

Technical Report # 1220

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

Reading Assessments:

Grade 5

Cheng-Fei Lai

Bitnara Jasmine Park

Daniel Anderson

Julie Alonzo

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 Reading Assessments: Grade 5

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 4 different G-theory analyses for passage reading fluency (PRF) to investigate 8 different test forms. The first facet in the analysis, *form*, was generally counterbalanced across occasions. The second facet was *occasion*.

For the first PRF analysis, data were analyzed for teacher 17 and test forms 10, 14, and 15, were examined in a partially counterbalanced design. The second analysis explored data from teacher 19 and examined forms 8, 11, 12, and 13 in a partially counterbalanced design. The third analysis came from teacher 20 and examined forms 8, 9, 10, and 12 in a partially counterbalanced design. Finally, the fourth analysis explored data from teacher 21 and examined test forms 9 and 13 in a fully 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 10, 14, and 15 varied by the teacher and occasion. However, during the analysis the data were analyzed for forms 10, 14, and 15 on the first occasion and forms 10, 14, and 15 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 5 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 eight forms are presented in Table 4.

Table 1
Descriptive Statistics for Grade 5 Passage Reading Fluency Measures: Session 1

Test Form	<i>n</i>	Min	Max	<i>M</i>	<i>SD</i>
PRF5.8.1	40	81	243	159.85	38.96
PRF5.9.1	69	102	296	180.57	42.08
PRF5.10.1	59	110	273	185.64	43.59
PRF5.11.1	50	84	289	172.56	43.27
PRF5.12.1	70	76	273	171.89	40.30
PRF5.13.1	69	73	301	173.38	47.45
PRF5.14.1	39	115	261	187.26	42.54
PRF5.15.1	20	156	250	205.80	28.28

Table 2
Descriptive Statistics for Grade 5 Passage Reading Fluency Measures: Session 2

Test Form	<i>n</i>	Min	Max	<i>M</i>	<i>SD</i>
PRF5.8.2	46	87	308	177.17	45.36
PRF5.9.2	68	86	342	199.22	49.37
PRF5.10.2	87	69	338	198.92	49.06
PRF5.11.2	44	100	273	184.55	39.93
PRF5.12.2	46	94	354	184.35	49.06
PRF5.13.2	85	86	320	182.08	42.61
PRF5.14.2	85	92	294	193.24	38.70
PRF5.15.2	44	83	282	192.34	36.93
PRF5.16.2	19	174	295	221.21	35.18

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 moderate to strong test-retest reliability, with all but one form ranging from .88 to .94. One form (PRF5.14) had moderate test-retest reliability ($R = .54$).

Table 3
Test-retest Reliability Results

Test Form	PRF5.8. 2	PRF5.9. 2	PRF5.1 0.2	PRF5.1 1.2	PRF5.1 2.2	PRF5.1 3.2	PRF5.1 4.2	PRF5.1 5.2
PRF5.8.1	0.94							
PRF5.9.1		0.91						
PRF5.10.1			0.88					
PRF5.11.1				0.93				
PRF5.12.1					0.90			
PRF5.13.1						0.90		
PRF5.14.1							0.54	
PRF5.15.1								0.91

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 .85 to .98, with one exception (the correlation between forms 5.10 and 5.15 was moderate, at $R = .73$).

Table 4
Correlation between Alternate Forms of Grade 5 Passage Reading Fluency Measure

Test Form	PRF5.9.2	PRF5.10. 2	PRF5.11. 2	PRF5.12. 2	PRF5.13. 2	PRF5.14. 2	PRF5.15. 2
PRF5.8.2	0.94	0.95	0.98	0.94	0.97		
PRF5.9.2		0.92		0.93	0.86	0.94	
PRF5.10. 2				0.97	0.91	0.85	0.73
PRF5.11. 2				0.95	0.91	0.88	0.92
PRF5.12. 2					0.97		
PRF5.13. 2						0.89	0.92
PRF5.14. 2							0.92

G-study / D-study results. For the four Passage Reading Fluency analyses, 57%, 89%, 79%, and 86% of the variance was associated with the 17, 18, 13, and 11 persons included in the

analysis, 0% was associated with forms, and 0% was associated with occasion. The relative error variance was 71.19, 38.41, 18.30, and 50.43 for the first, second, third, and fourth analysis, respectively. The absolute variance was 135.35, 58.53, 34.47, and 109.83, respectively. The G-Coefficients were .90 for the first analysis, .98 for the second, .97 for the third and the fourth, while the phi coefficients were .83, .96, .95, and .94 for the first, second, third, and fourth analysis, respectively.

Passage Reading Fluency: Forms 10, 14, & 15 (teacher 17)

Grade 5 PRF: Forms 10, 14 & 15
Generalizability ANOVA Table

Facet	<i>df</i>	SS	MS	Variance	Proportion
Persons	16	70750.65	4421.915	665.797	0.574
Forms	2	3651.941	1825.971	26.705	0.023
Occasions	1	6010.676	6010.676	98.891	0.085
Person*Forms	32	10405.06	325.158	26.363	0.023
Person*Occasion	16	5990.49	374.406	33.991	0.029
Forms*Occasion	2	1730.529	865.265	34.873	0.03
Person*Forms*Occasions (Residual)	32	8717.804	272.431	272.431	0.235

Note. Analysis included 17 students, with 3 forms (10, 14 & 15) on 2 occasions.

Error Variances:

Relative, $\sigma^2(\delta_p)$	Absolute, $\sigma^2(\Delta_p)$
71.189	135.348

G-coefficients:

G: E_p^2	Phi: Φ
.903	.831

Grade 5 PRF: Forms 10, 14 & 15

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	493.255	273.162	199.797	163.115	141.106
2	313.069	169.801	122.046	98.168	83.841
3	253.007	135.348	96.129	76.519	64.753
4	222.975	118.121	83.17	65.694	55.209
5	204.957	107.785	75.395	59.200	49.482

Grade 5 PRF: Forms 10, 14 & 15

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	22.209	16.528	14.135	12.772	11.879
2	17.694	13.031	11.047	9.908	9.156
3	15.906	11.634	9.805	8.748	8.047
4	14.932	10.868	9.120	8.105	7.430
5	14.316	10.382	8.683	7.694	7.034

Grade 5 PRF: Forms 10, 14 & 15

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	332.786	179.575	128.504	102.969	87.648
2	183.389	98.285	69.917	55.733	47.223
3	133.590	71.189	50.388	39.988	33.748
4	108.690	57.640	40.624	32.116	27.011
5	93.750	49.512	34.765	27.392	22.968

Grade 5 PRF: Forms 10, 14 & 15

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	18.242	13.401	11.336	10.147	9.362
2	13.542	9.914	8.362	7.465	6.872
3	11.558	8.437	7.098	6.324	5.809
4	10.425	7.592	6.374	5.667	5.197
5	9.682	7.036	5.896	5.234	4.792

Grade 5 PRF: Forms 10, 14 & 15

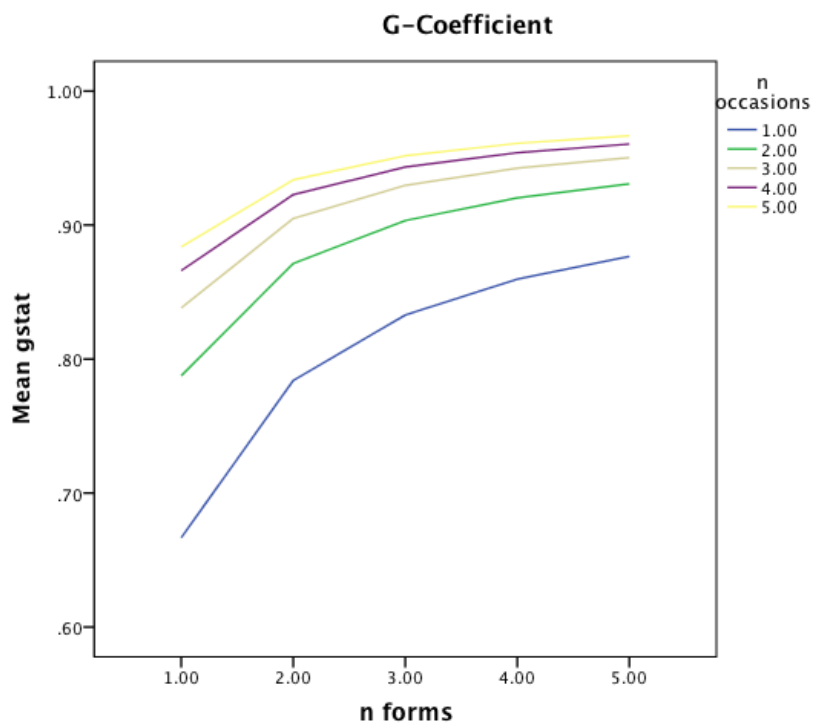
D-Study G Coefficients, Ep^2

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.667	0.788	0.838	0.866	0.884
2	0.784	0.871	0.905	0.923	0.934
3	0.833	0.903	0.930	0.943	0.952
4	0.860	0.920	0.942	0.954	0.961
5	0.877	0.931	0.950	0.960	0.967

Grade 5 PRF: Forms 10, 14 & 15

D-Study Phi Coefficients, Φ

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.574	0.709	0.769	0.803	0.825
2	0.680	0.797	0.845	0.872	0.888
3	0.725	0.831	0.874	0.897	0.911
4	0.749	0.849	0.889	0.910	0.923
5	0.765	0.861	0.898	0.918	0.931



 Passage Reading Fluency: Forms 8, 11, 12, and 13 (teacher 19)

Grade 5 PRF: Forms 8, 11, 12 & 13

Generalizability ANOVA Table

Facet	<i>df</i>	SS	MS	Variance	Proportion
Persons	17	214730	12631.18	1540.484	0.886
Forms	3	1184.972	394.991	9.479	0.005
Occasions	1	2738.778	2738.778	35.497	0.02
Person*Forms	51	6739.028	132.138	22.952	0.013
Person*Occasion	17	4443.722	261.395	43.79	0.025
Forms*Occasion	3	23.556	7.852	0	0
Person*Forms*Occasions (Residual)	51	4397.944	86.234	86.234	0.05

Note. Analysis included 18 students, with 4 forms (8, 11, 12 & 13) on 2 occasions.

Error Variances:

Relative, $\sigma^2(\delta_p)$		Absolute, $\sigma^2(\Delta_p)$
38.412		58.530

G-coefficients:

G: E_p^2		Phi: Φ
.976		.963

Grade 5 PRF: Forms 8, 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	197.952	115.191	87.604	73.811	65.535
2	138.619	77.417	57.017	46.816	40.696
3	118.842	64.826	46.821	37.818	32.417
4	108.953	58.530	41.723	33.319	28.277
5	103.020	54.753	38.664	30.620	25.793

Grade 5 PRF: Forms 8, 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	14.070	10.733	9.360	8.591	8.095
2	11.774	8.799	7.551	6.842	6.379
3	10.901	8.051	6.843	6.150	5.694
4	10.438	7.650	6.459	5.772	5.318
5	10.150	7.400	6.218	5.534	5.079

Grade 5 PRF: Forms 8, 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	152.976	87.964	66.293	55.458	48.957
2	98.383	54.930	40.445	33.203	28.857
3	80.186	43.918	31.829	25.784	22.158
4	71.087	38.412	27.521	22.075	18.808
5	65.628	35.109	24.936	19.850	16.798

Grade 5 PRF: Forms 8, 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	12.368	9.379	8.142	7.447	6.997
2	9.919	7.411	6.360	5.762	5.372
3	8.955	6.627	5.642	5.078	4.707
4	8.431	6.198	5.246	4.698	4.337
5	8.101	5.925	4.994	4.455	4.099

Grade 5 PRF: Forms 8, 11, 12 & 13

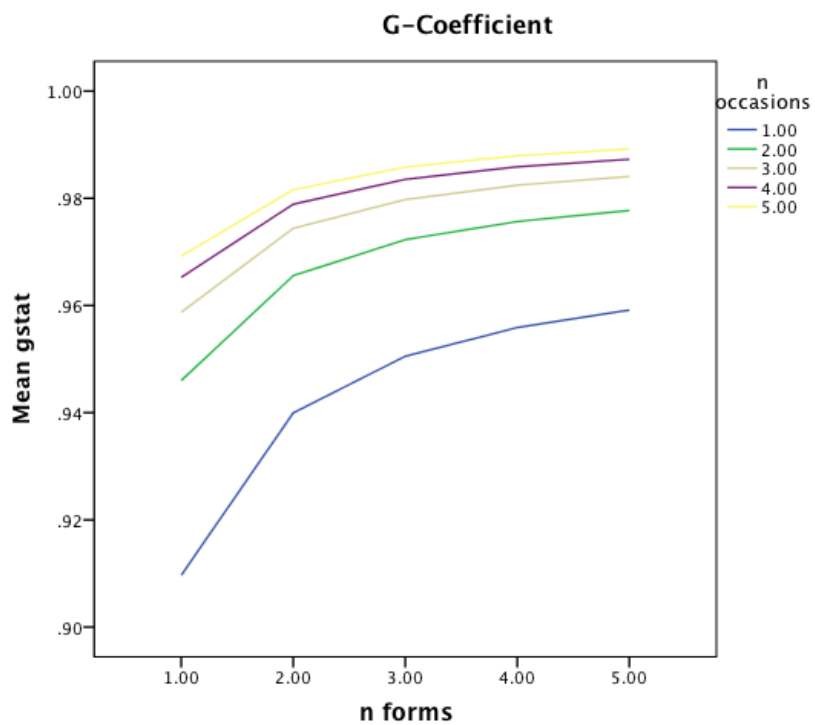
D-Study G Coefficients, Ep^2

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.910	0.946	0.959	0.965	0.969
2	0.940	0.966	0.974	0.979	0.982
3	0.951	0.972	0.980	0.984	0.986
4	0.956	0.976	0.982	0.986	0.988
5	0.959	0.978	0.984	0.987	0.989

Grade 5 PRF: Forms 8, 11, 12 & 13

D-Study Phi Coefficients, Φ

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.886	0.930	0.946	0.954	0.959
2	0.917	0.952	0.964	0.971	0.974
3	0.928	0.960	0.971	0.976	0.979
4	0.934	0.963	0.974	0.979	0.982
5	0.937	0.966	0.976	0.981	0.984



 Passage Reading Fluency: Forms 8, 9, 10, & 12 (teacher 20)

Grade 5 PRF: Forms 8, 9, 10 & 12

Generalizability ANOVA Table

Facet	<i>df</i>	SS	MS	Variance	Proportion
Persons	12	61363.9	5113.659	621.754	0.792
Forms	3	2289	763	26.736	0.034
Occasions	1	1098.5	1098.5	18.962	0.024
Person*Forms	36	3664.25	101.785	0	0
Person*Occasion	12	1756.75	146.396	9.461	0.012
Forms*Occasion	3	223.885	74.628	0	0
Person*Forms*Occasions (Residual)	36	3907.865	108.552	108.552	0.138

Note. Analysis included 13 students, with 4 forms (8, 9, 10 & 12) on 2 occasions.

Error Variances:

Relative, $\sigma^2(\delta_p)$		Absolute, $\sigma^2(\Delta_p)$
18.299		34.465

G-coefficients:

G: E_p^2		Phi: Φ
.971		.947

Grade 5 PRF: Forms 8, 9, 10 & 12

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	163.711	95.224	72.394	60.980	54.131
2	96.067	54.718	40.934	34.043	29.908
3	73.519	41.216	30.448	25.064	21.833
4	62.245	34.465	25.204	20.574	17.796
5	55.481	30.414	22.058	17.881	15.374

Grade 5 PRF: Forms 8, 9, 10 & 12

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	12.795	9.758	8.508	7.809	7.357
2	9.801	7.397	6.398	5.835	5.469
3	8.574	6.420	5.518	5.006	4.673
4	7.890	5.871	5.020	4.536	4.219
5	7.449	5.515	4.697	4.229	3.921

Grade 5 PRF: Forms 8, 9, 10 & 12

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	118.013	59.006	39.338	29.503	23.603
2	63.737	31.868	21.246	15.934	12.747
3	45.645	22.822	15.215	11.411	9.129
4	36.599	18.299	12.200	9.150	7.320
5	31.171	15.586	10.390	7.793	6.234

Grade 5 PRF: Forms 8, 9, 10 & 12

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	10.863	7.682	6.272	5.432	4.858
2	7.984	5.645	4.609	3.992	3.570
3	6.756	4.777	3.901	3.378	3.021
4	6.050	4.278	3.493	3.025	2.706
5	5.583	3.948	3.223	2.792	2.497

Grade 5 PRF: Forms 8, 9, 10 & 12

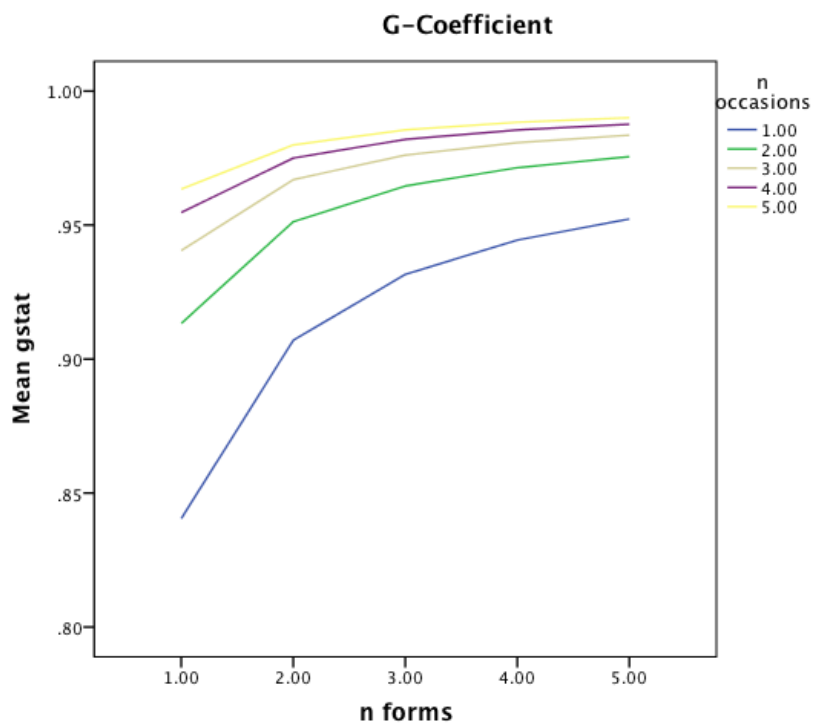
D-Study G Coefficients, E_p^2

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.840	0.913	0.940	0.955	0.963
2	0.907	0.951	0.967	0.975	0.980
3	0.932	0.965	0.976	0.982	0.986
4	0.944	0.971	0.981	0.985	0.988
5	0.952	0.976	0.984	0.988	0.990

Grade 5 PRF: Forms 8, 9, 10 & 12

D-Study Phi Coefficients, Φ

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.792	0.867	0.896	0.911	0.920
2	0.866	0.919	0.938	0.948	0.954
3	0.894	0.938	0.953	0.961	0.966
4	0.909	0.947	0.961	0.968	0.972
5	0.918	0.953	0.966	0.972	0.976



 Passage Reading Fluency: Forms 9 & 13 (teacher 21)

Grade 5 PRF: Forms 9 & 13

Generalizability ANOVA Table

Facet	<i>df</i>	SS	MS	Variance	Proportion
Persons	10	75655.14	7565.514	1840.95	0.855
Forms	1	1937.818	1937.818	63.464	0.029
Occasions	1	1375.364	1375.364	39.732	0.018
Person*Forms	10	1982.682	198.268	21.9	0.01
Person*Occasion	10	1579.136	157.914	1.723	0.001
Forms*Occasion	1	497.818	497.818	31.214	0.014
Person*Forms*Occasions (Residual)	10	1544.682	154.468	154.468	0.072

Note. Analysis included 11 students, with 2 forms (9 & 13) on 2 occasions.

Error Variances:

Relative, $\sigma^2(\delta_p)$		Absolute, $\sigma^2(\Delta_p)$
50.428		109.830

G-coefficients:

G: $E\rho^2$		Phi: Φ
.973		.944

Grade 5 PRF: Forms 9 & 13

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	312.500	198.932	161.076	142.148	130.791
2	176.977	109.830	87.447	76.256	69.541
3	131.803	80.129	62.904	54.292	49.124
4	109.216	65.278	50.633	43.310	38.916
5	95.664	56.368	43.270	36.720	32.791

Grade 5 PRF: Forms 9 & 13

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	17.678	14.104	12.692	11.923	11.436
2	13.303	10.480	9.351	8.732	8.339
3	11.481	8.951	7.931	7.368	7.009
4	10.451	8.079	7.116	6.581	6.238
5	9.781	7.508	6.578	6.060	5.726

Grade 5 PRF: Forms 9 & 13

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

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	178.091	99.995	73.964	60.948	53.138
2	89.907	50.428	37.269	30.689	26.741
3	60.512	33.906	25.037	20.603	17.942
4	45.815	25.645	18.922	15.560	13.543
5	36.996	20.688	15.252	12.534	10.903

Grade 5 PRF: Forms 9 & 13

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

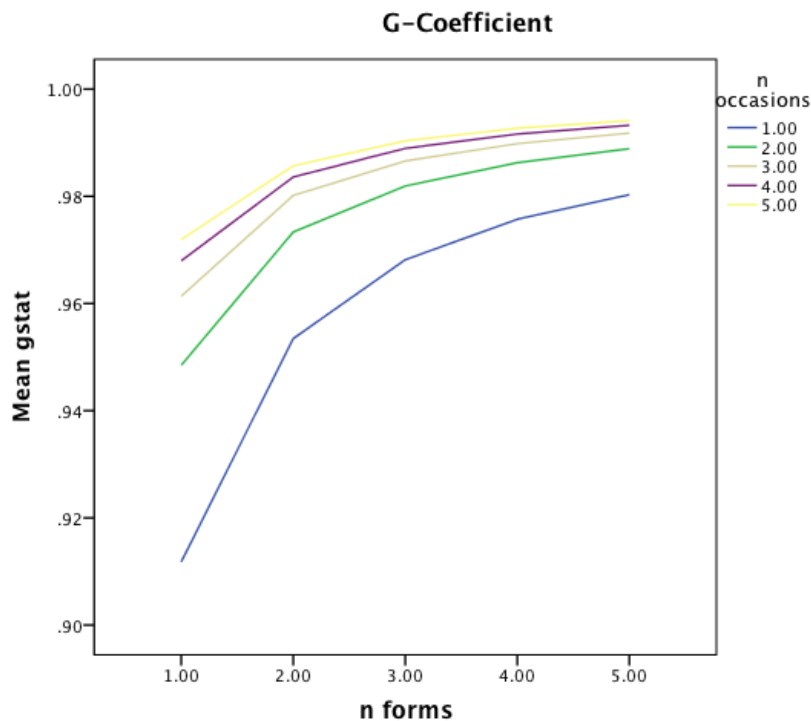
<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	13.345	10.000	8.600	7.807	7.290
2	9.482	7.101	6.105	5.540	5.171
3	7.779	5.823	5.004	4.539	4.236
4	6.769	5.064	4.350	3.945	3.680
5	6.082	4.548	3.905	3.540	3.302

Grade 5 PRF: Forms 9 & 13
D-Study G Coefficients, $E\rho^2$

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.912	0.948	0.961	0.968	0.972
2	0.953	0.973	0.980	0.984	0.986
3	0.968	0.982	0.987	0.989	0.990
4	0.976	0.986	0.990	0.992	0.993
5	0.980	0.989	0.992	0.993	0.994

Grade 5 PRF: Forms 9 & 13
D-Study Phi Coefficients, Φ

<i>n</i> forms	<i>n</i> occasions				
	1	2	3	4	5
1	0.855	0.902	0.920	0.928	0.934
2	0.912	0.944	0.955	0.960	0.964
3	0.933	0.958	0.967	0.971	0.974
4	0.944	0.966	0.973	0.977	0.979
5	0.951	0.970	0.977	0.980	0.982



Discussion

The test-retest and alternate form reliability results of this study provide moderate to strong evidence of the reliability of the easyCBM grade 5 PRF measures, with moderate to strong test-retest reliability and moderate to strong correlations between the alternate forms of the passage reading fluency measures.

The results of the G- and D-Theory analyses were generally mixed, with the first analysis (Teacher 17) displaying the poorest results and the second analysis (Teacher 19) displaying the best results. Overall, 57% - 89% of the total variance was associated with persons during the G-Study, while the predicted reliability for relative decisions for one form on one occasion ranged from .67 to .91. The standard errors were generally 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. Using .8 as the cutoff for acceptable reliability of relative decisions (Ep^2) the results generally suggest that one form on one occasion would be significant (with the exception of the analysis for Teacher 17). This finding is important because other measurement systems have recommended using 3 fluency forms and taking the median score to increase reliability (DibelsNext, 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
17	14 – 15 – 10	15 – 16 – 14 – 10
19	11 – 12 – 13 – 8	12 – 13 – 11 – 8
20	8 – 9 – 10 – 12	9 – 10 – 8 – 12
21	13 – 12 – 11 – 9	9 – 10 – 13 – 14

 Full test form administration order

Teacher	Passage Reading Fluency	
	Occasion 1	Occasion 2
17	14 – 15 – 10	15 – 14 – 10
19	11 – 12 – 13 – 8	12 – 13 – 11 – 8
20	8 – 9 – 10 – 12	9 – 10 – 8 – 12
21	13 – 9	9 – 13
