

Validating the ASER Testing Tools: Comparisons with Reading Fluency Measures and the Read India Measures ¹

Shaher Banu Vagh

This paper presents evidence for the psychometric properties of the ASER testing tools. It focuses on evaluating their validity and providing a detailed comparison of the ASER reading test with fluency measures that are an adaptation of the Early Grade Reading Assessment (EGRA).

Introduction to ASER

The Annual Status of Education Report (ASER), a nationwide survey of reading and math achievement of children from rural India, has been conducted annually since 2005. ASER provides basic and critical information about rural Indian children's foundational reading skills and basic math ability. Given its scale and comprehensive coverage, it is a path breaking initiative as it is the only nationwide survey, albeit rural, which assesses the learning achievement of children in the age group 5-16. The survey is conducted each year in the middle of the academic year (October to November) and the findings are made public, for most states, in the same school year (mid-January). The availability of results in the same school year is a tremendous feat for such a large survey, which enhances its potential as a tool to inform educational practice and policy.

The ASER test inference is about a child's level of foundational reading skills (letter identification, word decoding, etc.) and basic math ability (number recognition, subtraction, and division). The content of the ASER-reading test, i.e. the selection of words, length of sentences and paragraphs, and use of vocabulary is aligned to Grade 1 and Grade 2 level state textbooks and the ASER-math test is aligned to Grade 1, 2, 3, and 4 level state textbooks. The tests are orally and individually administered and require about 10 minutes of administration time. They are designed as criterion-referenced tests that categorize children on an ordinal scale indexing mastery in the basic skills of reading and number operations.

The tests are designed to understand what students can do and the skills they have mastered. For instance the ASER-reading test classifies children at the 'nothing', 'letter', 'word', 'paragraph' (grade 1 level text), and 'story' (grade 2 level text) level based on defined performance criteria or cut-off scores that allow examiners to classify children as masters or non-masters of any given level. For example, the inability to correctly identify 4 out of 5 letters classifies the child at the 'nothing' level. The ASER math test classifies children at the 'nothing', 'single digit recognition', 'double digit recognition', 'subtraction with carry over', and 'division' level (see www.asercentre.org for testing tools and the annual reports for test administration details).

The ASER testing tools have several advantages: they are simple, quick, cost-effective, and easy to train examiners to administer. All of these are desirable features (Wagner, 2003) as it makes it feasible to conduct a survey of the scale and scope of the ASER (assessing about 700,000 children every year) and make results available in a timely manner, which has the potential to inform educational practice and

¹ The complete technical report is available on request from ASER Centre; email: shahervagh@gmail.com

policy. However, several pertinent questions have been raised about the ASER testing tools in relation to their content and their statistical properties. Specifically, how robust are the ASER-reading and ASER-math testing tools? Do the ASER testing tools provide reliable and valid findings? The present report, therefore, aims to address these critical questions about the Hindi language and the math testing tools.

Defining Reliability and Validity

The traditional notion of reliability is the *consistency* with which a test measures any given skill and thereby enables us to *consistently* distinguish between individuals with regards to the ability or skill being measured. In other words, if it were possible and feasible to test children repeatedly using the same test, a reliable test would yield a *consistent* score across the repeated measurements. However, given that the ASER tests assess achievement of mastery rather than the relative standing of children in relation to their peers, reliability in this case is “the consistency of the decision-making process across repeated administrations of the test” (Swaminathan, Hambleton, & Algina, 1974). Hence, reliability in this case does not refer to ‘test reliability’ but rather to the ‘reliability of decisions’ (Huynh, 1976 as cited in Traub & Rowley, 1980) or ‘decision consistency’ (Swaminathan, Hambleton, & Algina, 1974) as the assessment here is about mastery or non-mastery of a level of reading or math.

Validity, on the other hand, indicates whether the test measures what it purports to measure, i.e. how well children’s performance on a test support the conclusions we make about a specific ability or skill. For instance is the inference based on the ASER-reading test about children’s *mastery or non-mastery of basic reading ability* valid? Is the inference based on the ASER-math test about children’s *mastery or non-mastery of basic math ability* valid? Specifically, validity is an evaluation of a test inference and not of the test per se. A test can be put to different uses such as, examining average school performance or making diagnostic decisions about individual students. Each of these uses or inferences “has its own degree of validity, [and] one can never reach the simple conclusion that a particular test “is valid”” (Cronbach, 1971, p.447). Another way to think about reliability and validity is that when playing darts, consistently hitting the same spot, irrespective of position on the target board is akin to reliability and consistently hitting bull’s eye or any other target of interest, is akin to validity. Reliability then is a necessary but not sufficient condition for validity.

Reliability of the ASER Testing Tools

Given that the ASER tests are criterion-referenced tests the evaluation of reliability is evaluated in two ways: (a) as a measure of agreement between decisions made in repeated test administrations, and (b) as a measure of agreement between raters in assigning a mastery level i.e. inter-rater reliability.

A simple method of estimating agreement across repeated test administrations and inter-rater reliability is to examine the association between the ratings of the two administrations or the two examiners for the same group of children by estimating a correlation coefficient. This method, however, tends to overestimate reliability as it merely evaluates *association* and not *agreement*. Specifically, it does not provide information about the agreement in the category to which children were assigned by repeated administrations or by the two independent examiners and it does not take into account agreement due to mere chance, thus correlations provide a limited picture. Instead, estimating a Cohen’s kappa coefficient provides a more accurate estimate of reliability as it estimates agreement between repeated

test administrations or between raters beyond agreement due to chance². The Cohen's kappa estimate varies from 0 to 1. A value of 0 indicates no improvement over chance and a value of 1 indicates maximal increase or perfect agreement³. There are two methods of estimating Kappa: the *simple* Kappa evaluates for exact agreement, whereas the *weighted* Kappa assigns a higher weight to close misses (e.g. a rank of 2 vs. a rank of 3) than to misses that are further apart (e.g. a rank of 1 vs. a rank of 5) based on the assumption that a difference of adjacent ranks is less critical than a difference that is farther apart⁴. In other words, categorizing a 'story' level child to be at the 'para' level is less incorrect than categorizing a 'story' level child to be at the 'nothing' level.

The reliability of decisions across repeated test administrations was examined for a group of 540 children who were assessed by the same team of examiners on two testing occasions. The test-retest correlation coefficients for the ASER-reading test for all children from Grades 1-5 is .95 and for the ASER-math test is .90. More importantly the average Cohen's kappa estimate for decision consistency across repeated test administrations for the ASER-reading test is .76 and for the ASER-math test is .71. These estimates suggest 'substantial' level of agreement⁵.

The inter-rater reliability estimated using Cohen's Kappa for a group of 590 children is .64 for the ASER-reading test and .65 for the ASER-math test on average, also indicating 'substantial' agreement. The average and median weighted Kappa across all pairs of examiners is .82 and .81 respectively for the ASER-reading test and is .79 and .80 for the ASER-math test indicating 'almost perfect' agreement for the ASER-reading test and 'substantial' agreement for the ASER-math test.

Validating the ASER Testing Tools

As part of an evaluation of Pratham's Read India program that was carried out by the Abdul Latif Jameel Poverty Action Lab (JPAL), the ASER-reading and ASER-math tests were administered along with a battery of tests of basic and advanced reading and math ability. Several other tests included : (a) the Fluency Battery, a test of early reading ability, which was adapted from the Early Grade Reading Assessment (USAID, 2009) and the Dynamic Indicators of Basic Early Literacy Skills (University of Oregon Center on Teaching and Learning, 2002); (b) the Read India (RI) Literacy test, which is a paper-and-pencil test assessing basic and advanced reading and writing ability; and (c) the Read India (RI) Math test, which is also a paper-and-pencil test assessing basic and advanced math ability. Several rounds of data were collected: (1) an initial pilot study (Pilot 1) with 256 children from Grades 1-5, (2) a second pilot study (Pilot 2) conducted with 412 children from Grades 1-5, (3) a baseline evaluation conducted in two districts in Bihar (n=8092) and Uttarakhand (n=7237) with children aged 5-16 from Grades 1-8, and (4) a

² Swaminathan, Hambleton, & Algina, 1974; Cohen, 1960

³ The formula for Cohen's Kappa is $\kappa = \frac{p_o - p_c}{1 - p_c}$ where p_o is the observed proportion of agreement between

raters and p_c is the expected proportion of agreement, i.e. agreement due to chance. The coefficient kappa of 0 occurs when observed agreement can be exactly accounted for by chance and the coefficient kappa of 1 occurs when there is complete agreement between raters. Kappa can yield a negative value when there is less observed agreement than is expected by chance.

⁴ Cohen, 1960, 1968

⁵ Landis & Koch, 1977

midline evaluation in Bihar conducted with 4807 children of ages 5-16 from Grades 1-8 who were assessed on the ASER-reading and -math tests and the Fluency Battery. Data from all tests are available for Pilot 1, Pilot2, and the Bihar baseline study. For the Uttarakhand baseline⁶ and the Bihar midline studies⁷ data from the ASER-reading, ASER-math, and the Fluency Battery are available.

The Tools and Measures

The Fluency Battery: The assessment of fluency is based on the premise that the ability to read fluently, i.e. with sufficient speed and accuracy is important to read well and to comprehend text. In fact, the fluent decoding of letters, letter combinations, words in list form, and words in connected text are important and robust correlates of early reading ability and comprehension. The automaticity of these lower-level skills ensures that limited cognitive resources such as attention and memory can be freed and allocated to the higher-level skills of meaning-making (LaBerge & Samuels, 1974; Perfetti, 1977, 1985). Hence, fluency measures, which are orally administered tests, are widely used to assess children's early reading ability in English and several other languages.

The Fluency Battery was adapted from the Early Grade Reading Assessment (USAID, 2009) and the Dynamic Indicators of Basic Early Literacy Skills (University of Oregon Center on Teaching and Learning, 2002). It comprises 8 subtests:

1. Akshar (letter) Reading Fluency (ARF): indicates the speed and accuracy with which children read aloud randomly arranged akshars of the Hindi alphasyllabary in a span of one minute. The score is the number of akshars named correctly in one minute.
2. Barakhadi Reading Fluency (BRF): indicates the speed and accuracy with which children read aloud randomly arranged consonant-vowel (CV) akshar units in one minute. All units were represented by a single consonant /k/ so as not to confound this task with the Akshar Reading Fluency subtest. The score is the number of barakhadi units named correctly in one minute.
3. Word Reading Fluency (WRF): indicates the speed and accuracy with which children read aloud a list of one- and two-syllable words in one minute. The score is the number of words read correctly in one minute.
4. Nonword Reading Fluency (NWRF): indicates the speed and accuracy with which children read aloud a list of two-syllable nonwords in one minute. The score is the number of nonwords read correctly in one minute.
5. Grade 1 Level Passage Reading Fluency (two passages, PRF): indicates the speed and accuracy with which children read aloud two Grade 1 level passages comprising 4 sentences and 21 words. The score is an average of the two passages and indexes the number of words read correctly in one minute.
6. Grade 2 level Passage Reading Fluency (two passages, SRF): indicates the speed and accuracy with which children read aloud two Grade 2 level passages comprising 6 sentences and 59-63 words. The score is an average of the two passages and indexes the number of words read correctly in one minute.

⁶ Data on the Read India Literacy and Math tests for the Uttarakhand sample are not used as several concerns remain about the credibility of these data.

⁷ The Read India Literacy and Math tests were not administered for the Bihar midline evaluation.

7. Grade 1 level Passage Comprehension Questions (PCOMP): comprises two comprehension questions for each of the two Grade 1 level passages. The score is the number of questions answered correctly.
8. Grade 2 level Passage Comprehension Questions (SCOMP): comprises four comprehension questions for each of the two Grade 2 level passages. The score is the number of questions answered correctly.

The content of the Fluency Battery was drawn from prior ASER reading tests as the material has been extensively evaluated and piloted to ensure their grade and content appropriateness for the population of interest. There was no overlap in test content of the ASER reading tests and the Fluency Battery. Scores for the fluency reading subtests represent number of units (akshars/ words/nonwords) read accurately in one minute and scores for the reading comprehension subtest represent number of questions correctly answered. Total administration time for the Fluency Battery is about 10 minutes. The median Cronbach's alpha estimates across the 5 samples ranged from .92 to .94 with a median Cronbach's alpha estimate of .93. Test-retest reliability coefficients for the subtests of the Fluency Battery ranged from .83 to .98. Since the association between the Fluency Battery sub-tests was high for all the 5 samples (r s range from .81 to .94), a single composite score of all the fluency sub-tests was created by taking an average.

The Literacy and Math Written Tests (Read India tests): A more traditional format of written tests of reading and math were developed to assess higher level reading, writing, and math skills. These tests were drawn from extensively piloted Urdu reading and math tests for use in Pakistan (Andrabi, Das, Khwaja, Farooqi, & Zajonc, 2002) and from the math tests of the TIMSS (<http://nces.ed.gov/timss/>). A few math items were also drawn from the ASER. The format of the questions on the Pakistan Urdu test were used as a reference to develop test items for Hindi, e.g. format of the reading vocabulary items, cloze sentences, maze passage⁸, etc. Care was taken to ensure that the content and item formats were appropriate for use in Hindi and aligned to the Bihar and Uttarakhand language and math curriculum, the two states that are part of the Read India evaluation study.

In order to ensure that the test items were appropriate for grade level, separate tests with overlapping content were designed for grades 1-2 and grades 3-5. The RI Literacy test for grades 1-2 requires about 15 minutes to administer and for grades 3-5 requires about 20 minutes to administer. The RI Math test for grades 1-2 requires about 10 minutes to administer and for grades 3-5 requires about 20 minutes to administer. Reliability based on internal consistency was estimated in two ways: (1) treating each item on the test as an individual item, and (2) treating each questions category on the test as an individual item thereby reducing the count of the number of items on the test. For the tests for Grades 1-2, the median Cronbach's alpha estimate for the first approach was .93 for the RI Literacy test and .93 for the RI Math test. The median Cronbach's alpha estimate for the second approach was .86 for the RI Literacy test and .86 for the RI Math test. For the tests for Grades 3-5, the median Cronbach's alpha estimate for

⁸ The Cloze and Maze test involve text where a few choice words are replaced by a blank. The Cloze is a fill-in-the-blank format where children are required to generate the missing word and Maze is a multiple-choice format where children are required to choose the target word from several choices.

the first approach was .93 for the RI Literacy test and .94 for the RI Math test. The median Cronbach's alpha estimate for the second approach was .88 for the RI Literacy test and .90 for the RI Math test (see Abdul Latif Jameel Poverty Action Lab, Pratham, & ASER, 2009 for a detailed evaluation of all Read India tests).

Analytic Plan

To assess *concurrent* validity⁹, we estimated the degree of association between the ASER-reading test and the Fluency Battery and the RI Literacy test using Spearman correlation coefficients. We expected the ASER-reading tests to be strongly correlated with both tests but we expected correlations of higher magnitude between the ASER-reading test and the Fluency Battery than with the RI Literacy test as the former two tests share a common inference about children's basic reading ability and are in the oral format. The association of the ASER-reading test and the RI Literacy test also helps us understand the relationship between literacy tests in the oral and written format. For the ASER-math test we estimated the degree of association between the ASER-math test and the RI Math test using Spearman correlation coefficients.

To assess *convergent-discriminant* validity¹⁰, we evaluated the differences in the estimated degree of association of the ASER-reading test with the other tests of literacy versus with the tests of math. Similarly, we evaluated the differences in the estimated degree of association of the ASER-math test with the test of math versus the tests of literacy. These were estimated separately for the 5 samples – Pilot 1, Pilot 2, the Bihar baseline, the Uttarakhand baseline, and the Bihar midline sample.

Comparing Performance on the ASER and the Fluency Battery: A Closer Look

Since the ASER-reading test and the Fluency Battery are tests of early reading ability, two additional sets of analysis were conducted for these two tests to better understand the appropriateness of the cut-off criteria used for the ASER-reading test. First, the fluency rates were examined for children classified at different reading levels based on the ASER-reading test. Second, the percentage of children on the Fluency Battery who read less than 3 akshars/words and more than 3 akshars/words was calculated.

⁹ Concurrent validity involves examining the strength of the association between performance on the new tests (ASER tests) and performance on already established standardized tests that share a common inference and which serve as the criterion tests. Strong and positive correlations between the two tests indicate strong evidence of concurrent validity.

¹⁰ Convergent-discriminant validity indicates that there is a stronger association among related constructs (e.g. reading tests) than among less related constructs (e.g. math test and reading test). In addition, a test of basic reading ability is expected to correlate more strongly with another test of basic reading ability than with a test of advanced reading ability. However, the association between tests of reading and math are also expected to be high as they draw on children's underlying cognitive ability (unless if the sample studied has a specific learning difficulty). Hence, the differences in the strength of the associations between tests of reading and tests of math are typically very small.

These percentages were calculated for each reading level of the ASER-reading test thus permitting an evaluation of decisions based on a short test such as the ASER that has only 5 items (akshars) on the akshar reading subtest and 5 items (words) on the word reading subtest versus the 52 akshars on the Akshar Reading Fluency subtest and 52 words on the Word Reading Fluency subtest, albeit with a time limit of 1 minute. These sets of analysis allow evaluating agreement between decisions across tests that are administered independently yet designed to assess the same abilities or skills.

Results

The concurrent validity coefficients range from .90 to .94 and indicate that the ASER-reading test is very highly correlated with the Fluency Battery. In addition, these coefficients indicate that the ASER-reading test is more strongly correlated with the Fluency Battery than it is with the ASER-math test. This pattern is evident across all the validity studies, i.e. pilot 1, pilot 2, the Bihar baseline, the Uttarakhand baseline, and the Bihar midline. In addition, the correlation coefficients indicate that the ASER-reading test is more strongly correlated with the RI Literacy test than with the math tests except for Grade 1-2 for the Bihar baseline study. In the latter case, the ASER-reading test is more strongly associated with the ASER-math test than with the RI Literacy test.

How did children at the different ASER-reading levels perform on the Fluency Battery?

The descriptive statistics for fluency rates for children at the different ASER-reading levels presented in Table 1 indicates that reading fluency rates increase with the increasing ASER-reading levels. In other words, children categorized at the Grade 2 story reading level (level 5) read the akshars, barakhadi, words, nonwords, and words in connected text with greater speed and accuracy than children classified at any of the lower levels of reading on the ASER-reading test. For instance, the fluency rates for akshars averaged across the four samples are about 2 for children at the 'nothing' level, about 17 for children at the 'akshar' level, 32 for children at the 'word' level, 45 for children at the 'para' level, and 62 for children at the 'story' level. These increasing fluency rates with higher ASER-reading levels are reflected in the strong validity coefficients between the ASER-reading test and the Fluency Battery noted earlier.

Given that the ASER-reading levels are mutually exclusive categories, children classified at the 'akshar' level are seen to demonstrate competency at the akshar level but not at the word level, and so on. It follows then that children at the 'nothing' level should perform poorly on the akshar reading fluency subtest and children at the 'akshar' level should perform poorly on the word reading fluency subtest and so on. Average performances presented in Table 1 substantiate this claim. For instance, averaging across the 4 samples, children classified at the 'nothing' level demonstrate akshar fluency rates of 2 akshars, children classified at the 'akshar' level demonstrate word fluency rates of 3 words, children classified at the 'word' level demonstrate Grade 1 level oral fluency rates of 25 words, children classified at the Grade 1 passage level demonstrate Grade 2 level oral fluency rates of 44 words¹¹.

¹¹ Much variation is noted for the fluency rates as some children demonstrate high fluency rates despite being categorized at lower levels of reading. A few instances of misclassification referred to as decision inconsistency is to be expected. However, the percentage of these misclassifications is on the low side (see next set of analysis). This warrants further examination if the ASER-reading tests are to be used for diagnostic purposes or for making decisions at the individual rather than the group level.

The ASER akshar and word reading subtests are extremely short tests that comprise only 5 items. As a result it is possible that children can be misclassified due to item sampling error. To evaluate the efficacy of such a short test the percentage of children who identified no akshars/words, who identified less than 4 akshars/words and who identified >4 akshars/words on the Akshar and Word Reading Fluency subtests was calculated. This enabled comparing children's performance on the ASER akshar and word reading subtests with performance on the akshar and word reading fluency subtests that comprise all the akshars of the Hindi alphasyllabary and a substantially larger number of words.

Results presented in Table 2 indicate that of the children classified at the 'nothing' level, 82% of the children in Uttarakhand, 94% of the children in the Bihar baseline study, and 95% of the children in the Bihar midline study could not correctly identify 4 or more akshars on the Akshar reading fluency subtest. Of the children classified at the 'akshar' level 96% of the children in Uttarakhand, 80% of the children in the Bihar baseline study, and 85% of the children in the Bihar midline study could in fact correctly identify 4 or more akshars.

Of the children classified at the 'word' level 98% of the children in Uttarakhand, 87% of the children in the Bihar baseline study, and 96% of the children in the Bihar midline study did correctly read 4 or more words correctly. This is a high level of consistency across the two tests. However there are children who were classified at the 'nothing' level who correctly read more than 3 akshars in one minute on the Akshar Reading Fluency subtest and there were children classified at the 'akshar' level who correctly read more than 3 words (Table 2). Further examination of the fluency rates for these decision inconsistencies indicates that although the children categorized at the 'nothing' level read 4 or more akshars correctly on the Akshar Reading Fluency subtest, they demonstrated low rates of fluency in comparison to their counterparts who were categorized at the 'akshar' level (Table 3 presents descriptives and Figures 1a-1c present score distributions). Similarly, children categorized at the 'akshar' level read 4 or more words correctly on the Word Reading Fluency subtest, but they demonstrated low rates of fluency in comparison to their counterparts who were categorized at the 'word' level (Table 4 presents descriptives and Figures 2a-2c present score distributions).

Table 1: Descriptive statistics for the Fluency Battery for children classified at different ASER-reading levels

ASER-Reading Level	Fluency Battery								
	Akshar Reading Fluency	Barakhadi Reading Fluency	Word reading fluency	Nonword reading fluency	Grade 1 level passage reading fluency	Grade 2 level passage reading fluency	Oral reading fluency	Grade 1 level comprehension	Grade 2 level comprehension
<i>Uttarakhand Baseline</i>									
Nothing Level (n=1775)	1.90 (4.2)	0.59 (2.1)	0.09 (0.60)	0.03 (0.52)	0.17 (1.43)	0.12 (1.17)	0.15 (1.18)	0.06 (0.32)	0.03 (0.27)
Akshar Level (n=1726)	22.65 (12.34)	8.46 (10.64)	4.41 (5.70)	1.43 (2.85)	7.23 (11.89)	6.67 (9.86)	6.95 (10.68)	0.97 (1.21)	0.87 (1.76)
Word Level (n=470)	37.40 (13.19)	28.50 (17.19)	17.90 (11.44)	8.18 (6.73)	29.98 (22.27)	26.37 (17.66)	28.17 (19.58)	2.50 (1.19)	3.91 (2.11)
Para Level (n=847)	45.05 (14.75)	39.47 (18.27)	28.14 (14.37)	13.28 (8.75)	52.49 (27.71)	44.60 (22.51)	48.54 (24.36)	3.12 (1.02)	5.27 (2.06)
Story Level (n=2361)	64.63 (21.33)	68.28 (22.28)	67.22 (25.68)	38.28 (18.36)	116.20 (42.39)	101.69 (36.36)	108.95 (38.15)	3.61 (0.65)	6.59 (1.44)
<i>Bihar Baseline</i>									
Nothing Level (n=4078)	1.10 (4.88)	0.84 (5.37)	0.39 (4.26)	0.28 (3.71)	0.38 (6.56)	0.38 (5.51)	0.38 (5.92)	0.02 (0.26)	0.03 (0.38)
Akshar Level (n=1564)	15.41 (13.54)	9.33 (13.25)	3.88 (8.46)	2.22 (7.36)	4.11 (11.05)	4.41 (11.17)	4.26 (10.53)	0.29 (0.81)	0.32 (1.08)
Word Level (n=592)	30.86 (14.87)	26.85 (19.42)	18.23 (14.38)	11.30 (11.33)	22.30 (21.16)	20.63 (19.17)	21.46 (18.92)	1.50 (1.43)	1.48 (1.98)
Para Level (n=836)	45.50 (20.15)	47.20 (23.69)	37.63 (21.43)	23.44 (16.04)	53.78 (34.19)	46.19 (28.99)	49.99 (29.00)	2.76 (1.28)	3.82 (2.42)
Story Level (n=1796)	63.37 (24.48)	71.16 (28.05)	62.82 (27.99)	40.93 (20.93)	93.04 (48.00)	78.87 (37.53)	85.67 (39.30)	3.49 (0.86)	5.98 (2.05)
<i>Bihar Midline</i>									
Nothing Level (n=2193)	1.03 (2.32)	0.73 (2.53)	0.07 (1.67)	0.03 (0.88)	0.11 (3.08)	0.09 (2.89)	0.10 (2.96)	0.02 (0.18)	0.02 (0.23)
Akshar Level (n=1135)	12.64 (10.67)	5.41 (6.25)	2.10 (3.27)	0.70 (1.62)	2.43 (5.35)	2.16 (4.63)	2.29 (4.87)	0.52 (0.93)	0.25 (0.89)
Word Level (n=337)	30.61 (11.54)	21.42 (14.07)	12.93 (8.55)	7.07 (5.31)	18.88 (15.46)	16.72 (11.98)	17.80 (13.37)	1.78 (1.35)	2.50 (2.11)
Para Level (n=595)	42.70 (13.17)	40.74 (18.20)	28.27 (13.69)	15.37 (8.95)	46.33 (25.27)	39.66 (23.89)	42.99 (23.04)	2.82 (1.09)	4.05 (2.30)
Story Level (n=1352)	65.58 (21.60)	73.34 (24.22)	62.89 (25.35)	35.51 (17.00)	104.10 (41.39)	90.13 (36.67)	97.11 (37.96)	3.51 (0.76)	6.15 (1.91)
<i>Pilot 1</i>									
Nothing Level (n=34)	2.21 (2.10)	2.18 (4.07)	0.32 (1.09)	0.06 (0.34)	0.03 (0.17)	0 (0)	0.01 (0.09)	0.03 (0.17)	0 (0)
Akshar Level (n=66)	15.63 (10.48)	9.35 (11.69)	3.56 (5.25)	1.35 (2.59)	5.16 (12.26)	4.95 (9.73)	5.05 (10.47)	0.50 (0.98)	0.45 (1.21)
Word Level (n=23)	29.68 (11.55)	28.65 (16.03)	17.18 (12.40)	7.17 (5.44)	30.73 (28.75)	25.06 (19.03)	27.90 (23.12)	2.30 (1.11)	2.70 (2.38)
Para Level (n=40)	45.79 (14.51)	47.34 (19.57)	32.11 (15.49)	16.67 (8.11)	54.70 (27.70)	45.15 (21.34)	49.52 (22.95)	3.13 (0.69)	4.65 (1.76)
Story Level (n=93)	55.39 (17.48)	64.76 (18.21)	56.83 (19.01)	30.32 (11.79)	94.13 (34.10)	79.33 (27.94)	87.22 (28.91)	3.41 (0.77)	5.85 (2.06)

Table 2: Percentage of decision consistencies and inconsistencies across the Fluency Battery and the ASER-reading test.

	Uttarakhand Baseline					Bihar Baseline					Bihar Midline				
	Nothing Level	Akshar Level	Word Level	Para Level	Story Level	Nothing Level	Akshar Level	Word Level	Para Level	Story Level	Nothing Level	Akshar Level	Word Level	Para Level	Story Level
	(n=1775)	(n=1726)	(n=470)	(n=847)	(n=2361)	(n=4078)	(n=1564)	(n=592)	(n=836)	(n=1796)	(n=4078)	(n=1564)	(n=592)	(n=836)	(n=1796)
<i>Akshar reading fluency</i>															
0-3 akshars	81.98%	3.65%	0.00%	0.00%	0.00%	93.55%	20.46%	3.38%	2.03%	0.50%	94.80%	15.42%	0.59%	0.17%	0.00%
> 3 akshars	18.03%	96.35%	100.00%	100.00%	100.00%	6.45%	79.54%	96.62%	97.97%	99.50%	5.20%	84.58%	99.41%	99.83%	100.00%
<i>Word reading fluency</i>															
0-3 words	99.38%	57.99%	2.34%	0.71%	0.04%	98.62%	72.18%	13.17%	2.87%	0.72%	99.59%	80.17%	4.45%	0.17%	0.07%
>3 words	0.62%	42.00%	97.66%	99.29%	99.96%	1.37%	27.81%	86.82%	97.13%	99.28%	0.41%	19.82%	95.55%	99.83%	99.93%

Table 3: Descriptive statistics of akshar fluency rates for children whose akshar fluency rates are 4 or more and were categorized at the 'nothing' level or 'akshar' level on the ASER-reading test

		M	SD
Uttarakhand Baseline	Nothing level (n=320)	8.76	6.21
	Akshar level (n=1663)	23.46	11.84
Bihar Baseline	Nothing level (n=263)	13.75	13.91
	Akshar level (n=1244)	19.19	12.66
Bihar Midline	Nothing level (n=114)	7.69	7.02
	Akshar level (n=960)	14.67	10.39

Table 4: Descriptive statistics of akshar and word fluency rates for children whose word fluency rates are 4 or more and were categorized at the 'akshar' level or 'word' level on the ASER-reading test

		Akshar fluency rates		Word fluency rates	
		M	SD	M	SD
Uttarakhand Baseline	Akshar level (n=725)	31.24	10.79	9.06	6.17
	Word level (n=459)	37.66	13.15	18.28	11.32
Bihar Baseline	Akshar level (n=435)	26.73	13.46	12.94	11.9
	Word level (n=514)	33.28	13.55	20.9	13.56
Bihar Midline	Akshar level (n=225)	23.93	9.2	7.2	4.1
	Word level (n=332)	31.24	11.22	13.45	8.39

Figure 1: Distribution of akshar reading fluency rates for children whose akshar fluency rates are 4 or more and were categorized at the 'nothing' level or 'akshar' level on the ASER-reading test

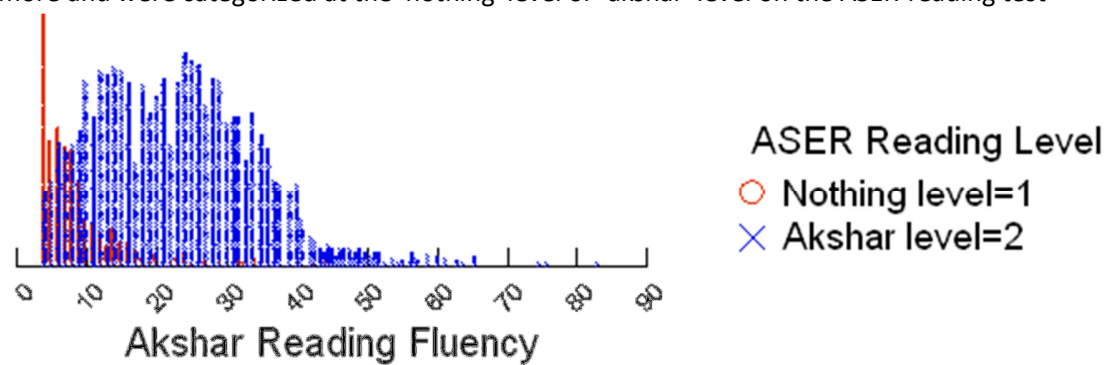


Figure 1a: Uttarakhand Baseline

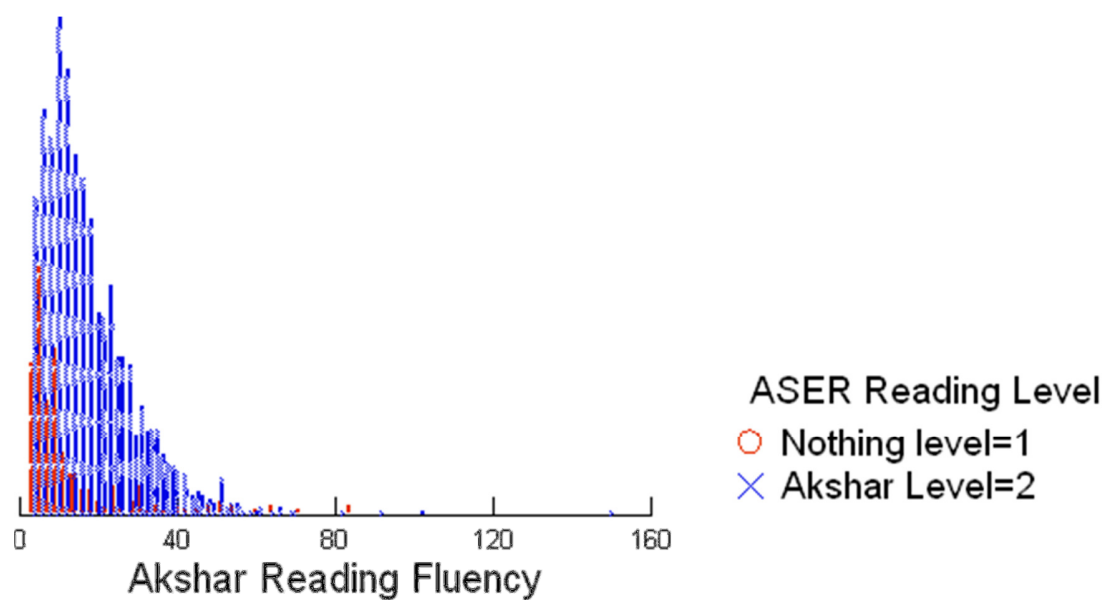


Figure 1b: Bihar Baseline

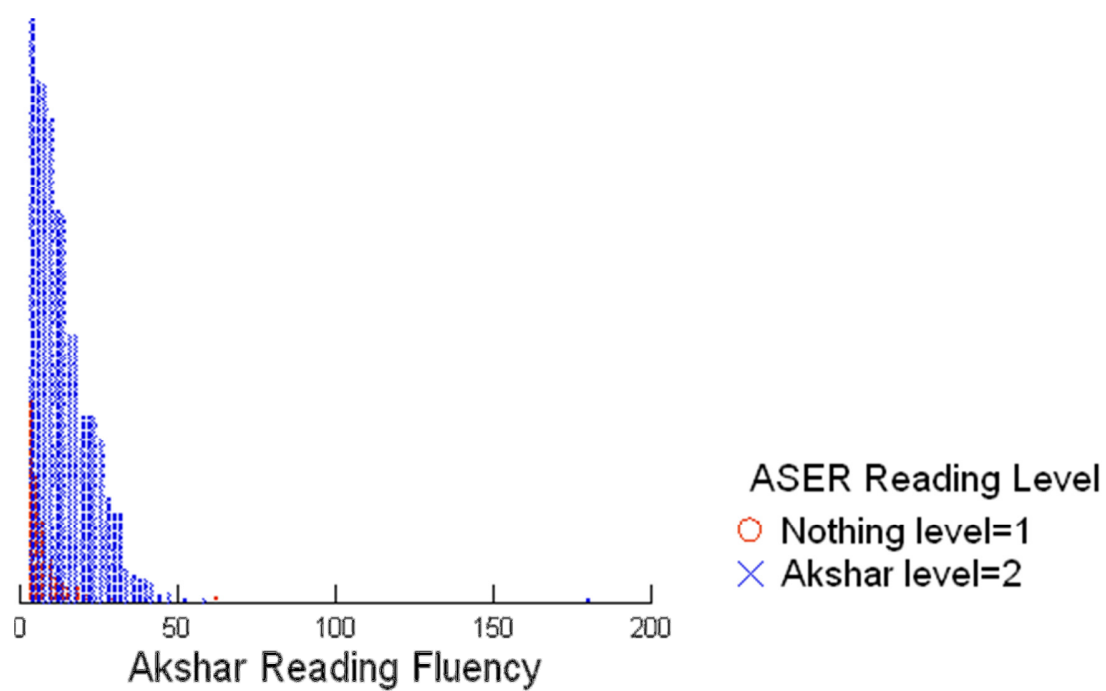


Figure 1c: Bihar Midline

Figure 2: Distribution of scores of akshar and word fluency rates for children whose word fluency rates are 4 or more and were categorized at the 'akshar' level or "word" level on the ASER-reading test

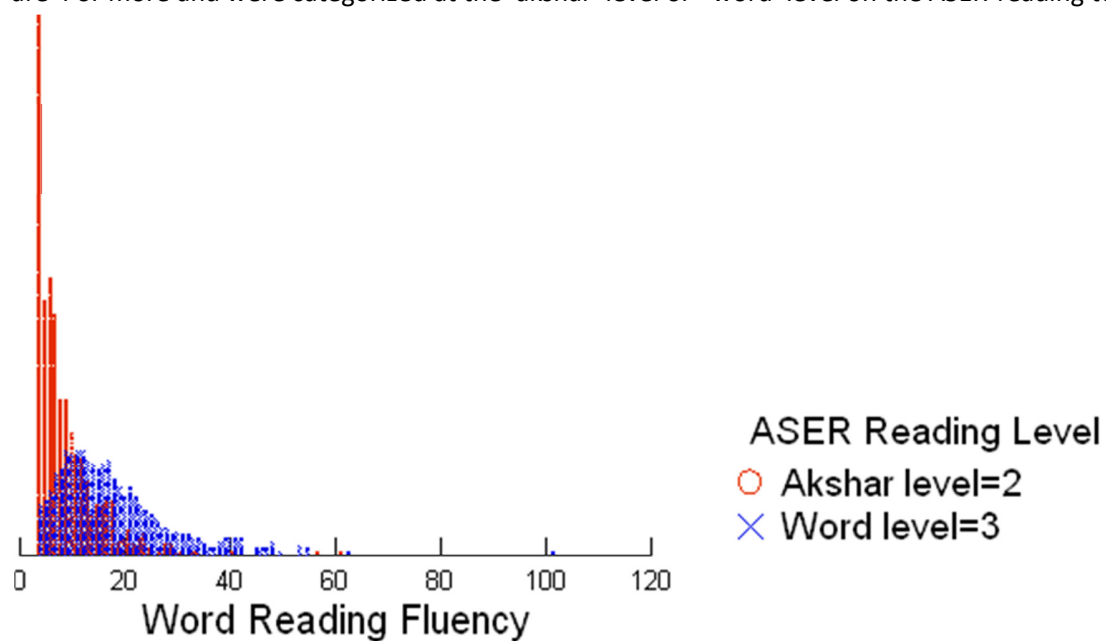


Figure 2a: Uttarakhand Baseline

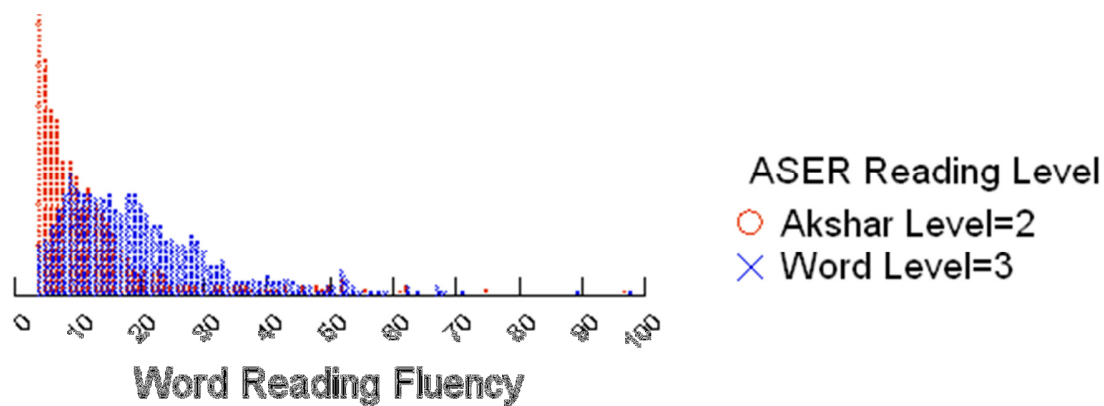


Figure 2b: Bihar Baseline

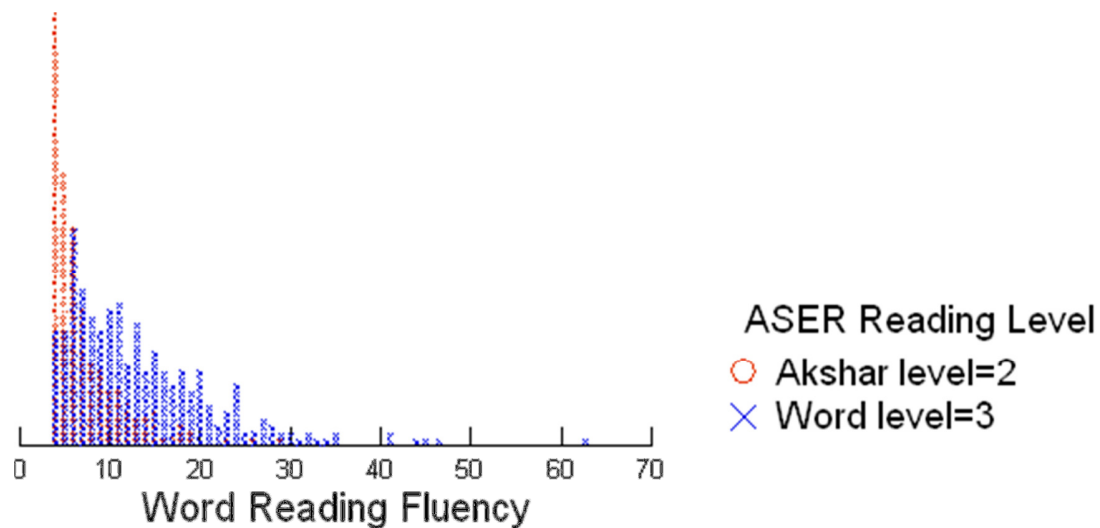


Figure 2c: Bihar Midline

Discussion

The ASER-reading and ASER-math tests are simple, quick, easy to administer and used primarily to obtain school-level and district-level data about children's foundational reading skills and basic math ability. The findings based on a series of studies reported in the present paper provide favorable empirical evidence for the reliability and validity of these tests. Specifically, the findings indicate substantial reliability of decisions across repeated measurements, satisfactory inter-rater reliability and favorable evidence for concurrent and convergent-discriminant validity.

Compelling evidence for the validity of the ASER tests is illustrated in (a) the very strong associations of the ASER-reading test with the concurrently administered Fluency Battery, which like the ASER-reading test assesses foundational reading skills, (b) the stronger association of the ASER-reading test with the Fluency Battery than with the RI Literacy test, which unlike the Fluency Battery also assesses advanced

reading and writing ability, (c) the stronger association of the ASER-reading test with the Fluency Battery and the RI Literacy test than with the math tests, and (d) the stronger association of the ASER-math test with the RI Math test than with the tests of literacy.

Additional comparisons of the decision consistency between the ASER-reading test and the Fluency Battery indicate that there is a high level of consistency across the two tests at the 'nothing', 'akshar', and 'word' level. Although there were some inconsistencies with children at the 'nothing' level correctly reading 4 or more 'akshars' on the Akshar Reading Fluency subtest and with children at the 'akshar' level correctly reading 4 or more words on the Word Reading Fluency subtest, the respective fluency rates were clustered at the lower end of the continuum. Moreover, given that the ASER reading levels are mutually exclusive categories it implies that children who demonstrate competency at the akshar level do not demonstrate competency at the word or any other higher level. As a result, the fluency rates of children at the akshar level are bound to be lower than the fluency rates of children who are classified at the word or higher level. This expectation is supported by the data and is in keeping with the viewpoint that fluency in reading words in connected text requires fluency at the levels of smaller units such as letters (akshars) and letter combinations (barakhadi) (Foulin, 2005, Wolf & Katzir-Cohen, 2001). Consequently, an important instructional implication of this finding is that children categorized at the 'akshar' level are demonstrating 'minimal' mastery as opposed to 'complete' mastery of akshar knowledge and need to further improve their akshar knowledge if they are to successfully decode words in list form or connected text. Similarly, children classified at the 'word' level are demonstrating 'minimal' mastery of their decoding knowledge and need to further improve their decoding skills in order to fluently read and comprehend words in connected text.

Finally, although the association between the ASER-reading test and the Fluency Battery is very strong and they both assess foundational reading skills, the decision to use any one of the tests should be based on considerations of the purpose of testing and the nature of information desired. The ASER-reading test provides information about children's reading levels in mutually exclusive ordinal ranks, whereas the Fluency Battery provides information about children's reading in terms of fluency at different levels of reading (akshars/words read correctly in one minute). Hence, both tests have their merits depending on the purpose of assessment. For instance, using the Fluency Battery along with the ASER-reading test in the evaluation of the Read India intervention program enables the assessment of children's progress in reading within and across reading levels. On the other hand, using the ASER-reading test for the nationwide ASER survey provides a reliable and valid snapshot of children's foundational reading skills in a simple, quick, cost-effective manner with results that are easy for policy makers, educators, and parents to understand and which are available in the same academic year.

References

- Abdul Latif Jameel Poverty Action Lab [J-PAL], Pratham, & ASER, (2009). *Evaluating READ INDIA: the development of tools for assessing Hindi reading and writing ability and math skills of rural Indian children in grades 1-5*. Unpublished manuscript: J-PAL, Chennai, India.
- Andrabi, Das, Khwaja, Farooqi, & Zajonc. (2002). Test feasibility survey Pakistan: Education Sector.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20, 37-46.
- Cohen, J. (1968). Weighted kappa: Nominal scale agreement with provision for scaled disagreement or partial credit. *Psychological Bulletin*, 71, 213-220.
- Cronbach, L.J. (1971). Test validation. In R.L.Thorndile (Ed.), *Educational measurement* (2nd ed.). Washington, DC: American Council on Education.
- Foulin, J.N. (2005). Why is letter-name knowledge such a good predictor of learning to read? *Reading and Writing*, 18, 129-155.
- LaBerge, F., & Samuels, S.J. (1974). Toward a theory of automatic information processing in reading. *Cognitive Psychology*, 6, 293-323.
- Landis, J.R., and Koch, G.G. (1977). The measurement of observer agreement for categorical data. *Biometrics*, 33, 159-174.
- Perfetti, C.A. (1977). Literacy comprehension and fast decoding. Some psycholinguistic prerequisites for skilled reading comprehension. In J.T.Guthrie (Ed.), *Cognition, curriculum, and comprehension* (pp. 20-41). Newark, DE: International Reading Association.
- Perfetti, C.A. (1985). *Reading ability*. London: Oxford.
- Pratham (2005). Annual Status of Education Report (ASER). Retrieved July 1, 2009 from the World Wide Web: <http://asercentre.org/asersurvey/aser05.php>
- Pratham (2006). Annual Status of Education Report (ASER). Retrieved July 1, 2009 from the World Wide Web: <http://asercentre.org/asersurvey/aser06.php>
- Pratham (2007). Annual Status of Education Report (ASER). Retrieved July 1, 2009 from the World Wide Web: <http://asercentre.org/asersurvey/aser07.php>
- Pratham (2008). Annual Status of Education Report (ASER). Retrieved July 1, 2009 from the World Wide Web: <http://asercentre.org/asersurvey/aser08/pdfdata/aser08.pdf>
- Swaminathan, H., Hambleton, R.K., & Algina, J. (1974). Reliability of criterion-referenced tests: a decision-theoretic formulation. *Journal of Educational Measurement*, 11(4), 263-267.
- Traub, R.E., & Rowley, G.L. (1980). Reliability of test scores and decisions. *Applied Psychological Measurement*, 4(4), 517-545.
- Trends in International Mathematics and Science Study [TIMSS] (2003). *Mathematics Items: Grade 4*. Retrieved February 1, 2008 from the World Wide Web: http://nces.ed.gov/TIMSS/pdf/TIMSS4_Math_Items.pdf
- USAID (2009). Early grade reading assessment (EGRA). Retrieved July 1, 2009 from the World Wide Web: <http://www.eddataglobal.org/documents/index.cfm?fuseaction=pubDetail&ID=95>
- University of Oregon Center on Teaching and Learning. (2002). *Dynamic indicators of basic early literacy skills [DIBELS]: Analysis of reading assessment measures*. Retrieved February 1, 2008 from the World Wide Web: http://dibels.uoregon.edu/techreports/dibels_5th_ed.pdf

- Wagner, D.A. (2003). Smaller, quicker, cheaper: alternative strategies for literacy assessment in the UN Literacy Decade. *International Journal of Educational Research*, 39, 293-309.
- Wolf, M., & Katzir-Cohen, T. (2001). Reading fluency and its intervention. *Scientific Studies of Reading*, 5(3), 211-239.