power, NPP = Negative Predictive Power, OCC = Overall Correct Classification. Students in Grades 11 and 12 are classified according to Grade 10 criteria.

Note that when using either the .85 or .70 cut-points the negative predictive power is above .85; yet, when the .85 cut-point is used, the specificity and positive predictive power are relatively low. The consequence of a low specificity value is that many students are required to take one or more additional tasks; in the present sample this would result in between 40% and 61% of students identified as false

positives and required to take the Diagnostic tasks. Conversely, if a .70 cut-point is used, this error rate range reduces from 40%-61% down to 24% -44%. Coupled with a false positive reduction is an increase in the positive predictive power and the overall correct classification. Although there is some loss of precision in the sensitivity, the negative predictive power maintains a high value to ensure that students who are identified as not at-risk have a high likelihood of being successful on end of year outcomes (i.e., 40th percentile or greater on the SAT-10).

Differential Accuracy of Prediction

An additional component of checking the validity of cut-points and scores on the assessments involved testing differential accuracy of the regression equations across different demographic groups. This procedure involved a series of logistic regressions predicting success on the SAT-10 test (i.e., at or above the 40th percentile). The independent variables included a variable that represented whether students were identified as not at-risk (PLS \geq .70; coded as '1') or at-risk (PLS $<$.70; coded as '0') on the combination of screening task scores, a variable that represented a selected demographic group, as well as an interaction term between the two variables. A statistically significant interaction term would suggest that differential accuracy in predicting end-of-year performance existed for different groups of individuals based on the risk status determined by the screening assessment. For the combination of FAIR-FS screening task scores, differential accuracy was separately tested for Black and Latino students as well as for students identified as English Language Learners (ELL) and students who were eligible for Free/Reduced Price Lunch (FRL).

When testing for differential accuracy between Black and White students (Table 12), a significant effect for the interaction between the PLS cut-point and minority status existed in grade 4 ($p = .003$). This finding indicated that for the sample tested at the winter assessment period, White students with a PLS above the cut-point had a 92% chance of being at or above the 40th percentile on the SAT-10 compared to Black students above the cut-point on the PLS who had a 76% chance of being at or above the 40th percentile on the SAT-10. This translates into a 16% advantage in success for White students in grade 4, but we should note that replication will be needed across multiple administrations with a larger sample to evaluate the extent to which this phenomenon continues to exist.

When testing for differential accuracy between Hispanic and White students (Table 13), a significant effect for the interaction between the PLS cut-point and minority status existed in grades 8 and 10 ($p = .015$ and $.02$, respectively). This finding indicated that for the sample tested at the winter, White students in grade 8 with a PLS above the cut-point had an 87% chance of being at or above the 40th percentile on the SAT-10 compared to Hispanic students above the cut-point on the PLS who had an 89% chance of being at or above the 40th percentile on the SAT-10. This translates into a 3% advantage in success for Hispanic students in grade 8. Similarly, White students in grade 10 with a PLS above the cut-point had an 82% chance of being at or above the 40th percentile on the SAT-10 compared to Hispanic students above the cut-point on the PLS who had an 86% chance of being at or above the 40th percentile on the SAT-10. This translates into a 4% advantage in success for Hispanic students in grade 10. The findings from these two grades should be interpreted with caution as the mean difference in

expected probability scores is quite small; thus, replication will be needed across multiple administrations with a larger sample to evaluate the extent to which this phenomenon continues to exist.

Table 12. Differential Accuracy for FAIR-FS Screening Tasks by Grade: Black-White (BW)

Grade	Parameter	df	Estimate	SE	χ^2	<i>p</i> -value
3	Intercept	1	-0.33	0.28	1.39	0.239
	PLS	1	3.69	0.65	32.12	<.001
	BW	1	-0.19	0.34	0.32	0.573
	PLS *BW	1	-1.31	0.77	2.86	0.091
4	Intercept	1	-0.66	0.33	3.98	0.046
	PLS	1	3.05	0.48	41.02	<.001
	BW	1	0.53	0.40	1.70	0.192
	PLS *BW	1	-1.78	0.60	8.83	0.003
5	Intercept	1	-0.31	0.27	1.37	0.243
	PLS	1	3.06	0.48	40.88	<.001
	BW	1	-0.33	0.34	0.91	0.340
	PLS *BW	1	-0.64	0.61	1.09	0.296
6	Intercept	1	-0.41	0.17	6.01	0.014
	PLS	1	2.62	0.34	59.29	<.001
	BW	1	-0.48	0.26	3.34	0.068
	PLS *BW	1	-0.85	0.57	2.22	0.137
7	Intercept	1	-0.31	0.18	2.98	0.085
	PLS	1	3.10	0.44	48.81	<.001
	BW	1	-0.14	0.28	0.25	0.615
	PLS *BW	1	-0.94	0.62	2.28	0.131
8	Intercept	1	-0.10	0.17	0.34	0.563
	PLS	1	1.97	0.29	46.72	<.001

	BW	1	-0.39	0.26	2.21	0.137
	PLS *BW	1	-0.09	0.49	0.04	0.849
9	Intercept	1	0.28	0.22	1.62	0.203
	PLS	1	2.31	0.42	30.23	<.001
	BW	1	-0.25	0.33	0.59	0.442
	PLS *BW	1	-0.38	0.59	0.42	0.517
10	Intercept	1	0.55	0.23	5.48	0.019
	PLS	1	0.99	0.30	11.05	0.001
	BW	1	-0.71	0.32	4.90	0.027
	PLS *BW	1	0.53	0.44	1.43	0.233

Note. PLS cut-off is .70. Estimates based on .85 cut-off approximate .70 results. PLS scores are based on student performance at the winter administration.

Table 13. Differential Accuracy for Screening Tasks by Grade: Hispanic-White (HW)

Grade	Parameter	df	Estimate	SE	χ^2	p-value
3	Intercept	1	-0.33	0.28	1.39	0.239
	PLS	1	3.69	0.65	32.12	<.001
	HW	1	-0.55	0.31	3.07	0.080
	PLS*HW	1	-1.32	0.70	3.60	0.058
4	Intercept	1	-0.66	0.33	3.98	0.046
	PLS	1	3.05	0.48	41.02	<.001
	HW	1	0.29	0.37	0.60	0.439
	PLS*HW	1	-0.56	0.55	1.04	0.307
5	Intercept	1	-0.31	0.27	1.37	0.243
	PLS	1	3.06	0.48	40.88	<.001
	HW	1	-0.39	0.30	1.63	0.202
	PLS*HW	1	-0.48	0.54	0.80	0.371
6	Intercept	1	-0.41	0.17	6.01	0.014
	PLS	1	2.62	0.34	59.29	<.001
	HW	1	-0.47	0.21	5.15	0.023
	PLS*HW	1	0.66	0.51	1.65	0.199
7	Intercept	1	-0.31	0.18	2.98	0.085
	PLS	1	3.10	0.44	48.82	<.001
	HW	1	-0.19	0.23	0.68	0.408
	PLS*HW	1	-0.37	0.55	0.44	0.509
8	Intercept	1	-0.10	0.17	0.34	0.563
	PLS	1	1.97	0.29	46.72	<.001

	HW	1	-0.72	0.22	10.20	0.001
	PLS*HW	1	0.98	0.40	5.96	0.015
9	Intercept	1	0.28	0.22	1.62	0.203
	PLS	1	2.31	0.42	30.23	<.001
	HW	1	-0.01	0.29	0.00	0.974
	PLS*HW	1	-0.59	0.52	1.28	0.258
10	Intercept	1	0.55	0.23	5.48	0.019
	PLS	1	0.99	0.30	11.05	0.001
	HW	1	-0.67	0.29	5.18	0.023
	PLS*HW	1	0.95	0.41	5.41	0.020

Note. PLS cut-off is .70. Estimates based on .85 cut-off approximate .70 results. PLS scores are based on student performance at the winter administration.

When testing for differential accuracy between ELL and non-ELL students (Table 14), a significant effect for the interaction between the PLS cut-point and ELL status existed in grade 5 ($p = .01$). This finding indicated that for the sample tested at the winter, non-ELL students with a PLS above the cut-point had a 90% chance of being at or above the 40th percentile on the SAT-10 compared to ELL students above the cut-point on the PLS who had a 61% chance of being at or above the 40th percentile on the SAT-10. This translates into a 29% advantage in success for non-ELL students in grade 5, but we should note that replication will be needed across multiple administrations with a larger sample to evaluate the extent to which this phenomenon continues to exist.

Table 14. Differential Accuracy for FAIR-FS Screening Tasks by Grade: English Language Learners (ELL)

Grade	Parameter	df	Estimate	SE	χ^2	p -value
3	Intercept	1	-0.42	0.12	12.44	<.001
	PLS	1	2.36	0.20	133.00	<.001
	ELL	1	-1.27	0.30	17.82	<.001
	PLS*ELL	1	0.71	0.66	1.15	0.284
4	Intercept	1	-0.10	0.14	0.57	0.450

	PLS	1	2.09	0.21	99.96	<.001
	ELL	1	-1.00	0.30	11.23	<.001
	PLS*ELL	1	0.24	0.89	0.07	0.788
5	Intercept	1	-0.50	0.13	14.52	<.001
	PLS	1	2.72	0.21	168.37	<.001
	ELL	1	-0.38	0.24	2.46	0.117
	PLS*ELL	1	-1.41	0.54	6.68	0.010
6	Intercept	1	-0.47	0.10	22.46	<.001
	PLS	1	2.46	0.21	134.43	<.001
	ELL	1	-1.37	0.25	29.01	<.001
	PLS*ELL	1	-0.63	0.79	0.63	0.426
7	Intercept	1	-0.08	0.11	0.59	0.441
	PLS	1	2.47	0.24	108.98	<.001
	ELL	1	-1.56	0.27	34.34	<.001
	PLS*ELL	1	-0.01	0.74	0.00	0.991
8	Intercept	1	-0.14	0.10	1.70	0.192
	PLS	1	2.11	0.18	134.92	<.001
	ELL	1	-1.74	0.28	40.22	<.001
	PLS*ELL	1	1.37	0.76	3.28	0.070
9	Intercept	1	0.29	0.13	5.04	0.025
	PLS	1	1.93	0.22	80.32	<.001
	ELL	1	-0.59	0.34	3.00	0.083
	PLS*ELL	1	-1.23	0.91	1.81	0.178
10	Intercept	1	0.20	0.13	2.49	0.114

PLS	1	1.54	0.18	75.19	<.001
ELL	1	-1.16	0.35	11.19	0.001
PLS*ELL	1	-0.63	0.59	1.12	0.291

Note. PLS cut-off is .70. Estimates based on .85 cut-off approximate .70 results. PLS scores are based on student performance at the winter administration.

When testing for differential accuracy between FRL and non- FRL students (Table 15), a significant effect for the interaction between the PLS cut-point and FRL status existed in grade 10 ($p = .002$). This finding indicated that for the sample tested at the winter, non- FRL students with a PLS above the cut-point had a 91% chance of being at or above the 40th percentile on the SAT-10 compared to FRL students above the cut-point on the PLS who had a 75% chance of being at or above the 40th percentile on the SAT-10. This translates into a 16% advantage in success for non-FRL students in grade 10, but we should note that replication will be needed across multiple administrations with a larger sample to evaluate the extent to which this phenomenon continues to exist.

Table 15. Differential Accuracy for Screening Tasks by Grade: Free or Reduced Price Lunch (FRL)

Grade	Parameter	df	Estimate	SE	χ^2	p -value
3	Intercept	1	0.59	0.32	3.56	0.059
	PLS	1	3.11	0.75	17.16	<.001
	FRL	1	-1.45	0.34	18.57	<.001
	PLS*FRL	1	-0.65	0.78	0.70	0.403
4	Intercept	1	1.00	0.41	5.83	0.016
	PLS	1	1.58	0.54	8.63	0.003
	FRL	1	-1.50	0.43	11.99	0.001
	PLS*FRL	1	0.66	0.58	1.29	0.257
5	Intercept	1	-0.17	0.34	0.24	0.623
	PLS	1	2.77	0.47	34.72	<.001
	FRL	1	-0.50	0.36	1.99	0.159
	PLS*FRL	1	-0.22	0.51	0.19	0.664

6	Intercept	1	-0.54	0.19	7.67	0.006
	PLS	1	2.95	0.38	61.37	<.001
	FRL	1	-0.27	0.22	1.53	0.216
	PLS*FRL	1	-0.57	0.45	1.64	0.200
7	Intercept	1	0.29	0.21	1.79	0.180
	PLS	1	2.63	0.44	36.49	<.001
	FRL	1	-0.90	0.24	13.97	0.000
	PLS*FRL	1	-0.10	0.51	0.04	0.836
8	Intercept	1	-0.01	0.19	0.00	0.948
	PLS	1	2.22	0.30	55.92	<.001
	FRL	1	-0.64	0.22	8.52	0.004
	PLS*FRL	1	0.19	0.37	0.26	0.611
9	Intercept	1	0.45	0.21	4.71	0.030
	PLS	1	1.99	0.33	36.53	<.001
	FRL	1	-0.37	0.25	2.13	0.144
	PLS*FRL	1	-0.13	0.42	0.10	0.752
10	Intercept	1	0.08	0.18	0.22	0.642
	PLS	1	2.21	0.27	65.32	<.001
	FRL	1	-0.10	0.23	0.18	0.675
	PLS*FRL	1	-1.08	0.35	9.64	0.002

Note. PLS cut-off is .70. Estimates based on .85 cut-off approximate .70 results. PLS scores are based on student performance at the winter administration.

Construct Validity

Construct validity describes how well scores from an assessment measure the construct it is intended to measure. Components of construct validity include convergent validity, which can be evaluated by testing relations between a developed assessment and another related assessment, and discriminant

validity, which is can be evaluated by correlating scores from a developed assessment with an unrelated assessment. The goal of the former is to yield a high association which indicates that the developed measure converges, or is empirically linked to, the intended construct. The goal of the latter is to yield a lower association, which indicates that the developed measure is unrelated to a particular construct of interest. Reading and language skills tend to have moderate associations between them; thus, the expectation of the FAIR-FS Vocabulary Knowledge, Word Recognition, and Syntactic Knowledge Tasks would be that stronger associations with reading comprehension would be observed compared to more moderate associations with each other. Correlation results are reported in Table 16.

Table 16. Bivariate Associations among FAIR-FS Tasks.

Grade	Measure	Reading Comprehension		Word Recognition	
			Vocabulary		Syntax
3	Reading Comprehension	1.00			
	Vocabulary Knowledge	.60	1.00		
	Word Recognition	.42	.37	1.00	
	Syntax Knowledge	.48	.38	.30	1.00
4	Reading Comprehension	1.00			
	Vocabulary Knowledge	.42	1.00		
	Word Recognition	.43	.30	1.00	
	Syntax Knowledge	.52	.35	.29	1.00
5	Reading Comprehension	1.00			
	Vocabulary Knowledge	.58	1.00		
	Word Recognition	.40	.37	1.00	
	Syntax Knowledge	.57	.44	.31	1.00
6	Reading Comprehension	1.00			
	Vocabulary Knowledge	.54	1.00		
	Word Recognition	.48	.36	1.00	
	Syntax Knowledge	.58	.45	.36	1.00

7	Reading Comprehension	1.00			
	Vocabulary Knowledge	.46	1.00		
	Word Recognition	.45	.38	1.00	
	Syntax Knowledge	.60	.44	.42	1.00
8	Reading Comprehension	1.00			
	Vocabulary Knowledge	.49	1.00		
	Word Recognition	.49	.40	1.00	
	Syntax Knowledge	.59	.44	.46	1.00
9	Reading Comprehension	1.00			
	Vocabulary Knowledge	.53	1.00		
	Word Recognition	.55	.53	1.00	
	Syntax Knowledge	.63	.58	.54	1.00
10	Reading Comprehension	1.00			
	Vocabulary Knowledge	.50	1.00		
	Word Recognition	.49	.51	1.00	
	Syntax Knowledge	.59	.55	.57	1.00

References

- Cain, K., & Nash, H. M. (2011). The influence of connectives on young readers' processing and comprehension of text. *Journal of Educational Psychology, 103*(2), 429-441.
- Carlisle, J. F., Stone, C. A. (2005). Exploring the role of morphemes in word reading. *Reading Research Quarterly, 40*, 428-449.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (Second ed.). Hillsdale: Lawrence
- Crosson, A. C., Lesaux, N. K. (2013). Connectives: Fitting another piece of the vocabulary instruction puzzle. *The Reading Teacher, 67*(3), 193-200.
- Ehri, L. C., Nunes, S., Stahl, S., & Willows, D. (2001). Systematic phonics instruction helps students learn to read: Evidence from the National Reading Panel's meta-analysis. *Review of Educational Research, 71*, 393-447.
- Ehri, L. C., Nunes, S., Willows, D., Schuster, B., Yaghoub-Zadeh, Z., & Shanahan, T. (2001). Phonemic awareness instruction helps children learn to read: Evidence from the national Reading Panel's meta-analysis. *Reading Research Quarterly, 36*, 250-287.
- Florida Department of Education (2009-2011). Florida Assessments for Instruction in Reading (FAIR). Tallahassee, FL: Author.
- Foorman, B. R. (2009). Text difficulty in reading assessment. In E.H. Hiebert (Ed.), *Reading more, reading better* (pp. 231-247.) New York, NY: Guilford.
- Foorman, B. R., & Connor, C. (2011). Primary reading. In M. Kamil, P. D. Pearson, & E. Moje (Eds.), *Handbook on reading research* (Vol. 4, pp. 136–156). New York, NY: Taylor and Francis.
- Foorman, B. R., Petscher, Y., & Bishop, M. D. (2012). The incremental variance of morphological knowledge to reading comprehension in grades 3-10 beyond prior reading comprehension, spelling, and text reading efficiency. *Learning and Individual Differences, 22*, 792-798.
- Foorman, B. R., Koon, S., Petscher, Y., Mitchell, A., & Truckenmiller, A. (2014). Relations among syntax, vocabulary, decoding fluency, and reading comprehension in grades 4-10: A bi-factor approach. Manuscript submitted for publication.
- Foorman, B. R., & Petscher, Y. (2010a). *Summary of the predictive relationship between the FAIR and the FCAT in grades 3-10*. Tallahassee, FL: Florida Center for Reading Research.
- Foorman, B. R., & Petscher, Y. (2010b). *The unique role of the FAIR Broad Screen in predicting FCAT Reading Comprehension*. Tallahassee, FL: Florida Center for Reading Research.
- Meade, A. W. (2010). A taxonomy of effect sizes for the differential functioning of items and scales. *Journal of Applied Psychology, 95*, 728-743.

- Muthén, B., & Muthén, L. (2008). *Mplus User's Guide*. Los Angeles, CA: Muthén and Muthén.
- National Early Literacy Panel. (2008). *Developing early literacy: Report of the National Early Literacy Panel*. Washington, DC: National Institute for Literacy. Retrieved from: <http://lincs.ed.gov/publications/pdf/NELPReport09.pdf>
- National Institute of Child Health and Human Development. (2000). *Report of the National Reading Panel. Teaching children to read: Reports of the subgroups* (NIH Publication No. 00-4754). Washington, DC: U.S. Department of Health and Human Services.
- National Research Council (1998). *Preventing reading difficulties in young children*. Committee on the Prevention of Reading Difficulties in Young Children, Committee on Behavioral and Social Science and Education, C.E. Snow, M.S. Burns, & P. Griffin, eds. Washington, D.C.: National Academy Press.
- Nelson, J., Perfetti, C., Liben, D., Liben, M. (2012). Measures of text difficulty: Testing their predictive value for grade levels and student performance. Retrieved from: http://www.ccsso.org/Documents/2012/Measures%20ofText%20Difficulty_final.2012.pdf
- Nunnally, J. C., & Bernstein, I. H. (1994). *Psychometric theory* (3rd Ed.). New York: McGraw-Hill.
- Perfetti, C., & Stafura, J. (2014). Word knowledge in a theory of reading comprehension. *Scientific Studies of Reading*, 18, 22-37. DOI: 10.1080/10888438.2013.827687
- Petscher, Y., & Foorman, B. R. (2011). *Summary of the predictive relationship between the FAIR and the FCAT in grades 3-10*. Tallahassee, FL: Florida Center for Reading Research.
- RAND Reading Study Group (2002). *Reading for understanding*. Santa Monica, CA: RAND Corporation.
- Rayner, K., Foorman, B. R., Perfetti, C. A., Pesetsky, D., & Seidenberg, M. S. (2001). How psychological science informs the teaching of reading. *Psychological Science in the Public Interest*, 2(2), 31-74.
- Schatschneider, C., Petscher, Y., & Williams, K. M. (2008). How to evaluate a screening process: The vocabulary of screening and what educators need to know (pg. 304-317). In L. Justice & C. Vukelic (Eds.). *Every moment counts: Achieving excellence in preschool language and literacy instruction*. New York: Guilford Press.
- Sireci, S. G., Thissen, D., & Wainer, H. (1991). On the reliability of testlet-based tests. *Journal of Educational Measurement*, 28, 237-247.
- Stocking, M. L., & Lord, F. M. (1983). Developing a common metric in item response theory. *Applied Psychological Measurement*, 7, 201-210.
- Streiner, D. L. (2003). Diagnosing tests: Using and misusing diagnostic and screening tests. *Journal of Personality Assessment*, 81, 209-219.
- Ziegler, J. C., Stone, G. O., & Jacobs, A. M. (1997). What's the pronunciation for _OUGH and the spelling for /u/? A database for computing feedforward and feedback inconsistency in English. *Behavior Research Methods, Instruments, & Computers*, 29, 600-618.