

Technical Report # 1219

**An Examination of Test-Retest, Alternate Form Reliability,
and Generalizability Theory Study of the easyCBM Passage**

Reading Fluency Assessments:

Grade 4

Julie Alonzo

Cheng-Fei Lai

Daniel Anderson

Bitnara Jasmine Park

Gerald Tindal

University of Oregon



behavioral research & teaching

Published by

Behavioral Research and Teaching
University of Oregon • 175 Education
5262 University of Oregon • Eugene, OR 97403-5262
Phone: 541-346-3535 • Fax: 541-346-5689
<http://brt.uoregon.edu>

Note: Funds for this data set used to generate this report come from a federal grant awarded to the UO from the U.S. Department of Education, Institute for Education Sciences: Reliability and Validity Evidence for Progress Measures in Reading. U.S. Department of Education, Institute for Education Sciences. R324A100014. June 2010 - June 2014. The opinions expressed are those of the authors and do not represent views of the Institute or the U.S. Department of Education.

Copyright © 2012. Behavioral Research and Teaching. All rights reserved. This publication, or parts thereof, may not be used or reproduced in any manner without written permission.

The University of Oregon is committed to the policy that all persons shall have equal access to its programs, facilities, and employment without regard to race, color, creed, religion, national origin, sex, age, marital status, disability, public assistance status, veteran status, or sexual orientation. This document is available in alternative formats upon request.

Abstract

This technical report is one in a series of five describing the reliability (test/retest and alternate form) and G-Theory / D-Study research on the easyCBM reading measures, grades 1-5. Data were gathered in the spring of 2011 from a convenience sample of students nested within classrooms at a medium-sized school district in the Pacific Northwest. Due to the length of the results, we present results of each grade level's analysis in its own technical report, sharing a common abstract, introduction, and methods section, while differing in the results and conclusions.

**An Examination of Test-Retest, Alternate Form Reliability, and Generalizability Theory
Study of the easyCBM Passage Reading Fluency Assessments: Grade 4**

Progress monitoring assessments are a key component of many school improvement efforts, including the Response to Intervention (RTI) approach to meeting students' academic needs. In an RTI approach, teachers first administer a screening or benchmarking assessment to identify students who need supplemental interventions to meet grade-level expectations, then use a series of progress monitoring measures to evaluate the effectiveness of the interventions they are using with the students. When students fail to show expected levels of progress (as indicated by "flat line" scores or little improvement on repeated measures over time), teachers use this information to help them make instructional modifications with the goal of finding an intervention or combination of instructional approaches that will enable each student to make adequate progress toward achieving grade-level proficiency on content standards. In such a system, it is critical to have reliable measures that assess the target construct and are sensitive enough to detect improvement in skill over short periods of time.

Conceptual Framework: Curriculum-Based Measurement and Progress Monitoring

Curriculum-based measurement (CBM), long a bastion of special education, is gaining support among general education teachers seeking a way to monitor the progress their students are making toward achieving grade-level proficiency in key skill and content areas. By definition, CBM is a formative assessment approach. By sampling skills related to the curricular content covered in a given year of instruction yet not specifically associated with a particular textbook, CBMs provide teachers with a snapshot of their students' current level of proficiency in a particular content area as well as a mechanism for tracking the progress students make in gaining desired academic skills throughout the year. Historically, CBMs have been very brief

individually administered measures (Deno, 2003; Good, Gruba, & Kaminski, 2002), yet they are not limited to the one minute timed probes with which many people associate them.

In one of the early definitions of CBM, Deno (1987) stated that “the term curriculum-based assessment, generally refers to any approach that uses direct observation and recording of a student’s performance in the local school curriculum as a basis for gathering information to make instructional decisions...The term curriculum-based measurement refers to a specific set of procedures created through a research and development program ... and grew out of the *Data-Based Program Modification* system developed by Deno and Mirkin (1977)” (p. 41). He noted that CBM is distinct from many teacher-made classroom assessments in two important respects: (a) the procedures reflect technically-adequate measures (“they possess reliability and validity to a degree that equals or exceeds that of most achievement tests” (p. 41), and (b) “growth is described by an increasing score on a standard, or constant task. The most common application of CBM requires that a student’s performance in each curriculum area be measured on a single global task repeatedly across time” (p. 41).

In the three decades since Deno and his colleagues introduced CBM, *progress monitoring probes* as they have come to be called, have increased in popularity, and they are now a regular part of many schools’ educational programs (Alonzo, Tindal, & Ketterlin-Geller, & 2006). However, CBMs – even those widely used across the United States – often lack the psychometric properties expected of modern technically-adequate assessments. Although the precision of instrument development has advanced tremendously in the past 30 years with the advent of more sophisticated statistical techniques for analyzing tests on an item by item basis rather than relying exclusively on comparisons of means and standard deviations to evaluate comparability of alternate forms, the world of CBMs has not always kept pace with these statistical advances.

A key feature of assessments designed for progress monitoring is that alternate forms must be as equivalent as possible to allow meaningful interpretation of student performance data across time. Without such cross-form equivalence, changes in scores from one testing occasion to the next are difficult to attribute to changes in student skill or knowledge. Improvements in student scores may, in fact, be an artifact of the second form of the assessment being easier than the form that was administered first. The advent of more sophisticated data analysis techniques (such as the Rasch modeling used in the development of the easyCBM progress monitoring and benchmarking assessments) has made it possible to increase the precision with which we develop and evaluate the quality of assessment tools.

In this technical report, we provide the results of a series of studies to evaluate the technical adequacy of the easyCBM progress monitoring assessments in reading, designed for use with students in Grades 1 - 5. This assessment system was developed to be used by educators interested in monitoring the progress their students make in acquiring skills in the constructs of early literacy (phonemic awareness, phonics), and both word and passage reading fluency. Specifically, we conducted traditional test-retest and alternate form reliability analyses of the easyCBM reading measures. In addition to these more traditional analyses, we applied generalizability theory – a more modern approach to reliability that parses out sources of error variance. As part of the methods section, we briefly outline the purpose and application of generalizability theory.

The easyCBM™ Progress Monitoring Assessments

The online easyCBM™ progress monitoring assessment system, launched in September 2006 as part of a Model Demonstration Center on Progress Monitoring, was initially funded by the Office of Special Education Programs (OSEP). At the time this technical report was

published, there were 92,925 teachers with easyCBM accounts, representing schools and districts spread across every state in the country. During the 2010-2011 school year, the system had an average of 1200 new accounts registered each week, and the popularity of the system continues to grow. In the month of November 2011, alone, 5945 new teachers registered for accounts, with almost 2 million students active on the system at the end of December 2011. The online assessment system provides both universal screener assessments for fall, winter, and spring administration and multiple alternate forms of a variety of progress monitoring measures designed for use in K-8 school settings.

As part of state funding for Response to Intervention (RTI), states need technically-adequate measures for monitoring progress. Given the increasing popularity of the easyCBM online assessment system, it is imperative that a thorough analysis of the measures' technical adequacy be conducted and the results shared with research and practitioner communities. This technical report addresses that need directly, providing the results of a series of studies examining the technical adequacy of the 2009 / 2010 version of the individually-administered easyCBM assessments in reading.

Methods

Data for these analyses were gathered in the spring of 2011 from a convenience sample of students in a mid-sized school district in the Pacific Northwest. Teams of trained research assistants from the University of Oregon administered a battery of easyCBM assessments to students in participating classrooms. Data were gathered in two separate sessions, one week apart. Each day, students were administered a series of alternate forms of grade-appropriate easyCBM assessments in one-on-one settings. Assessors followed standardized administration protocols for all assessments. The assessments were counter-balanced to enable examination of

order effect as well as alternate form reliability, with selected forms repeated across testing sessions, to allow for test-retest analyses. All assessments were administered in the order displayed in Appendix A.

Test-Retest and Alternate Form Reliability

We used bivariate correlations to calculate the test-retest and alternate form reliability of the measures included in this study. These analyses were completed, in part, as a requisite step to the generalizability theory (G-Theory) analyses. That is, the G-Theory analyses treated each form as a random observation from the universe of possible forms. The G-Theory analyses thus assume form equivalence during the d-study prophecy estimations (i.e., the model assumes each form contributes an equal amount to the measurement process, and that any successive forms will likewise contribute an equal amount). The comparability of forms had to first be established to ensure there were no egregious departures.

Generalizability Theory

For our generalizability theory study (G-Study) we calculated the variances associated persons and two facets: forms and occasions. We then conducted decision studies (D-Studies) to help determine the necessary conditions for reliable measurement. In this section we first provide an overview of G- and D-Studies for the two-facet design for readers who may be unfamiliar with the technique. Readers familiar with G-Theory may want to skip this section and proceed to the *G-Theory analyses* section.

G-Theory overview. G-theory designs can be crossed or nested. A crossed design is one that includes students being administered *the same test forms* on both occasions, while a nested design includes students being administered *different test forms* on both occasions. G-studies are usually followed up with decision studies (D-study analyses), which provide the number of

levels needed to obtain adequate measurement for each facet. For example, to obtain reliable estimates of students' ability, should students be administered 1, 2, 3, 4, or 5 forms during any one occasion? Similarly, does increasing the number of occasions increase the reliability of the estimate, and at what point is a reliable estimate obtained? The results of the G-study are analogous to an analysis of variance (ANOVA), while the results of the D-study are similar to a Spearman-Brown prophecy analysis. Ideally, most of the variance in the G-theory analysis would be associated with persons, and administering students one test form on one occasion would result in sufficiently reliable estimates for the D-study.

Absolute and relative error variances are produced during the D-study. The absolute error variance is the sum of all variance components minus the variance uniquely associated with persons. That is

$$\sigma_{\Delta}^2 = \frac{\sigma_F^2}{n'_F} + \frac{\sigma_O^2}{n'_O} + \frac{\sigma_{pF}^2}{n'_p n'_F} + \frac{\sigma_{pO}^2}{n'_p n'_O} + \frac{\sigma_{FO}^2}{n'_F n'_O} + \frac{\sigma_{pFO}^2}{n'_p n'_F n'_O} \quad (1)$$

where σ_{Δ}^2 = absolute error variance,

σ_F^2 = variance associated with forms,

σ_O^2 = variance associated with occasions,

σ_{pF}^2 = variance associated with the interaction between persons and forms,

σ_{pO}^2 = variance associated with the interaction between persons and occasions,

σ_{FO}^2 = variance associated with the interaction between forms and occasions,

σ_{pFO}^2 = variance associated with the interaction between persons, forms, and occasions, and

all n 's represent the number of factors contributing to the variance component. The single quotation mark on each n represents a value that can be changed to obtain estimates of the variance with different numbers contributing to the variance estimate – for example, increasing the number of test forms or testing occasions. Each of these variance components is produced

from the G-study and is reported for the observed n 's. The final variance term (person by form by occasion interaction) is generally interpreted as the residual.

The square root of the absolute variances can be interpreted as the “absolute” standard error of measurement (SEM). Absolute variances are generally used to make criterion/domain-referenced decisions (Shavelson & Webb, 2006), or within-student decisions (Hintze, Owen, Shapiro, & Daly, 2000). Relative error variances are used to make normative decisions (i.e., relative to the other persons tested, what is the standard error?). According to Brennan (2001), the square root of the relative error variances can be interpreted essentially identically to the SEM in classical test theory. The relative error variances will nearly always be lower than the absolute variance because only variance components including persons are included. For the two-facet design the relative error variance is defined as

$$\sigma_{\delta}^2 = \frac{\sigma_{pF}^2}{n'_F} + \frac{\sigma_{pO}^2}{n'_O} + \frac{\sigma_{pFO}^2}{n'_F n'_O} \quad (2)$$

where σ_{δ}^2 = relative error variance, and all other terms are defined as above. In this paper, we present both the variances and their corresponding square root, which places the value back onto the scale of the measure. For ease of interpretation, we call the square root of the variances the absolute or relative standard error of the measures. Although the analogy is not direct, the interpretation is similar enough that these terms can be used to facilitate understanding. Just as with classical test theory, the SEMs can be used to construct confidence intervals, as in

$$95\% \text{ CI} = X_{pFO} \pm 1.96(\text{SEM}) \quad (3)$$

where X_{pFO} is the score X for person p on form F on occasion O . One of the added benefits of G-theory is the potential to construct both absolute and relative confidence intervals depending on the decision to be made.

Two types of coefficients are generally produced during the D-study analyses: Generalizability or G-coefficients (Ep^2), which are analogous to coefficient alpha in classical test theory (Brennan, 2001) and phi coefficients (Φ), which are an index of the dependability of the measurement process. Just as with the variance components, these two coefficients correspond to absolute (phi) and relative (g) decisions. The phi index of dependability for absolute decisions is given by

$$\Phi = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_\Delta^2} \quad (4)$$

where all terms are defined as above. In contrast, the g-coefficient for relative decisions is given by

$$Ep^2 = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_\delta^2} \quad (5)$$

where all terms are defined as above. Note that the only difference between equations 4 and 5 is the variance component in the denominator, with the phi-coefficient using the absolute error variance term and the g-coefficient using the relative error variance term.

For each analysis, plots can be produced detailing the change in Ep^2 or Φ with increasing the number of testing occasions and forms administered within each occasion. These are generally displayed as line graphs, with each line representing a different n' of Facet 1 and the x-axis representing a different n' for Facet 2. The plot is simply a visual depiction of the change in reliability coefficients with a corresponding change in the measurement process.

In sum, the G-study provides further information on the sources of error in the measurement process while the D-study provides further information on potential ways that the measurement process could become more dependable. The coefficients to be interpreted depend

upon the use of the measurement tool. If decisions are being made relative to other students (e.g., benchmarking assessments), then the relative error variances and g-coefficients should be interpreted. In contrast, if within-student decisions are being made (e.g., progress-monitoring assessments) then the absolute variances and phi-coefficients should be interpreted.

G-Theory analyses. Data for this study were analyzed in a two-facet fully crossed design (i.e., all students in the analysis were included in both testing occasions and administered the same test forms). The test forms were often administered in a different order on the separate occasions to mitigate order effects. The forms themselves remained constant across occasions in all analyses. We conducted two G-theory analyses for each of the word reading fluency (WRF) and passage reading fluency (PRF) measure types. The first facet in the analysis, *form*, was generally counter-balanced across occasions. The second facet was *occasion*.

For the first PRF analysis, data were collapsed for Teachers 13 and 16 and test forms 14 and 16 were examined in a partially counterbalanced design. The second analysis included students instructed by teachers 14 and 15 and test forms 11, 12, and 13 were examined in a partially counterbalanced design. See Appendix A for the full administration order by teacher.

For all G-theory analyses, forms were analyzed in ascending order regardless of administration order. For example, for the first analysis for PRF, the order of administration for forms 14 and 16 varied by the teacher and occasion. However, during the analysis the data were analyzed for forms 14 and 16 on the first occasion and forms 14 and 16 on the second occasion. In other words, the analysis did not attempt to replicate the administration order because the counterbalanced design was intended to mitigate any order effects. All G-theory analyses were conducted using the SPSS macro produced by Mushquash and O'Connor (2006).

In our results section, we present the results of our G-Studies through an analysis of variance (ANOVA) table detailing the variance associated with each facet of the measurement process as well as all interactions among facets. We then present the error variances and G-coefficients for the design used before presenting the D-Study prophecy estimations results. The D-Study error variance estimates are also presented in their standard error form (i.e., $\sqrt{\sigma^2(\Delta_p)}$ and $\sqrt{\sigma^2(\delta_p)}$ for absolute and relative standard errors respectively), which places the error term back on the scale of the measure and can be used to construct confidence intervals for any individual student's score for any of the measurement designs investigated. Following the error variance estimates, the prophesized G- and Phi-coefficient estimates are presented. Finally a plot was produced for each analysis detailing the estimated change in Ep^2 (labeled on the y-axis as "Mean gstat") with increasing the number of testing occasions and forms administered within each occasion. Each line on the graph represents a different number of testing occasions, ranging from 1-5, while the x-axis represents the number of forms within any occasion. The plot is simply a visual depiction of the G-coefficients table for the corresponding analysis.

Results

The results of the grade 4 Passage Reading Fluency (PRF) measures are presented below. Descriptive statistics are presented in Tables 1 and 2. Test-retest reliability results are presented in Table 3. Correlations between each of the nine forms are presented in Table 4.

Table 1
Descriptive Statistics for Grade 4 Passage Reading Fluency Measures: Session 1

Test Form	<i>n</i>	Min	Max	<i>M</i>	<i>SD</i>
PRF4.8.1	24	27	210	120.58	39.72
PRF4.9.1	26	97	223	150.81	36.07
PRF4.10.1	30	114	249	155.13	32.30
PRF4.11.1	71	41	251	139.54	37.77
PRF4.12.1	50	27	257	135.28	45.17
PRF4.13.1	50	37	218	135.36	36.60
PRF4.14.1	51	97	230	148.67	32.85
PRF4.15.1	21	106	251	159.29	36.90
PRF4.16.1	51	97	254	152.45	32.68

Table 2
Descriptive Statistics for Grade 4 Passage Reading Fluency Measures: Session 2

Test Form	<i>n</i>	Min	Max	<i>M</i>	<i>SD</i>
PRF4.8.2	20	35	217	130.50	44.97
PRF4.9.2	28	111	255	168.89	39.64
PRF4.10.2	29	123	252	168.62	30.23
PRF4.11.2	74	50	255	155.50	39.81
PRF4.12.2	48	32	286	148.38	46.96
PRF4.13.2	74	43	244	145.27	37.63
PRF4.14.2	55	92	255	160.07	36.40
PRF4.15.2	27	85	224	156.07	31.14
PRF4.16.2	55	93	254	162.55	34.98

Test-Retest Reliability

To evaluate test-retest reliability, we correlated performance on each form of the PRF measure that was administered across the two testing sessions. Table 3 presents results of these analyses. Overall, we found a moderately strong test-retest reliability ranging from .86 to .96.

Table 3
Test-retest Reliability Results

Test Form	PRF 4.8.2	PRF 4.9.2	PRF 4.10.2	PRF 4.11.2	PRF 4.12.2	PRF 4.13.2	PRF 4.14.2	PRF 4.15.2	PRF 4.16.2
PRF 4.8.1	0.95								
PRF 4.9.1		0.96							
PRF 4.10.1			0.95						
PRF 4.11.1				0.95					
PRF 4.12.1					0.96				
PRF 4.13.1						0.93			
PRF 4.14.1							0.91		
PRF 4.15.1								-	
PRF 4.16.1									0.86

Alternate Form Reliability

Alternate form reliability was analyzed using bi-variate correlations. We present the correlations between the different forms of each measure in Table 4. We found a moderate to strong positive relationship between the alternate forms, with correlations ranging from .83 to .98.

Table 4
Correlation between Alternate Forms of Grade 4 Passage Reading Fluency Measure

Test Form	PRF4.11.2	PRF4.12.2	PRF4.13.2	PRF4.14.2	PRF4.15.2	PRF4.16.2
PRF4.8.2	0.95	0.97	0.93			
PRF4.9.2	0.98	0.97	0.97			
PRF4.10.2				0.95	0.92	0.94
PRF4.11.2		0.96	0.94	0.95		0.74
PRF4.12.2			0.95			
PRF4.13.2				0.92		0.70
PRF4.14.2					0.89	0.83
PRF4.15.2						0.88

G-study / D-study results. The results of the test-retest and alternate-form reliability analyses suggested acceptable form equivalence for subsequent G-Theory analyses. For the two Passage Reading Fluency analyses, 88% and 80% of the variance was associated with the 39 and 48 persons included in the analysis, 0% was associated with forms, and 0% was associated with occasion. The relative error variance was 30.00 and 66.75 for the first and second analysis, respectively. The absolute variance was 64.07 and 98.35, respectively. The G-Coefficients were .98 for the first analysis and .94 for the second analysis, while the phi coefficients were .96 and .91, respectively.

Passage Reading Fluency: Forms 11, 12, and 13 (teachers 14 and 15)

Grade 4 PRF: Forms 11, 12 & 13

Generalizability ANOVA Table

Facet	<i>df</i>	SS	MS	Variance	Proportion
Persons	38	317624.9	8358.549	1363.097	0.877
Forms	2	744.231	372.115	0	0
Occasions	1	8091.35	8091.35	65.518	0.042
Person*Forms	76	10200.44	134.216	31.151	0.02
Person*Occasion	38	4471.316	117.666	15.251	0.01
Forms*Occasion	2	759.906	379.953	7.898	0.005
Person*Forms*Occasions (Residual)	76	5465.427	71.914	71.914	0.046

Note. Analysis included 39 students, with 3 forms (11, 12 & 13) on 2 occasions.

Error Variances:

Relative, $\sigma^2(\delta_p)$ | Absolute, $\sigma^2(\Delta_p)$
 29.995 64.070

G-coefficients:

G: $E\rho^2$ | Phi: Φ
 .978 .955

Grade 4 PRF: Forms 11, 12 & 13

D-Study Absolute Error Variances, $\sigma^2(\Delta_p)$

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	191.733	111.442	84.678	71.297	63.268
2	136.251	75.913	55.801	45.744	39.711
3	117.757	64.070	46.175	37.227	31.858
4	108.510	58.149	41.362	32.968	27.932
5	102.962	54.596	38.474	30.413	25.577

Grade 4 PRF: Forms 11, 12 & 13

D-Study Absolute Standard Errors, $\sigma(\Delta_p)$

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	13.847	10.557	9.202	8.444	7.954
2	11.673	8.713	7.470	6.763	6.302
3	10.852	8.004	6.795	6.101	5.644
4	10.417	7.626	6.431	5.742	5.285
5	10.147	7.389	6.203	5.515	5.057

Grade 4 PRF: Forms 11, 12 & 13

D-Study Relative Error Variances, $\sigma^2(\delta_p)$

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	118.316	74.734	60.206	52.942	48.584
2	66.783	41.180	32.645	28.378	25.817
3	49.606	29.995	23.458	20.189	18.228
4	41.017	24.402	18.864	16.095	14.434
5	35.864	21.047	16.108	13.639	12.157

Grade 4 PRF: Forms 11, 12 & 13

D-Study Relative Standard Errors, $\sigma(\delta_p)$

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	10.877	8.645	7.759	7.276	6.970
2	8.172	6.417	5.714	5.327	5.081
3	7.043	5.477	4.843	4.493	4.269
4	6.404	4.940	4.343	4.012	3.799
5	5.989	4.588	4.013	3.693	3.487

Grade 4 PRF: Forms 11, 12 & 13

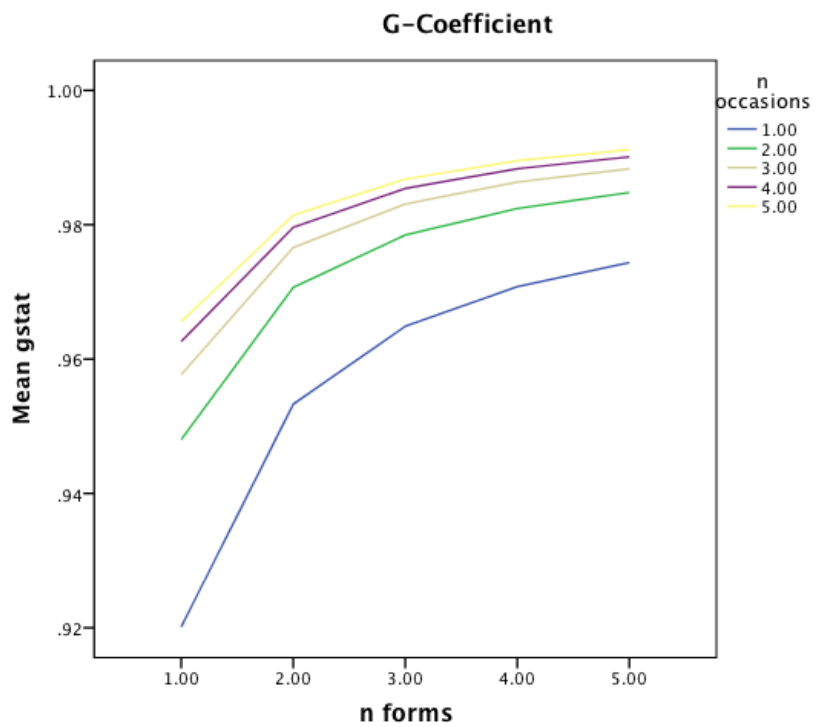
D-Study G Coefficients, Ep^2

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.920	0.948	0.958	0.963	0.966
2	0.953	0.971	0.977	0.980	0.981
3	0.965	0.978	0.983	0.985	0.987
4	0.971	0.982	0.986	0.988	0.990
5	0.974	0.985	0.988	0.990	0.991

Grade 4 PRF: Forms 11, 12 & 13

D-Study Phi Coefficients, Φ

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.877	0.924	0.942	0.950	0.956
2	0.909	0.947	0.961	0.968	0.972
3	0.920	0.955	0.967	0.973	0.977
4	0.926	0.959	0.971	0.976	0.980
5	0.930	0.961	0.973	0.978	0.982



 Passage Reading Fluency: Forms 14 and 16 (teachers 13 & 16)

Grade 4 PRF: Forms 14 & 16

Generalizability ANOVA Table

Facet	<i>df</i>	SS	MS	Variance	Proportion
Persons	47	209733.2	4462.408	1048.85	0.801
Forms	1	602.083	602.083	4.843	0.004
Occasions	1	5676.75	5676.75	58.353	0.045
Person*Forms	47	10757.42	228.881	50.22	0.038
Person*Occasion	47	7828.75	166.569	19.064	0.015
Forms*Occasion	1	36.75	36.75	0	0
Person*Forms*Occasions (Residual)	47	6036.75	128.441	128.441	0.098

Note. Analysis included 48 students, with 2 forms (14 & 16) on 2 occasions.

Error Variances:

Relative, $\sigma^2(\delta_p)$	Absolute, $\sigma^2(\Delta_p)$
66.752	98.350

G-coefficients:

G: E_p^2	Phi: Φ
.940	.914

Grade 4 PRF: Forms 14 & 16

D-Study: Absolute Error Variances, $\sigma^2(\Delta_p)$

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	260.921	157.992	123.682	106.527	96.234
2	169.169	98.350	74.744	62.941	55.859
3	138.585	78.469	58.431	48.412	42.400
4	123.293	68.529	50.275	41.147	35.671
5	114.117	62.565	45.381	36.789	31.633

Grade 4 PRF: Forms 14 & 16

D-Study: Absolute Standard Errors, $\sigma(\Delta_p)$

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	16.153	12.569	11.121	10.321	9.810
2	13.006	9.917	8.645	7.934	7.474
3	11.772	8.858	7.644	6.958	6.512
4	11.104	8.278	7.090	6.415	5.973
5	10.683	7.910	6.737	6.065	5.624

Grade 4 PRF: Forms 14 & 16

D-Study Relative Error Variances, $\sigma^2(\delta_p)$

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	197.725	123.973	99.388	87.096	79.721
2	108.395	66.752	52.871	45.931	41.767
3	78.618	47.679	37.366	32.209	29.115
4	63.729	38.142	29.613	25.349	22.790
5	54.796	32.420	24.961	21.232	18.994

Grade 4 PRF: Forms 14 & 16

D-Study Relative Standard Errors, $\sigma(\delta_p)$

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	14.061	11.134	9.969	9.333	8.929
2	10.411	8.170	7.271	6.777	6.463
3	8.867	6.905	6.113	5.675	5.396
4	7.983	6.176	5.442	5.035	4.774
5	7.402	5.694	4.996	4.608	4.358

Grade 4 PRF: Forms 14 & 16

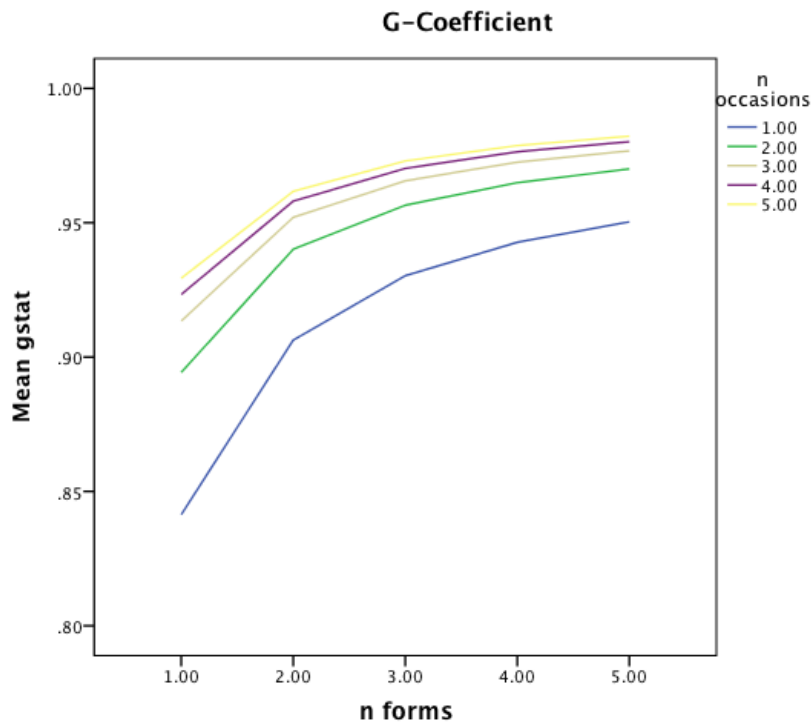
D-Study G Coefficients, Ep^2

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.841	0.894	0.913	0.923	0.929
2	0.906	0.940	0.952	0.958	0.962
3	0.930	0.957	0.966	0.970	0.973
4	0.943	0.965	0.973	0.976	0.979
5	0.950	0.970	0.977	0.980	0.982

Grade 4 PRF: Forms 14 & 16

D-Study Phi Coefficients, Φ

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.801	0.869	0.895	0.908	0.916
2	0.861	0.914	0.933	0.943	0.949
3	0.883	0.930	0.947	0.956	0.961
4	0.895	0.939	0.954	0.962	0.967
5	0.902	0.944	0.959	0.966	0.971



Discussion

The test-retest and alternate form reliability results of this study provide moderate to high evidence of the reliability of the easyCBM grade 4 PRF measures, with high test-retest reliability and moderate to high correlations between the alternate forms of the PRF measures.

The results of the G- and D-studies also increased the overall reliability evidence for the easyCBM reading measures. For the G-studies, the majority of variance was attributed to persons in every analysis, with 88% and 80% respectively. The standard errors were also quite low. It is important to note that the error variances and dependability coefficients reported in text in the results section are those of the corresponding *analysis* and not of a particular form. For example, an examination of the error variance or standard error tables will show a bolded number, which is the error for the analysis. However, if only one form were given on one occasion then the error is increased (as reported in the D-study tables). Thus, in a classroom where decisions are made

from one test form after one testing occasion, the error more closely resembles the one form on one occasion numbers reported in the D-study standard error tables.

Generally, increasing either facet (occasions or forms) resulted in a similar increase in the overall dependability. When examining the overall results, however, it is evident that using a single test form on a single occasion is generally sufficient for dependable measurement (i.e., $> .8$). This finding is important because other measurement systems have recommended using 3 fluency forms and taking the median score to increase reliability (Dibels*Next*, 2011) – a procedure that may appear unnecessary given the results of this study.

References

- Alonzo, J., Tindal, G., & Ketterlin-Geller, L.R. (2006). General outcome measures of basic skills in reading and math. In L. Florian (Ed.), *Handbook of Special Education*. Thousand Oaks, CA: Sage.
- Brennan, R. L. (2001). *Statistics for social science and public policy: Generalizability theory*. New York: Springer.
- Deno, S. L. (2003). Developments in curriculum-based measurements. *The Journal of Special Education, 37*, 184-192.
- Deno, S. (1987). Curriculum-based measurement. *Teaching Exceptional Children*. (Fall), 41-47.
- Deno, S. L., & Mirkin, P. M. (1977). *Data based program modification*. Minneapolis, MN: University of Minnesota Leadership Training Institute/Special Education.
- DibelsNext (2011). *Dibels Oral Reading Fluency*. Retrieved February 14, 2011, from https://www.mclasshome.com/wgenhelp/dnext/DIBELS_Next/Assessment_and_Scoring/DO_RF_Details.htm
- Good, R. H., Gruba, J., & Kaminski, R. A. (2002). Best practices in Using Dynamic Indicators of Basic Early Literacy Skills (DIBELS) in an Outcomes-Driven Model. In A. Thomas and J. Grimes (Eds.). *Best Practices in School Psychology IV* (pp.679-700). Washington, DC: National Association of School Psychologists.
- Hintze, J. M., Owen, S. V., Shapiro, E. S., and Daly, E. J. (2000). Research design and methodology section: Generalizability of oral reading fluency measures: Application of G theory to curriculum-based measurement. *School Psychology Quarterly, 15*, 52-68.
- Mushquash, C., & O'connor, B. P. (2006). SPSS and SAS programs for generalizability theory analyses. *Behavior Research Methods, 38*, 542-547.

Shavelson, R. J., & Webb, N. M. (2006). Generalizability theory. In Green, J. L., Camilli, G. & Elmore, P. B. (Eds.), *Complementary Methods for Research in Education*, (pp. 309-322). (3rd ed.) Washington, DC: AERA.

 Appendix A

 Full test form administration order

Teacher	Passage Reading Fluency	
	Occasion 1	Occasion 2
13	14 – 16 – 10	15 – 16 – 14 – 10
14	13 – 12 – 11 – 9	13 – 11 – 12 – 9
15	11 – 12 – 13 – 8	12 – 13 – 11 – 8
16	16 – 15 – 14 – 11	16 – 14 – 13 – 11

 Test forms used in G-Theory analyses

Teacher	Passage Reading Fluency	
	Occasion 1	Occasion 2
13	14 – 16	16 – 14
14	13 – 12 – 11	13 – 11 – 12
15	11 – 12 – 13	12 – 13 – 11
16	16 – 14	16 – 14
