



# **Equating Test Scores for Receptive Vocabulary Across Rounds and Cohorts in Ethiopia, India and Vietnam**

Juan Leon and Abhijeet Singh



# Equating Test Scores for Receptive Vocabulary Across Rounds and Cohorts in Ethiopia, India and Vietnam

Juan Leon and Abhijeet Singh

First published by Young Lives in May 2017

© Young Lives 2017

Printed on FSC-certified paper from traceable and sustainable sources.

## About Young Lives

Young Lives is an international study of childhood poverty, following the lives of 12,000 children in four countries (Ethiopia, India, Peru and Vietnam) over 15 years. [www.younglives.org.uk](http://www.younglives.org.uk)

Young Lives is core-funded by UK aid from the Department for International Development (DFID).

The views expressed are those of the author(s). They are not necessarily those of, or endorsed by, Young Lives, the University of Oxford, DFID or other funders.

Core-funded by



**Young Lives**, Oxford Department of International Development (ODID), University of Oxford,  
Queen Elizabeth House, 3 Mansfield Road, Oxford OX1 3TB, UK

Tel: +44 (0)1865 281751 • Email: [younglives@younglives.org.uk](mailto:younglives@younglives.org.uk)

# Contents

<b>The authors</b>	<b>4</b>
<b>Acknowledgements</b>	<b>4</b>
<b>1. Introduction</b>	<b>5</b>
<b>2. The Peabody Picture of Vocabulary Test - III (PPVT - III)</b>	<b>5</b>
<b>3. Methodology</b>	<b>6</b>
3.1. The Young Lives study	6
3.2. Why IRT scores instead of CTT scores?	7
3.3. The IRT model: the three-parameter model	7
3.4. Item fit	8
3.5. Differential Item Functioning (DIF)	9
3.6. Scores equating	10
3.7. Limitations of IRT scores	11
3.8. Vocabulary test approach for Round 4	11
<b>4. Results</b>	<b>12</b>
<b>5. Final remarks</b>	<b>18</b>
<b>References</b>	<b>19</b>
<b>Appendices</b>	<b>20</b>
Appendix A. Details of the STATA analysis performed	20
Appendix B. ICC for each item by country and main language	24
Appendix C. DIF analysis for each item by country and main language	44
Appendix D. Item analysis performed by country and main languages	64
Appendix E. Item parameter for all the equated scales estimated	84
Appendix F. ICC curves for IRT equating analysis performed with siblings in Round 4	102
Appendix G. Item parameters used to equate sibling scores with main survey sample	113

## The authors

**Juan León** has a PhD in Educational Theory and Policy and Comparative and International Education from Pennsylvania State University, United States. He has a Bachelor's degree in Economics and a diploma in Liberal Arts from the Pontificia Universidad Católica of Peru. Juan is an Associate Researcher at GRADE in Peru. He is also a Lecturer in the Department of Psychology at the Universidad Antonio Ruiz de Montoya in Lima.

**Abhijeet Singh** is a post-doctoral researcher in Economics Development at University College London, and was previously a Research Officer with Young Lives. His research applies micro-econometric methods to the study of salient issues of policy interest in developing countries, focusing especially on issues relating to the economics of education and child health and the delivery of public services.

## Acknowledgements

We wish to thank Yessenia Collahua for her research assistance.

## 1. Introduction

In longitudinal studies such as Young Lives, getting comparable measures of children's cognitive abilities over time is essential for identifying individual, household, and school-level factors that affect children's development. Few longitudinal studies that follow birth/age cohorts include comparable cognitive measures across waves, and those studies that are available are mainly from developed countries. Young Lives provides a unique opportunity to explore the development of value added or growth curve modelling analysis aimed at identifying variables at different levels and across time and space, associated with children's learning outcomes, in developing countries.

This Technical Note discusses the construction of cognitive scores that are comparable across rounds and age cohorts for Young Lives in Ethiopia, India (the states of Andhra Pradesh and Telangana) and Vietnam. Young Lives gathers information from children and their families through individual and household questionnaires, including different cognitive and achievement tests. The Peabody Picture Vocabulary Test (PPVT) is the one test that is common across rounds and cohorts. Therefore, this test was selected to build cognitive measures comparable across Rounds 2, 3 and 4, and age cohorts (the Younger Cohort, born in 2001/02, and Older Cohort, born in 1994/95) employing Item Response Theory (IRT) to achieve standardised cognitive measures. Scores were estimated using a three-parameter model which considers the item's difficulty, discrimination, and pseudo-guessing as parameters to estimate the individual's ability. The second step was to perform a Differential Item Functioning analysis (DIF) by cohort and survey round in order to identify possible item bias and correct it. The last step consisted of equating the scores of common items (anchor items) as a means of obtaining comparable PPVT scores across rounds and cohorts without cohort and round biases.

This note has five sections. The second section presents a brief description of the Peabody Picture of Vocabulary test. Then, Section 3 presents the methodology of analysis and Section 4 explains the main results. Section 5 provides some final remarks on the main findings of the analysis performed.

## 2. The Peabody Picture of Vocabulary Test - III (PPVT - III)

The Peabody Picture Vocabulary Test (PPVT) is a widely used test to measure receptive vocabulary. It was originally developed in English in 1959 and has been updated several times. This study used version III (Dunn and Dunn 1997) in Ethiopia, India, and Vietnam; this was the version available for Rounds 2, 3 and 4. The PPVT test is administered individually, orally, untimed, and norm-referenced, where the test taker selects the picture that best represents the meaning of a stimulus word presented orally by the examiner. Not all the items in the PPVT are expected to be administered. Instead, the examiner administers enough items to establish both a ceiling and a baseline. The rule to set the baseline is based on making one error or no errors, in a set of 12 items. The rule to set the ceiling, on the other hand, is based on making eight or more errors in a set of 12 items. Non-administered items below the baseline are automatically given a score of 1, given that they are expected to be easier, while items above the ceiling are given a score of 0, given that they are more difficult. The raw score is formed by all the items given a score of 1 (i.e. answered correctly or below the basal item).

## 3. Methodology

### 3.1. The Young Lives study

Young Lives is a longitudinal study of childhood poverty that examines the development of around 12,000 children, from two cohorts born in 1994 and 2001, over 15 years in Ethiopia, India (in the states of Andhra Pradesh and Telangana), Peru and Vietnam.

In order to identify which factors affect children's development, it is necessary to have comparable measures of children's cognitive abilities over time. Table 1 shows the measures of abilities and achievement administered by round and cohort. Because the PPVT is the one test that has been administered consistently across rounds and cohorts, we selected it in order to build comparable cognitive measures for Rounds 2, 3 and 4, and for both the Younger and Older Cohort. The PPVT, however, was administered only to the Younger Cohort in Round 4. This change was informed by the ceiling effects observed for the Older Cohort in Round 3 (see Cueto and Leon 2012).

It is important to note that the PPVT was originally developed to measure receptive vocabulary in English. The versions of the test administered in Ethiopia, India and Vietnam were therefore translated into the main languages in each country. However, it became evident that some items did not keep the same cognitive equivalence (or level of difficulty) as in the original tests; much less between languages. Therefore, in Round 4, a subsample of the original 204 items was selected for each country. Only Peru administered the full set of items of the PPVT since they use the Spanish version of the test. For Ethiopia and India, around one quarter of the total number of items were selected, while in Vietnam a third were selected. The selection criteria were: (i) adequate item fit using data from Round 2 and 3; (ii) items without DIF by round and cohort using data from Rounds 2 and 3; and (iii) items across the different range of item difficulty.

**Table 1.** Measures of abilities and achievement administered in Young Lives

Round	Cohort	Cognitive	Reading	Mathematics
Round 1	Younger Cohort	NA	NA	NA
	Older Cohort	Raven's Progressive Matrices for children	One item on reading One item on writing	One item on multiplication
Round 2	Younger Cohort	PPVT	NA	CDA
	Older Cohort	PPVT	One item on reading One on writing	One multiplication item and maths test
Round 3	Younger Cohort	PPVT	One item on reading One item on writing The Early Grade Reading Assessment (EGRA)	One multiplication item and maths test
	Older Cohort	PPVT	Cloze test of reading comprehension	Maths test
Round 4*	Younger Cohort	PPVT	Reading comprehension	Maths test
	Older Cohort	NA	Reading comprehension	Maths test

Notes: NA = Not administered.

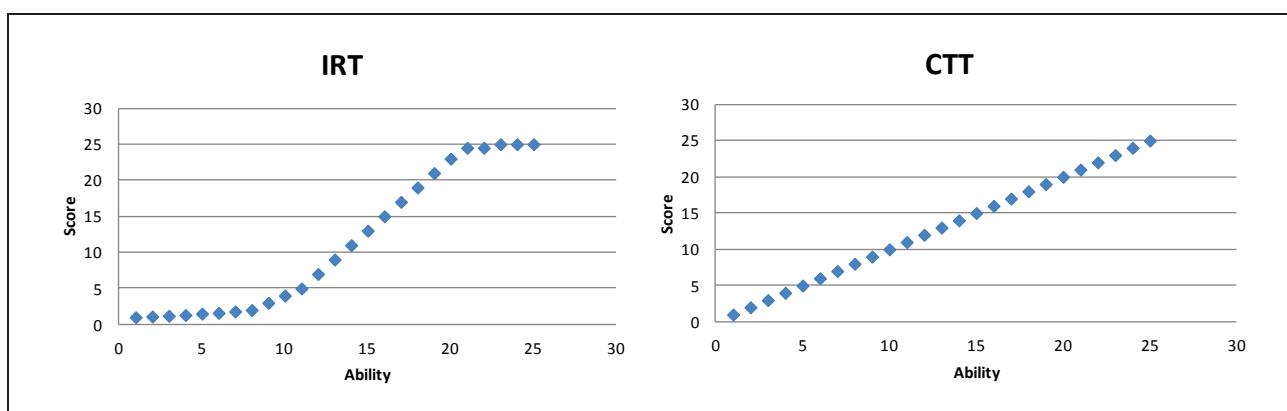
\*Round 4 considered PVVT for only 125 items in Peru, unlike other countries where a sub-sample of the original 204 items was administered.

### 3.2. Why IRT scores instead of CTT scores?

Unlike Classical Test Theory (CTT), Item Response Theory (IRT) is more focused on the item rather than the test. Moreover, the standard error of measurement in IRT is a function of the ability of individuals, therefore it varies at each level of ability. Thus, IRT estimates the probability of answering the item correctly through a logistic function based on the difference between the item difficulty and the individual's ability. The idea is that individuals with higher ability will have a greater probability of answering correctly easier items than difficult ones.

Figure 1 shows the relationship between the ability and the score in CTT as well as in IRT. In the case of CTT, we observe that the raw score increases in the same proportion as the ability, thus it follows a linear and monotonic trend. In contrast, in IRT we observe that as the ability increases the score does not increase in the same proportion, in other words the growth of scores is nonlinear. This implies that, under CTT, an individual will have the same ability if it changes from 10 to 15, or 20 to 25. However, under IRT an individual's ability will not be the same since it follows a different functional form that relies on the characteristics of the items.

**Figure 1.** Functional form between scores and ability, by theory



The main advantages of using the IRT statistical technique instead of CTT are: (i) the item parameters do not depend on an individual's ability, being invariant over different samples of examinees, and an individual's ability does not depend on the items presented, being invariant over different samples of items (i.e. principle of invariance); (ii) it allows comparisons of individual's ability from different populations if tested with instruments that have common items; and (iii) it allows the allocation of individuals' ability and items difficulty in the same scale or metric, creating an interval scale in logits for both scores. Thus, using this statistical technique, we were able to build comparable scores by cohort and round.

### 3.3. The IRT model: the three-parameter model

The IRT model relies on two main assumptions. First, the model assumes *local independence*, which means that the probability of answering an item correctly depends on an individual's ability only and not on his/her answer to other items. Second, it assumes *unidimensionality*: the model considers that only one latent trait is measurable across all items or that at least one dominant factor is observed behind the set of items tested. Of these

two assumptions, the latter is the most difficult to accomplish since different factors could be affecting the individual performance (for example, test anxiety).<sup>1</sup>

Under the IRT model used in this note, an individual's ability depends on three item parameters – item difficulty, item discrimination, and item guessing. Item difficulty refers to the proportion of individuals who get each item right. Item discrimination indicates how well an item discriminates between high and low achievers, while the guessing parameter refers to the chances that an individual has to get an item right. This parameter is mainly considered for multiple choice tests since these allow examinees to guess.<sup>2</sup> These parameters and the individual's ability level are part of the Item Characteristic Curve (ICC) that defines the probability that each individual has to get an item right given the item characteristics (difficulty, discrimination and guessing) and individual ability. The following equation represents the general ICC model:

$$P_i(\theta) = c_i + (1 - c_i) \frac{e^{a_i(\theta - b_i)}}{1 + e^{a_i(\theta - b_i)}} \quad i = 1, 2, \dots, n$$

$P_i(\theta)$  : the probability that an individual with ability  $\theta$  get item  $i$  right

$a_i$  : item discrimination

$b_i$  : the item difficulty

$c_i$  : guessing parameter

$n$  : the number of items in the test

$\theta$  : individual's ability parameter

The two-parameter model uses the same equation but assumes that the guessing parameter ( $c_i$ ) is equal to zero, while the one-parameter model not only assumes a guessing parameter ( $c_i$ ) of zero but also that the item discrimination ( $a_i$ ) is constant across items.

### 3.4. Item fit

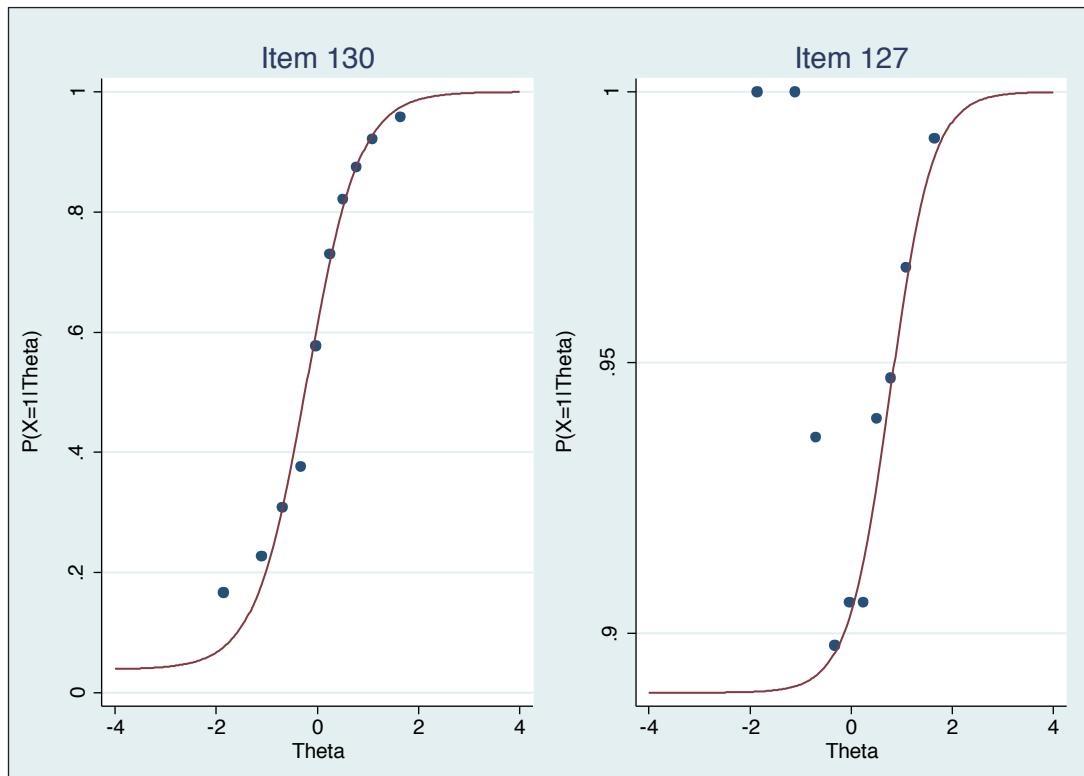
An item has *good fit* if the ICC shows that the proportion of children who answer an item correctly varies monotonically as a function of child's ability. As an example, Figure 2 shows an item with good fit (item 130) and an item with poor fit (item 127). We observe for the item with good fit that the proportion of children who correctly answer the item varies monotonically with the average child's ability, while the item with bad fit shows no correlation between the proportion of children who correctly answer an item with average child's ability.

---

1 For more information, see Cueto, Leon, Guerrero and Muñoz (2009).

2 As PPVT is a multiple choice test, it is necessary to consider a guessing parameter.

**Figure 2.** Item characteristic curves of items with good and bad fit

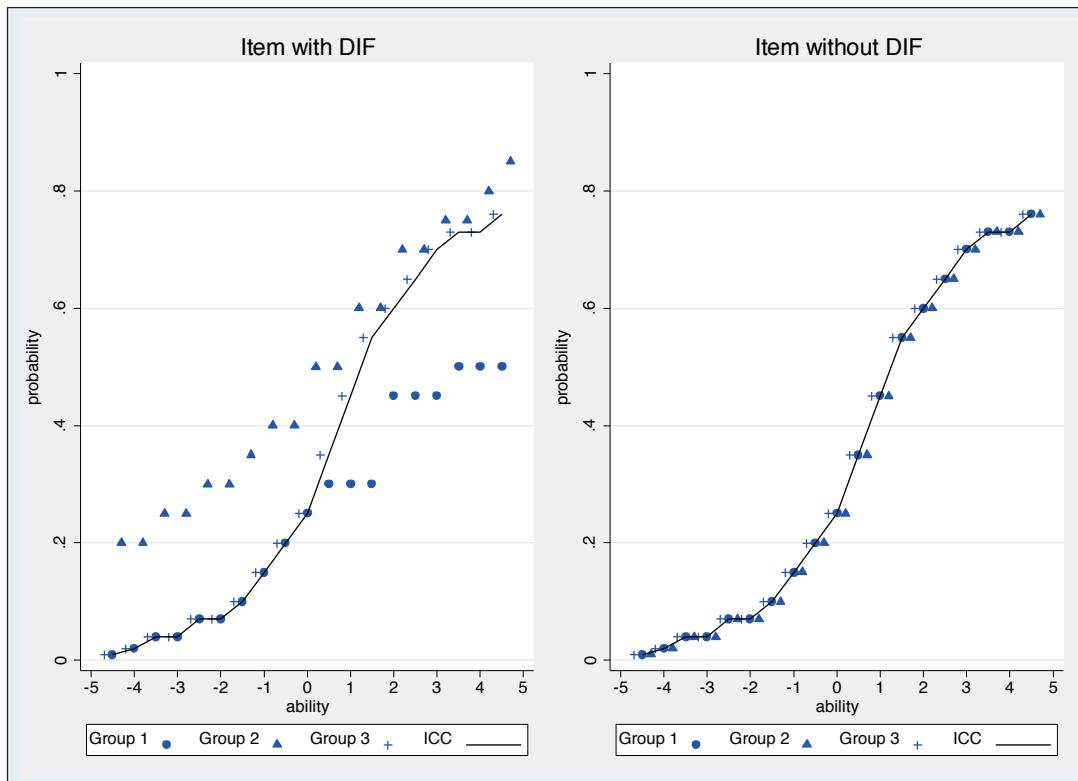


### 3.5. Differential Item Functioning (DIF)

An item is considered to have DIF if the probability of answering an item correctly differs across groups or memberships (for example, gender), controlling for level of ability (Hambleton and Swaminathan 1985; Dorans and Holland 1993; Linacre 2008). DIF analysis, however, could be sensitive to sample size since the standard errors of the item difficulty depend on the size of the groups that are being compared. Thus, large sample sizes could lead to the acceptance of even small differences between item difficulties as DIF. Therefore, it is necessary to use normalised standard errors in order to have better estimates of DIF between groups. The Educational Testing Services in the United States, as well as different scholars (Wright and Douglas 1976), suggest that for large sample sizes logit differences in item difficulty above 0.50 are signals of DIF between groups.

For this technical note, we used two approaches to check DIF. The first approach was graphically: we estimated the ICC for each item from the full sample and for each group (cohort and round). An item was considered with DIF if the ICC for any given group had a different shape than the ICC for the full sample, as seen in Figure 3.

**Figure 3.** Item characteristic curves of items with DIF and without DIF



The second approach was to calculate the Welch test using the one-parameter model, and an item was flagged with DIF if the difference between item difficulties across groups was statistically different at 5% according this test. Finally, in both analyses, we consider an item with DIF across groups if the number of children who took the item was equal or above 30.

### 3.6. Scores equating

As mentioned before, one of the main advantages of using IRT modelling is that it helps to build comparable scores using common items. Hambleton (1989) indicates that if we have different tests (with common items across them) and the items of those tests meet the IRT assumptions (good item fit indicators), then it is possible to estimate a score for each individual that is independent of the group of items that he/she answered. Thus, it is possible to use those PPVT items with adequate fit index as anchors in order to have a score that could be comparable across rounds and cohorts.

The main types of test equating are (Linacre 2008):

*Common item equating:* there are different examinees but common items across all tests forms. Two different type of analysis could be performed. First, the common and non-common items could be analysed simultaneously (for example, vertical test equating). Second, common items across all tests forms are analysed and calibrated in order to use them to adjust the mean and standard deviation of each test form.

*Common person equating:* there are different tests of the same subject (e.g. maths) but common examinees across tests. The average ability of the common examinees is used to adjust examinees' mean and standard deviations.

*Virtual item equating:* there are different examinees and different tests cover the same subject (e.g. maths). This type of equating involves identifying test pairs of items that cover the same subject and using them as pseudo-anchor items for the equating analysis.

For our analysis, we used the common item equating approach since we have the same test across cohort and rounds. It is not possible to use common person equating since having the scores of the same examinee at two different time points is similar to having different examinees.

Finally, the subsequent procedures for the equating analysis are to: (i) run the three-parameter model for the pool sample; (ii) identify those items with poor item fit, deleting them from our analysis; (iii) identify those items with DIF for all the groups, deleting them from the analysis; (iv) identify those items with the presence of DIF and consider them as different items; and (v) then run the three parameter model again using as anchor items those with the absence of DIF by round and cohort.

The item response analysis was carried out using the ado file openirt. This ado file was developed by Tristan Zajonc, who not only provided the STATA files to run the analysis but also provided technical assistance to interpret and improve the IRT analysis performed in STATA.<sup>3</sup>

### 3.7. Limitations of IRT scores

One limitation of IRT scores is that they are not comparable across languages. IRT scores are specific for each language and each scale is independent from each other. This caveat is because the PPVT test administered in Ethiopia, India and Vietnam corresponds to the English version, making it difficult to get to item cognitive equivalence across languages. Young Lives ensured the comparability of the items within each language, in order to have PPVT IRT scores comparable across rounds and age cohorts for each main language.

Therefore, we need to be careful when using IRT scores for analysis and be clear as to *what we could do and what not* with these scores. For example, we cannot use the IRT scores to compare the vocabulary level of children who took the test in Amharic to those who took it in Tigrinya. Instead, we can compare the vocabulary level between the children from the Younger and Older Cohort who took the test in Amharic or other main language. Also, we cannot use the PPVT standard scores or norms (those derived from the normalisation sample) since the population used for those norms (the United States) is completely different from the Young Lives study countries.

### 3.8. Vocabulary test approach for Round 4

The translation of the PPVT items led to changing difficulty levels of individual items and to 'disordering' the sets in Round 2 and Round 3. Given that the PPVT is one of the few cognitive longitudinal measures available in the Young Lives data, it was worthwhile to retain it, albeit with modification. With this end in mind, a subset of items which performed well in Rounds 2 and 3 were kept and administered again in Round 4 to the Younger Cohort, with the additional change that the rules to set the basal and ceiling item were also abandoned.<sup>4</sup>

---

3 See Appendix A for the steps followed in STATA to run the IRT analysis.

4 The basal set rule is one error, or no errors, in a set of 12 items, and the ceiling set rule is eight or more errors in a set of 12 items.

The criteria followed to select the subset of items within each country was: (i) IRT scores (3PL) were calculated for PPVT in Ethiopia (Amharic, Oromifa, and Tigrinya), India (Telugu) and Vietnam (Vietnamese), restricting the sample to the main languages within each country; (ii) ability cut-offs were identified that would split the full sample (combined Older Cohort and Younger Cohort, Rounds 2 and 3) into four equal bands of ability; (iii) items were sorted within these four bands based on their difficulty parameters; (iv) further, the items were sorted within these four bands based on their discrimination; (v) finally, roughly equal number of items across bins were selected, inspecting the item characteristic curves: specifically, items with bad fit, high guessing parameter, or zero variation were excluded.

## 4. Results

We estimated the three parameter IRT analysis for the pool sample for each of the main languages in the three countries: Amharic, Tigrinya and Oromifa, Telugu, and Vietnamese. This first analysis allows us to identify those items with poor fit that have to be dropped from each of the composite scores. Table 1 shows the percentage of items that were dropped because of poor item fit or DIF for all the comparison groups (round and cohort). The percentage of items dropped, on average, was around one third of the total items,<sup>5</sup> Oromifa being the language with the highest percentage (36%) of items dropped.

**Table 1.** Number of items dropped by language

Language	Total items	Items dropped	%
Amharic	204	32	16%
Tigrinya	204	53	26%
Oromifa	204	74	36%
Telugu	204	48	24%
Vietnamese	204	42	21%

Source: Young Lives, Main Survey Rounds 2, 3 and 4

However, one of the main concerns was the number of common items or anchor items dropped between Round 4 and the previous rounds. Table 2 shows that the percentage of anchor items dropped was less than 10%. These results indicate that we have enough anchor items to ensure an adequate equating across rounds and cohorts. Finally, those items that have a good item fit but have DIF were split and considered as a different item.<sup>6</sup>

**Table 2.** Number of anchor items dropped by language

Language	Total anchor items	Anchor items dropped	%
Amharic	55	0	0%
Tigrinya	55	4	7%
Oromifa	55	2	4%
Telugu	57	0	0%
Vietnamese	76	1	1%

Source: Young Lives, Main Survey Rounds 2, 3 and 4

5 See Appendix B and Appendix C for details of the ICC curves and item DIF analysis for all the country analyses performed.

6 See Appendix D for details of the items dropped and flagged with DIF.

Once items with poor fit and DIF for all groups were dropped, we ran the three parameter model again in order to get corrected IRT scores for each set of children. Table 3 shows the average mean scores for all the languages by cohort and round. IRT scores, for both cohorts, increase over time.

**Table 3.** Mean scores by language for each round and age cohort (standard deviation)

Language	Older Cohort		Younger Cohort		
	R2	R3	R2	R3	R4
Amharic	2.2 (0.97)	2.7 (1.16)	0.0 (1.00)	1.3 (1.35)	2.1 (1.29)
Oromifa	2.3 (1.12)	2.7 (1.16)	0.0 (1.00)	0.9 (1.01)	2.7 (1.16)
Tigrinya	2.8 (1.01)	3.4 (1.22)	0.0 (1.00)	1.5 (1.08)	2.4 (1.20)
Telugu	2.4 (1.05)	2.6 (0.99)	0.0 (1.00)	0.8 (0.96)	1.9 (0.98)
Vietnamese	3.4 (1.26)	3.6 (1.15)	0.0 (1.00)	1.7 (0.85)	3.0 (0.91)

Source: Young Lives, Main Survey Rounds 2, 3 and 4

Table 4 shows the increment in the IRT scores over time by cohort and language. We found that all the increments are statistically significant for both cohorts and languages; also, in the Younger Cohort, since we have three time points, we could estimate the increments between Rounds 2 and 3, and 3 and 4. Our results show that Amharic, Tigrinya and Vietnamese children have the highest increment between Rounds 2 and 3, while for Oromifa and Telugu children, the highest increment was between Rounds 3 and 4.

**Table 4.** Gap analysis for each age cohort

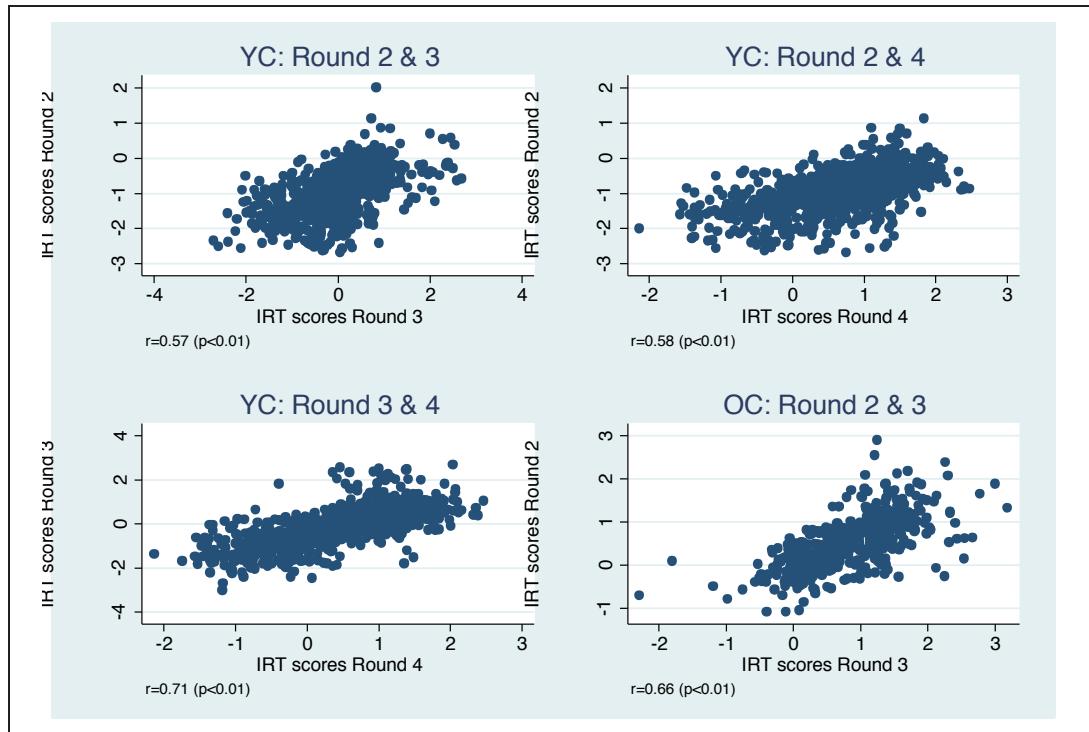
Language	Older Cohort		Younger Cohort	
	R3 - R2	R3 - R2	R4 - R3	R4 - R3
Amharic	0.42*	1.26*	0.88*	0.88*
Oromifa	0.36*	0.89*	1.85*	1.85*
Tigrinya	0.64*	1.47*	0.97*	0.97*
Telugu	0.21*	0.77*	1.17*	1.17*
Vietnamese	0.14*	1.71*	1.26*	1.26*

Notes: \* Mean scores differences between rounds are statistically significant at 5% according to the ttest for dependent or correlated samples.

Source: Young Lives, Main Survey Rounds 2, 3 and 4

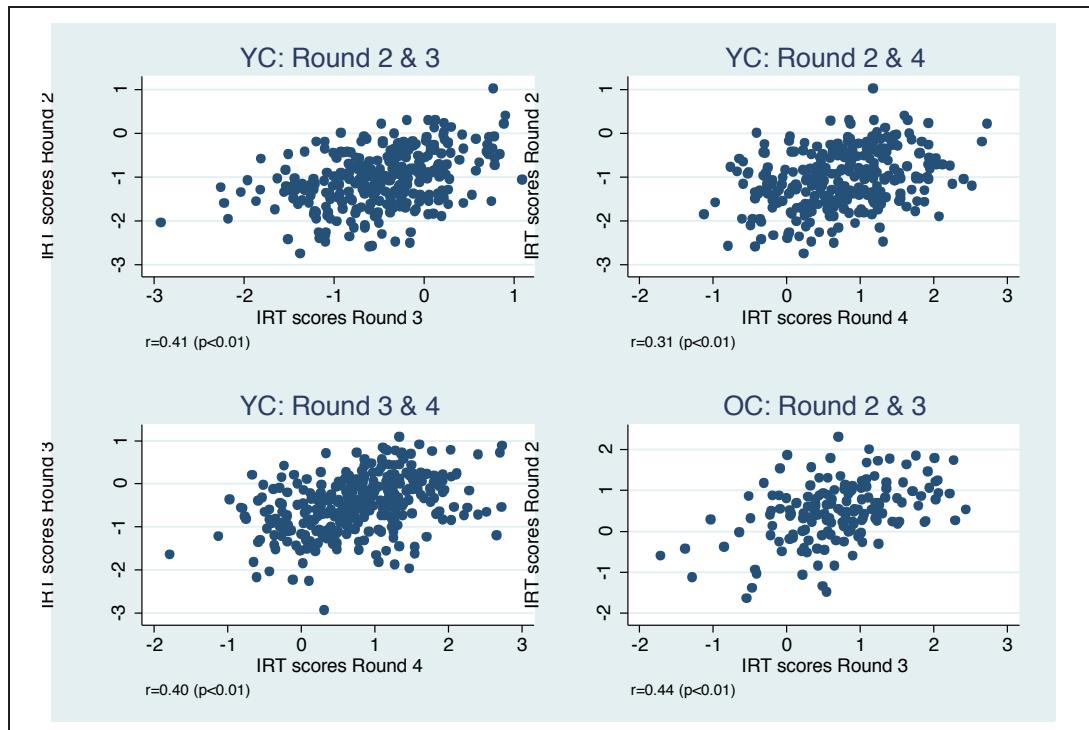
Figures 4 to 8 show the correlation of IRT scores between rounds for each language. Amharic, Telugu and Vietnam have the highest correlations, while Oromifa and Tigrinya show the lowest correlations between rounds for both age cohorts.

**Figure 4.** Scatterplots for IRT scores between rounds for Younger and Older Cohort – Amharic



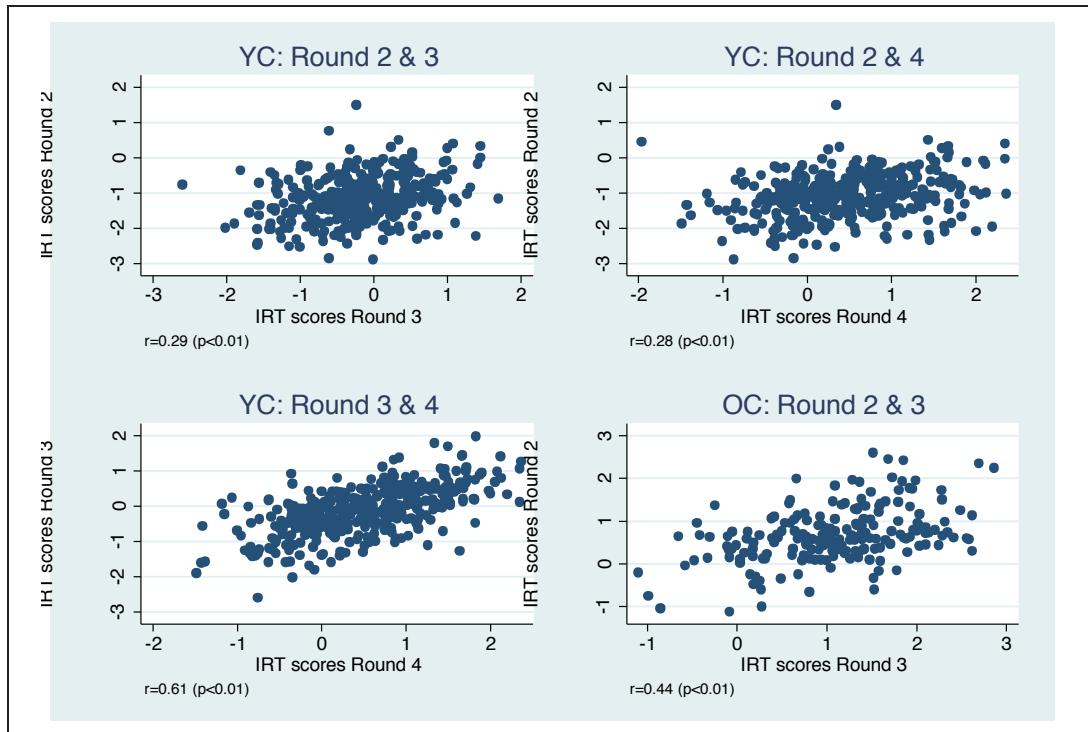
Source: Young Lives, Main Survey Rounds 2, 3 and 4

**Figure 5.** Scatterplots for IRT scores between rounds for Younger and Older Cohort – Oromifa



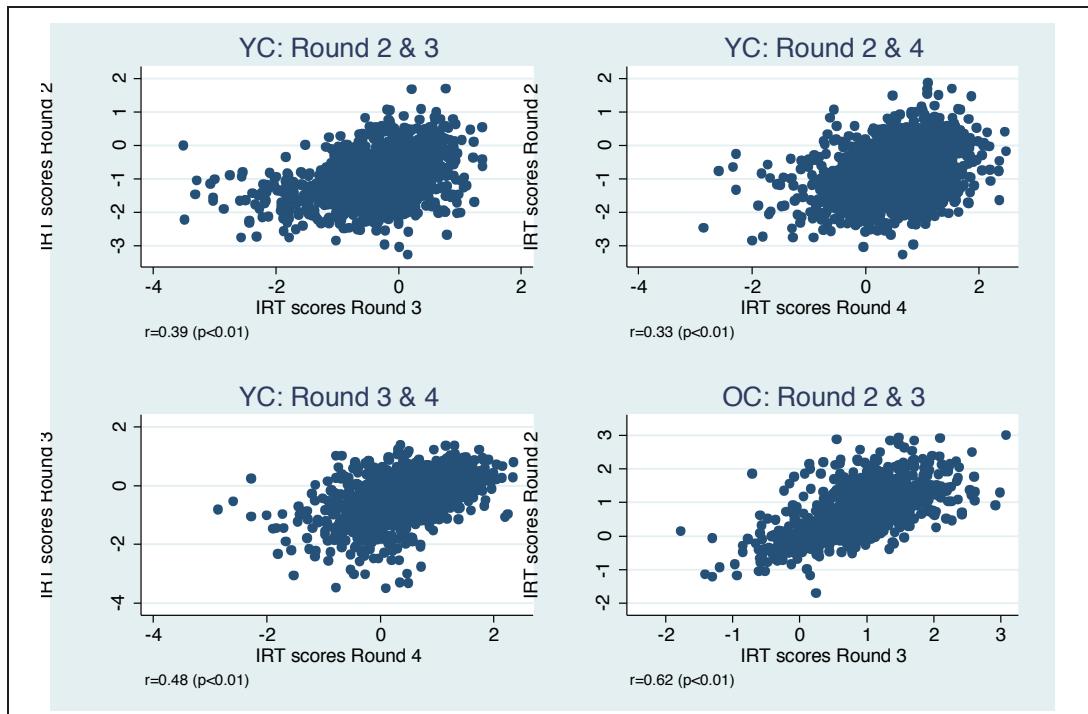
Source: Young Lives, Main Survey Rounds 2, 3 and 4

**Figure 6.** Scatterplots for IRT scores between rounds for Younger and Older Cohort – Tigrinya



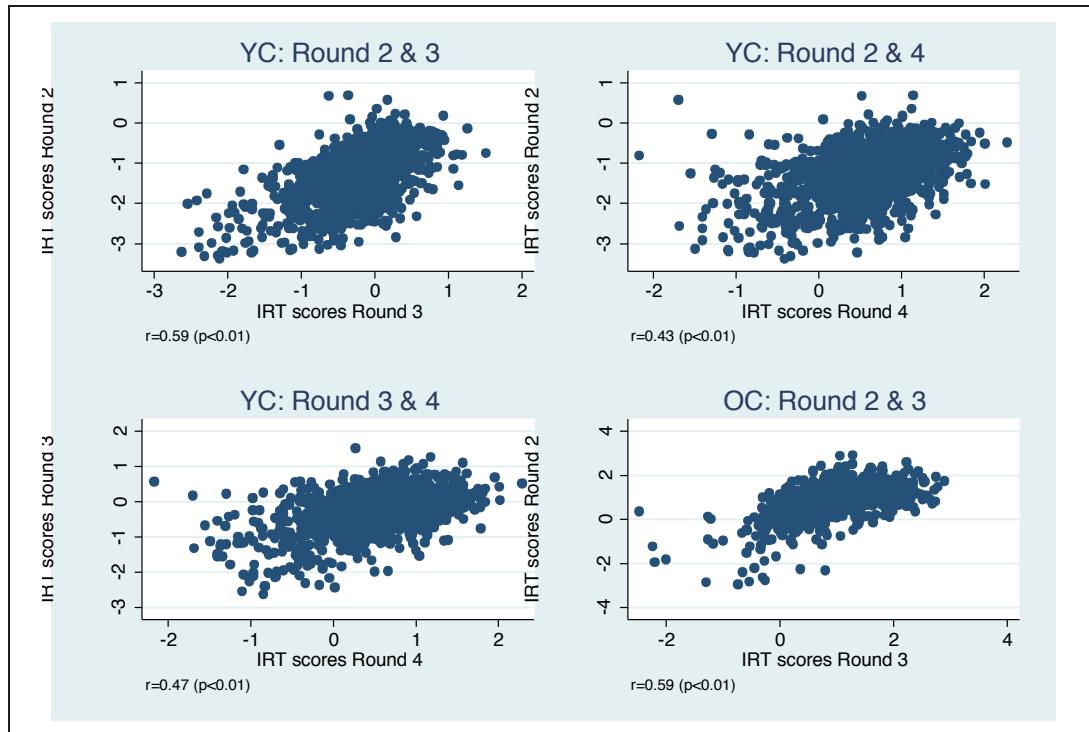
Source: Young Lives, Main Survey Rounds 2, 3 and 4

**Figure 7.** Scatterplots for IRT scores between rounds for Younger and Older Cohort – Telugu



Source: Young Lives, Main Survey Rounds 2, 3 and 4

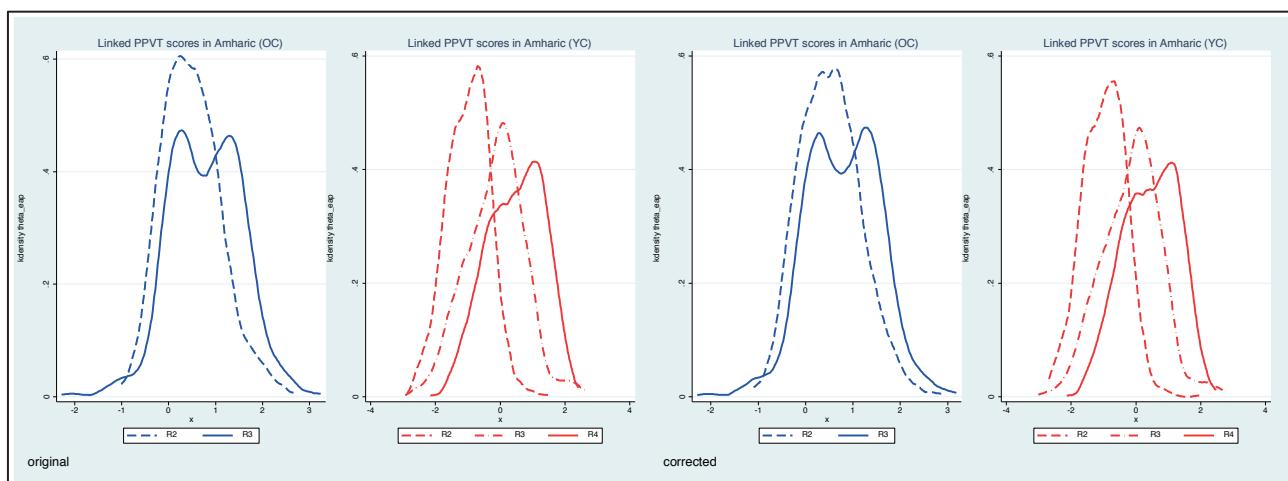
**Figure 8.** Scatterplots for IRT scores between rounds for Younger and Older Cohort – Vietnam



Source: Young Lives, Main Survey Rounds 2, 3 and 4

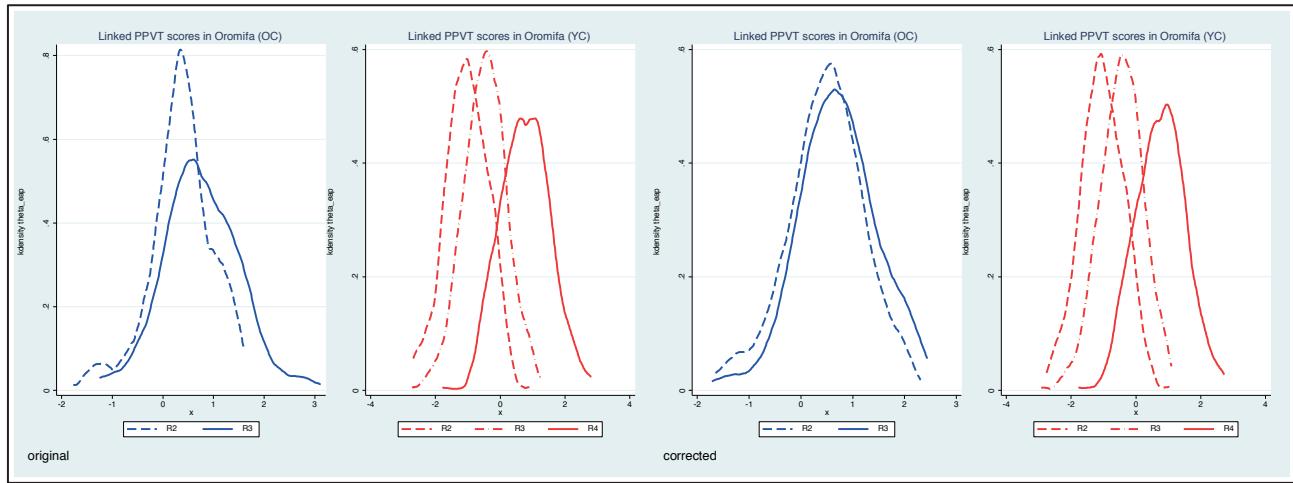
Then, we compared the distribution of the original and corrected scores. Figures 9 to 13 show that average scores for Younger Cohort children increase over time for all the main languages analysed, while the scores for the Older Cohort show some stagnation for Telugu and Vietnamese children; therefore, the average scores are fairly similar across rounds, confirming possible ceiling effects.

**Figure 9.** Original and corrected IRT scores for Amharic



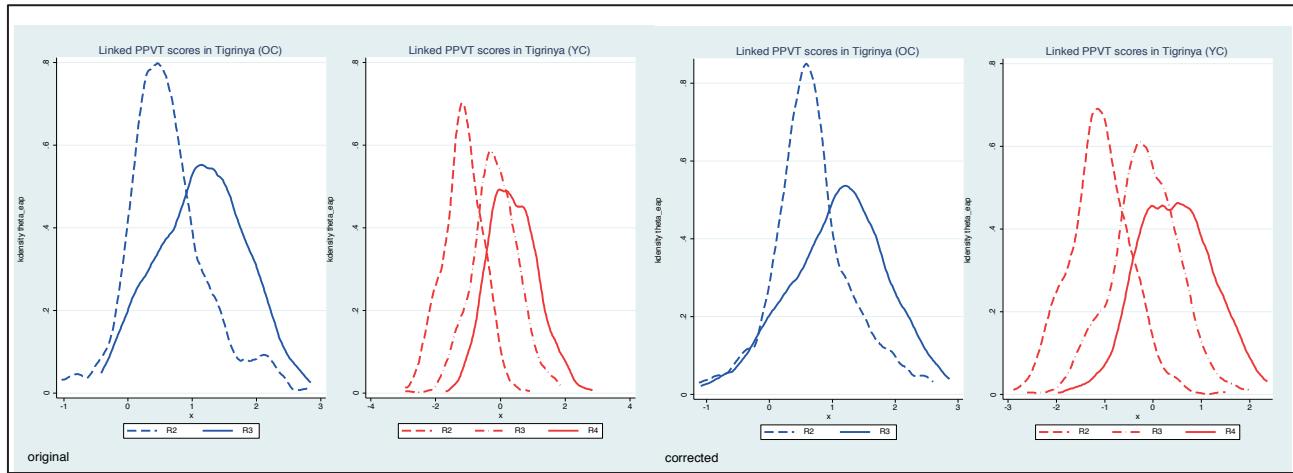
Source: Young Lives, Main Survey Rounds 2, 3 and 4

**Figure 10.** Original and corrected IRT scores for Oromifa



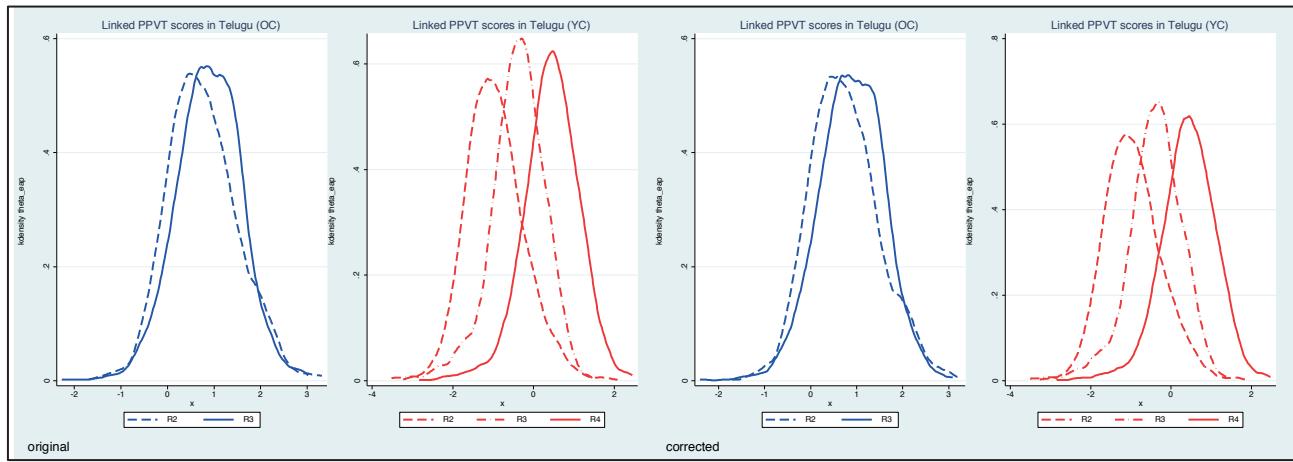
Source: Young Lives, Main Survey Rounds 2, 3 and 4

**Figure 11.** Original and corrected IRT scores for Tigrinya



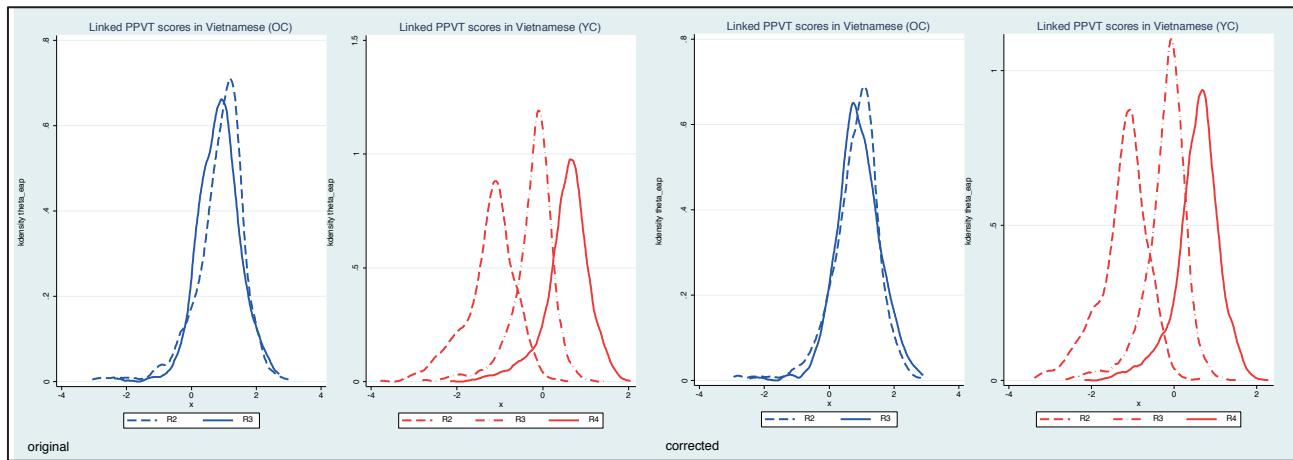
Source: Young Lives, Main Survey Rounds 2, 3 and 4

**Figure 12.** Original and corrected IRT scores for Telugu



Source: Young Lives, Main Survey Rounds 2, 3 and 4

**Figure 13.** Original and corrected IRT scores for Vietnamese



Source: Young Lives, Main Survey Rounds 2, 3 and 4

Finally, we equated the Round 4 PPVT scores of the siblings with the scores of the young lives children in order to have comparable measures. The analysis showed an adequate anchoring (see Appendix F) since none of the items administered to the siblings in Ethiopia (Amharic, Oromifa, and Tigrinya) and Vietnam showed a poor fit.<sup>7</sup>

## 5. Final remarks

This technical note gives details of the procedures followed to equate the PPVT scores for the main languages in three of the four Young Lives study countries. The main results are:

- Results confirm that the new approach followed in Round 4 was adequate, as most of the items selected had a good item fit and did not show DIF by groups (cohort or round) for all the main languages in each country.
- For some languages, such as Tigrigna and Oromifa, the number of items deleted by poor fit or DIF across all groups (rounds and cohort) was significant, with almost one third of items dropped from the final scale.
- It was possible to ensure an adequate equating of the PPVT scores for all the main languages across rounds and cohorts. Our results show that PPVT scores for the Younger and Older Cohort increase over time for all the main languages, and these increments are statistically significant.
- Results from the Younger Cohort show a curvilinear (convex) trend in the vocabulary acquisition for Amharic, Tigrigna, and Vietnamese children since the highest increment was between Rounds 2 and 3; while an exponential growth was observed for Oromifa and Telugu children as the scores increment is higher between Rounds 3 and 4.
- Vocabulary knowledge decreases over time across cohorts for all languages, as PPVT scores for the Older Cohort in Round 2 (at 12 years old) are higher than Younger Cohort children in Round 4 (also at 12 years old).

<sup>7</sup> Appendix G provides details of the item parameters used for the test equating.

# References

- Campbell, J.M. (1998) 'Review of the Peabody Picture Vocabulary Test – Third Edition', *Journal of Psychoeducational Assessment* 16.4: 334–8.
- Campbell, J.M., S.K. Bell and L.K. Keith (2001) 'Concurrent Validity of the Peabody Picture Vocabulary Test – Third Edition as an Intelligence and Achievement Screener for Low SES African American Children', *Assessment* 8.1: 85–94.
- Cueto, S., and J. Leon (2013) *Psychometric Characteristics of Cognitive Development and Achievement Instruments in Round 3 of Young Lives*, Technical Note 25, Oxford: Young Lives
- Cueto, S., J. Leon, G. Guerrero and I. Munoz (2009) *Psychometric Characteristics of Cognitive Development and Achievement Instruments in Round 2 of Young Lives*, Technical Note 15, Oxford: Young Lives
- Dorans, N.J., and P.W. Holland (1993) 'DIF detection and description: Mantel-Haenzel and standardization', In P.W. Holland and H. Wainer (eds.) *Differential Item Functioning*, 35–66, Hillsdale, NJ: Lawrence Erlbaum Associates.
- Dunn, L. and L. Dunn (1997) *Examiner's Manual for the PPVT-III. Form IIIA and IIIB*, Minnesota: AGS.
- Gray, S., E. Plante, R. Vance and M. Henrichsen (1999) 'The Diagnostic Accuracy of Four Vocabulary Tests Administered to Preschool-Age Children', *Language, Speech and Hearing Services in Schools* 30: 196–206.
- Hambleton, R.K. (1989) 'Principles and selected applications of item response theory', In R.L. Linn (ed.) *Educational Measurement*, 147-200, New York: Macmillan.
- Hambleton, R.K., and H. Swaminathan (1985) *Item Response Theory: Principle and Applications*, Boston, MA: Kluwer Nijhoff.
- Linacre, J.M. (2008) 'Winsteps: A Rasch analysis comptuer program' [Version 3.68], Chicago, IL. ([www.winsteps.com](http://www.winsteps.com))
- Wright B.D., and G.A. Douglas (1976) 'Rasch item analysis by hand', MESA Research Memorandum Number 21, Statistical Laboratory, Department of Education, Chicago: University of Chicago.

# Appendices

## Appendix A. Details of the STATA analysis performed

Appendix A presents details of the steps followed to obtain the PPVT equated scores for Young Lives. The steps were:

*Step 1: Appending data.* Our first step was to append the PPVT data (child item answers) for each country (Ethiopia, India and Vietnam) from the last three rounds for the Younger Cohort and Rounds 2 and 3 for the Older Cohort. The objective was to increase the sample size for the IRT analysis.

*Step 2: Generating datasets by language.* Once the data was appended for each country, we used the variable related with the *language test administration* to filter and select cases who took the test in the same language within each country. The languages selected by country were:

**Table A1.** Main languages chosen by country

Ethiopia	India	Vietnam
Amharic	Telugu	Viet
Oromifa		
Tigrinya		

Thus, we had a dataset for each of the languages above. These datasets comprised PPVT information for both children cohorts from Round 2 to 4.

*Step 3: Item scoring.* With the dataset for each language, we scored the items using a dummy coding, where the item took the value of 1 if the children answered the item correctly, and 0 otherwise. The following is an example for this item scoring using STATA:

### Item scoring

\* Generating answer codes for all items ( 1 to 204 )

```

gen ans_1=4
gen ans_2=3
gen ans_3=1
...
gen ans_201=4
gen ans_202=2
gen ans_203=3
gen ans_204=2

mvdecode ppvt1 - ppvt204, mv(77 79 88 99)

forval i = 1/204 {
    replace ppvt`i' = 0 if ppvt`i'!=ans_`i' & ppvt`i'!=.
    replace ppvt`i' = 1 if ppvt`i'==ans_`i'
}

```

The answer keys for the 204 items can be obtained from the child questionnaires available online.

*Step 4: IRT pool analysis.* Once we scored all the items, we ran the IRT pool analysis for each language using the openirt ado file. This IRT analysis was used to check for item fit and item bias by round and age cohort. The following is an example for each analysis performed in this step.

#### IRT analysis for Vietnam

##### STATA code

```
openirt, id(id) save_item_parameters("viet_items_ppvt.dta") save_trait_parameters("viet_traits_ppvt.dta")
item_prefix("ppvt") model("3PL") samplesize(500) burnin(500)
```

viet\_items\_ppvt.dta : this dataset will comprise the item difficulty, discrimination and guessing parameters.  
viet\_traits\_ppvt.dta : this dataset will comprise the child ability measures

#### Item fit analysis (ICC graphs) per item for Vietnam

##### STATA code

```
use "$output\ppv_full_viet.dta", clear
drop ans*
merge 1:1 id using "$output\viet_traits_ppvt_p.dta"

drop _merge
xtile perc_theta=theta_pv1,nq(10)
sort perc
by perc: egen mean_theta=mean(theta_pv1)

forval i=1/204{
    by perc: egen mean_ppvt`i'=mean(ppvt`i')
}

keep perc mean*

forval i=1/204{
ren mean_ppvt`i' id`i'
}

duplicates drop perc, force
reshape long id,i(perc) j(item)
ren id prop
ren ite id

merge m:1 id using "$output\viet_items_ppvt_p.dta"
drop _merge

cd "$graphs\
forval i=1/204 {
local j=(`i' - 1)*10+ 1
twoway (scatter prop mean_theta if id==`i', sort)(function c_pv1[`j']+ (1-c_pv1[`j'])/(1+exp(-1.7*a_pv1[`j']*(x-
b_pv1[`j'])))) if id==`i', range(-4 4)), ///
xtitle("Theta") ytitle("P(X=1|Theta)") title("Item `i'") legend(off)
graph save tmp`i', replace
}

ppv_full_viet.dta : It is the pool dataset for vietnam
```

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

**Item DIF analysis (ICC graphs by group per item) for Vietnam**

STATA code

```

use "$output\ppvt_full_viet.dta", clear
drop ans*
merge 1:1 id using "$output\viet_traits_ppvt.dta", nogen

xtile perc_thet=theta_pv1,nq(10)
sort perc
by perc: egen mean_theta=mean(theta_pv1)

forval i=1/204{
    by perc: egen mean_ppvt`i'=mean(ppvt`i')
}
levels of group, local(group)

forval i=1/204{
foreach n of local group{
    by perc: egen mean_item`i'_`n'= mean(ppvt`i') if group==`n'
}
}

keep perc mean* group

forval i=1/204{
ren mean_ppvt`i' id`i'
}

duplicates drop perc group, force
egen group_perc=group(group perc)
reshape long id mean_item,i(group_perc) j(item )
ren id prop
ren ite id

merge m:1 id using "$output\viet_items_ppvt.dta", nogen
sort id group perc_thet

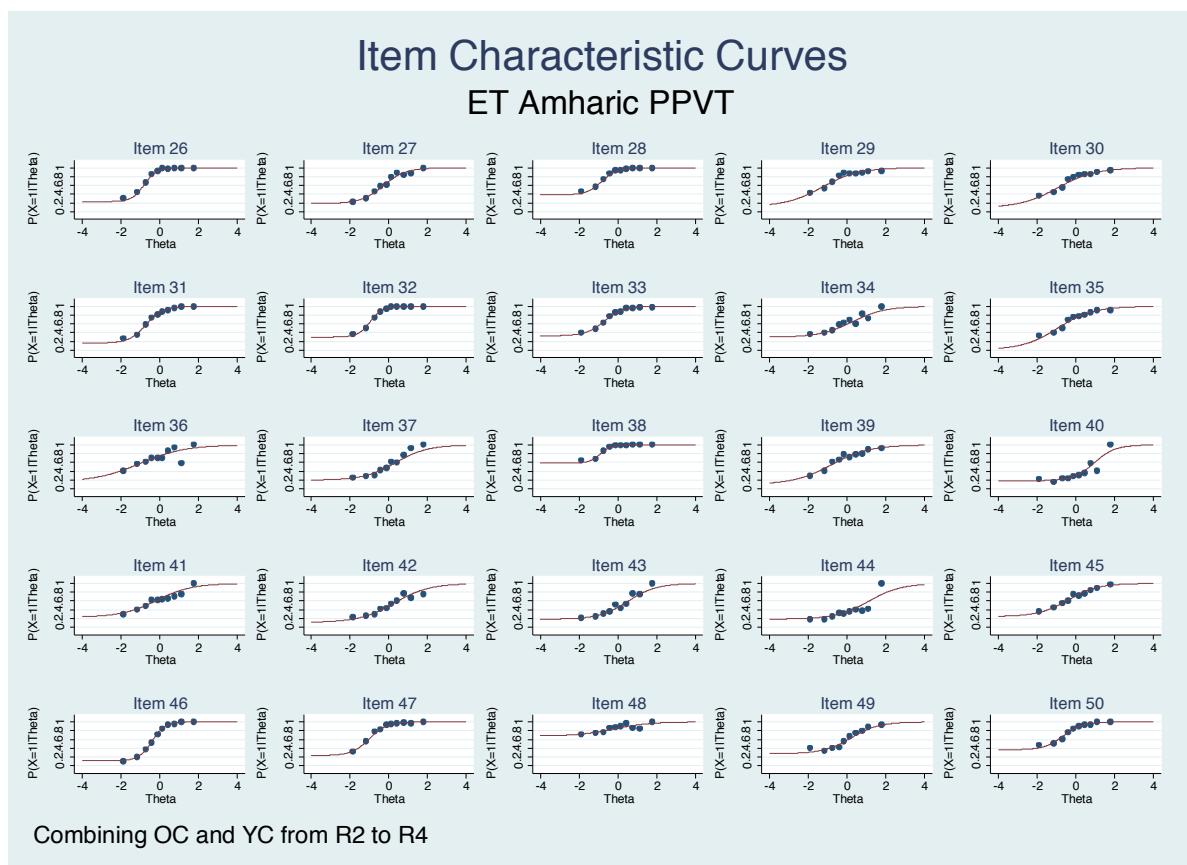
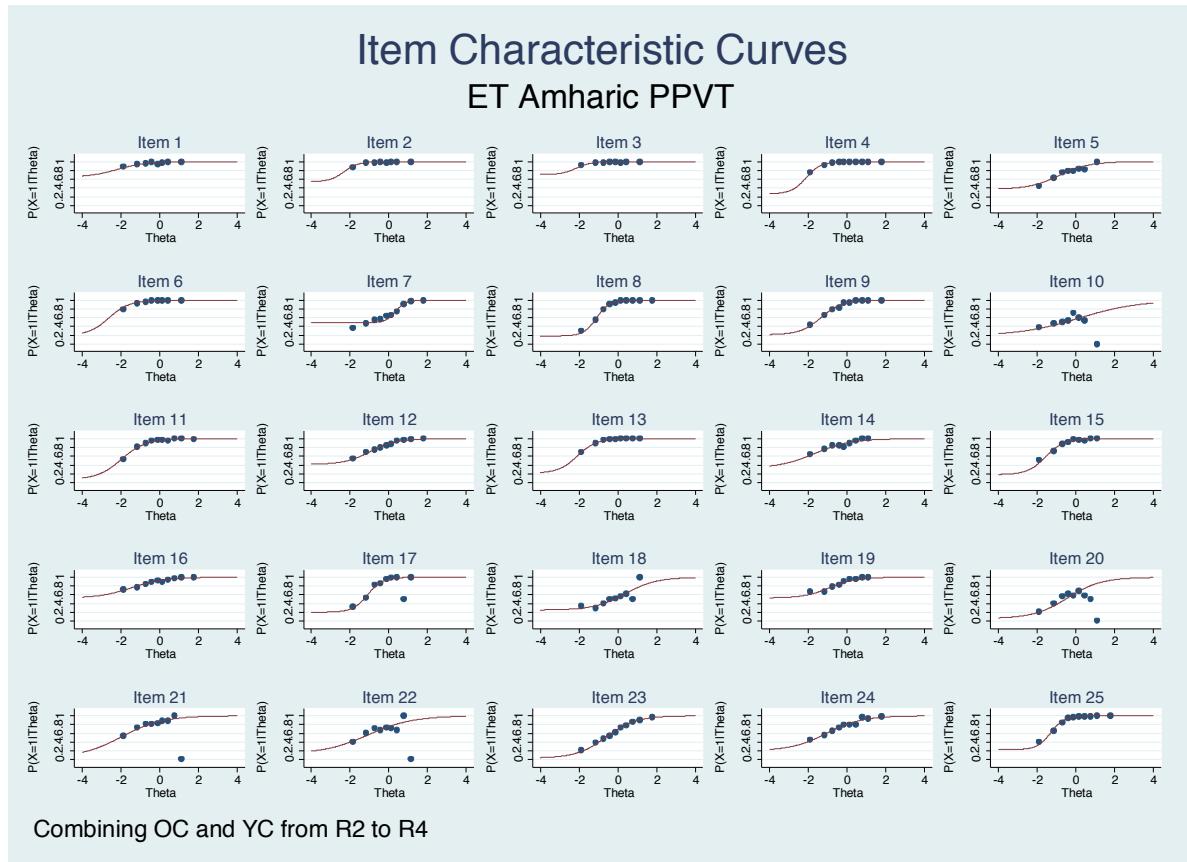
/* group:
1 R2OC
2 R3OC
3 R2YC
4 R3YC
5 R4YC
*/
forval i=1/204 {
local j=(`i' - 1)*48 + 1
twoway (scatter mean_item`i'_1 mean_theta if id=='`i' & group==1, sort msymbol(plus) mcolor(red))/*
*/(scatter mean_item`i'_2 mean_theta if id=='`i' & group==2, sort msymbol(triangle) mcolor(blue))/*
*/(scatter mean_item`i'_3 mean_theta if id=='`i' & group==3, sort msymbol(circle) mcolor(black))/*
*/(scatter mean_item`i'_4 mean_theta if id=='`i' & group==4, sort msymbol(lgx) mcolor(dknavy))/*
*/(scatter mean_item`i'_5 mean_theta if id=='`i' & group==5, sort msymbol(square) mcolor(green))/*
*/(function c_pv1[`j'] + (1-c_pv1[`j'])/(1+exp(-1.7*a_pv1[`j`]*(`x-b_pv1[`j`])))) if id=='`i', range(-4 4), ///
xtitle("Theta") ytitle("P(X=1|Theta)") title("Item `i'") legend(off)
graph save tmp`i', replace
}

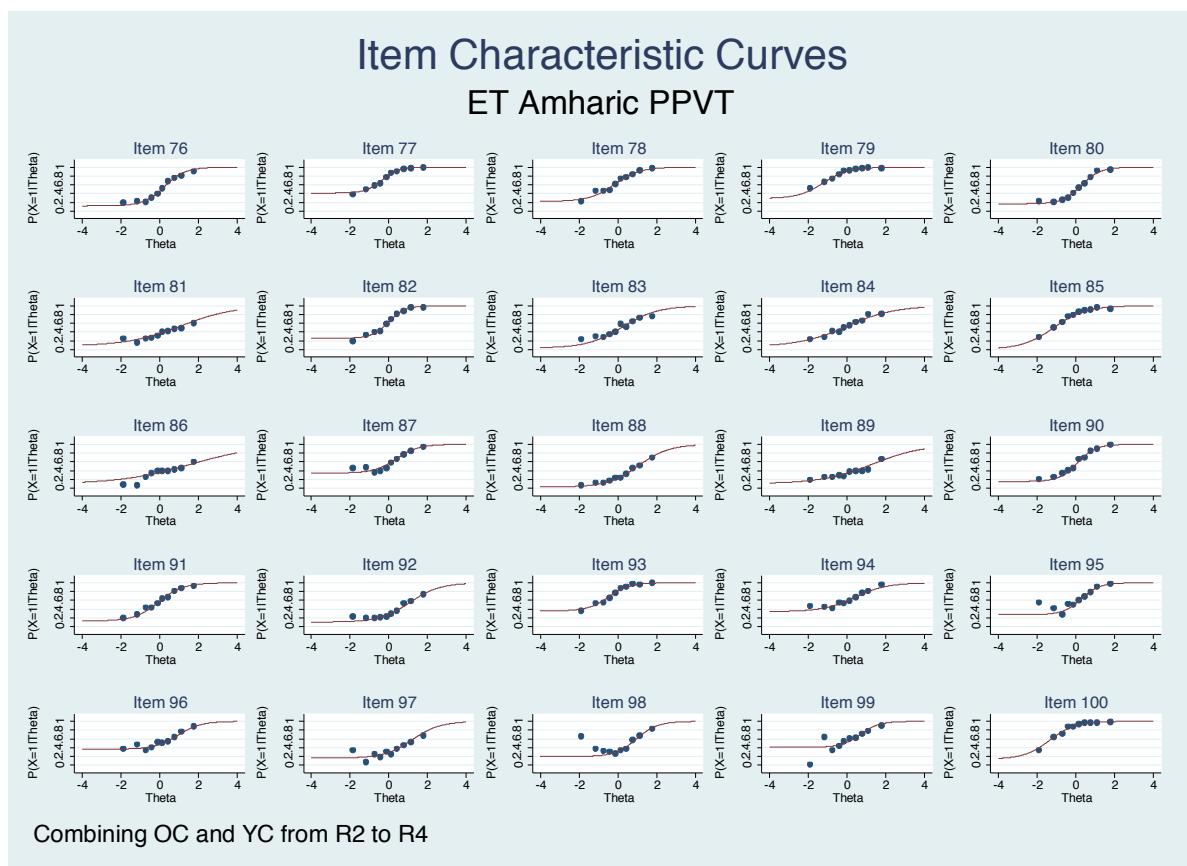
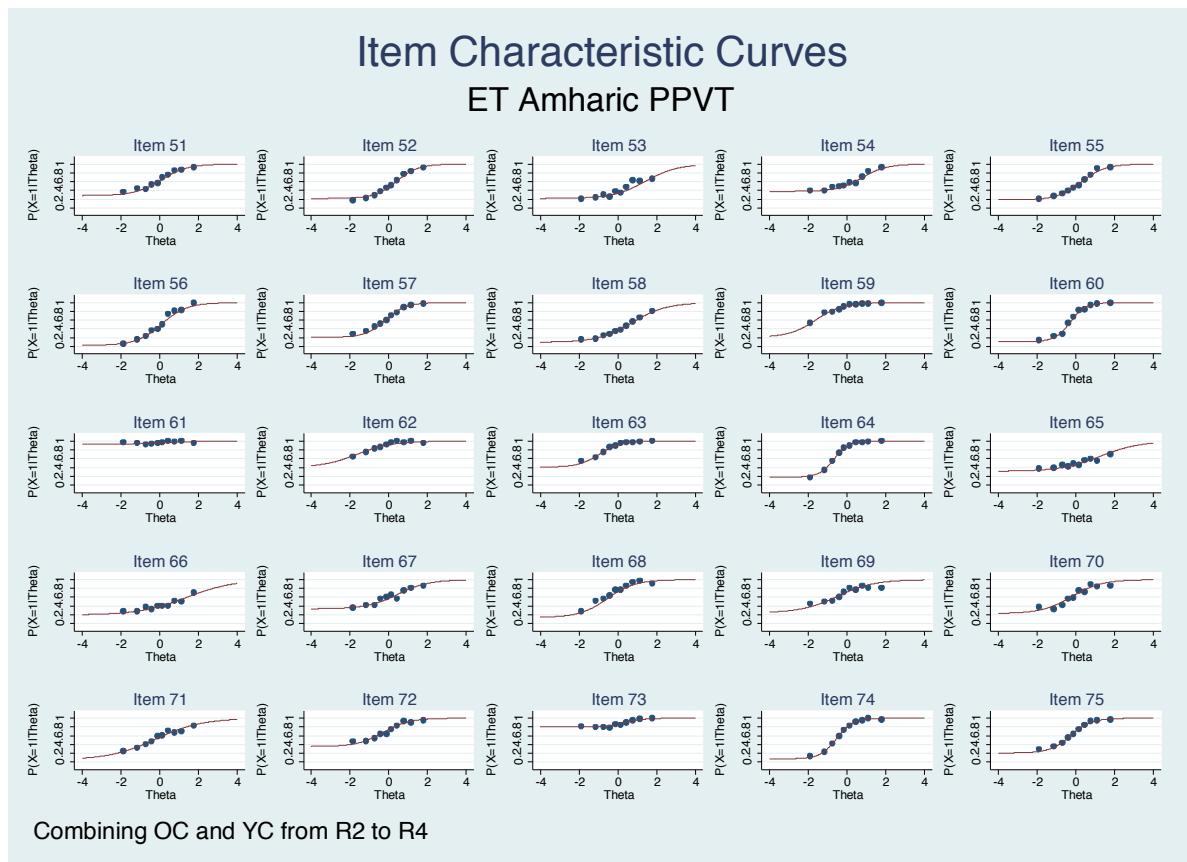
```

*Step 5: Item dropping and splitting.* Items with poor item fit (ICC) were dropped from the dataset, as well as those items that showed a different ICC for all groups under analysis. For items that showed a different ICC for one or two groups, a new item was generated (splitting) for each group with different behaviour, and children's answers for that group were deleted from the original item.

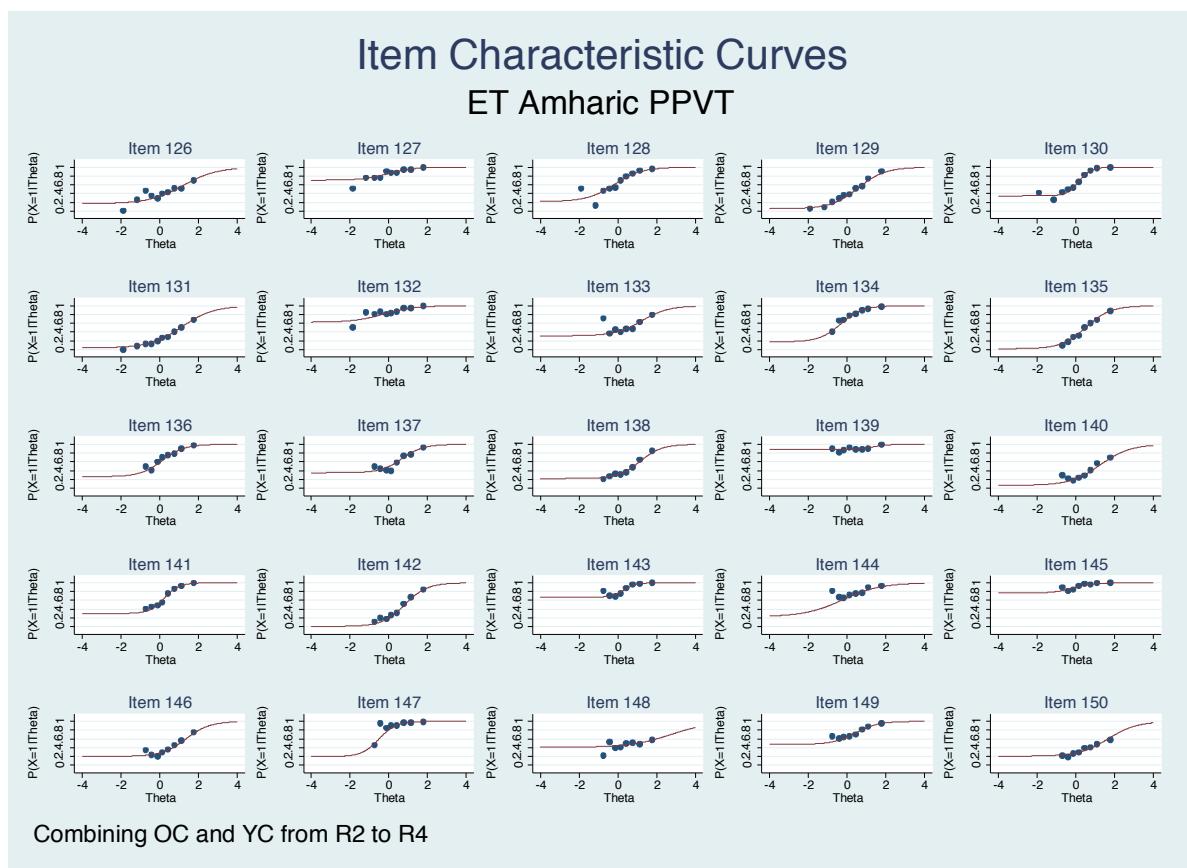
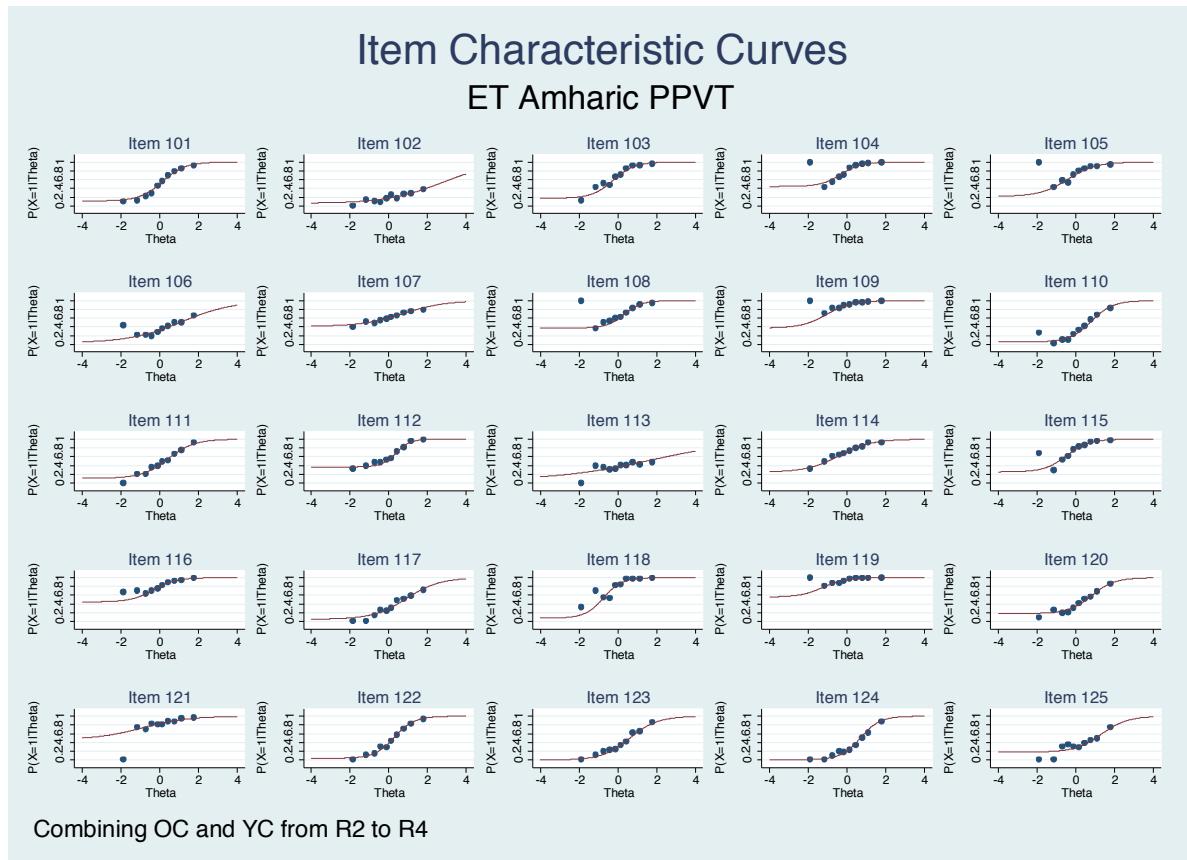
*Step 6: IRT pool analysis corrected.* Once all the items were corrected using the first IRT analysis and we kept the items with good fit and without DIF in each dataset; we reran the IRT analysis in order to get the corrected IRT scores. The codes used for this new analysis are similar to those followed in Step 3.

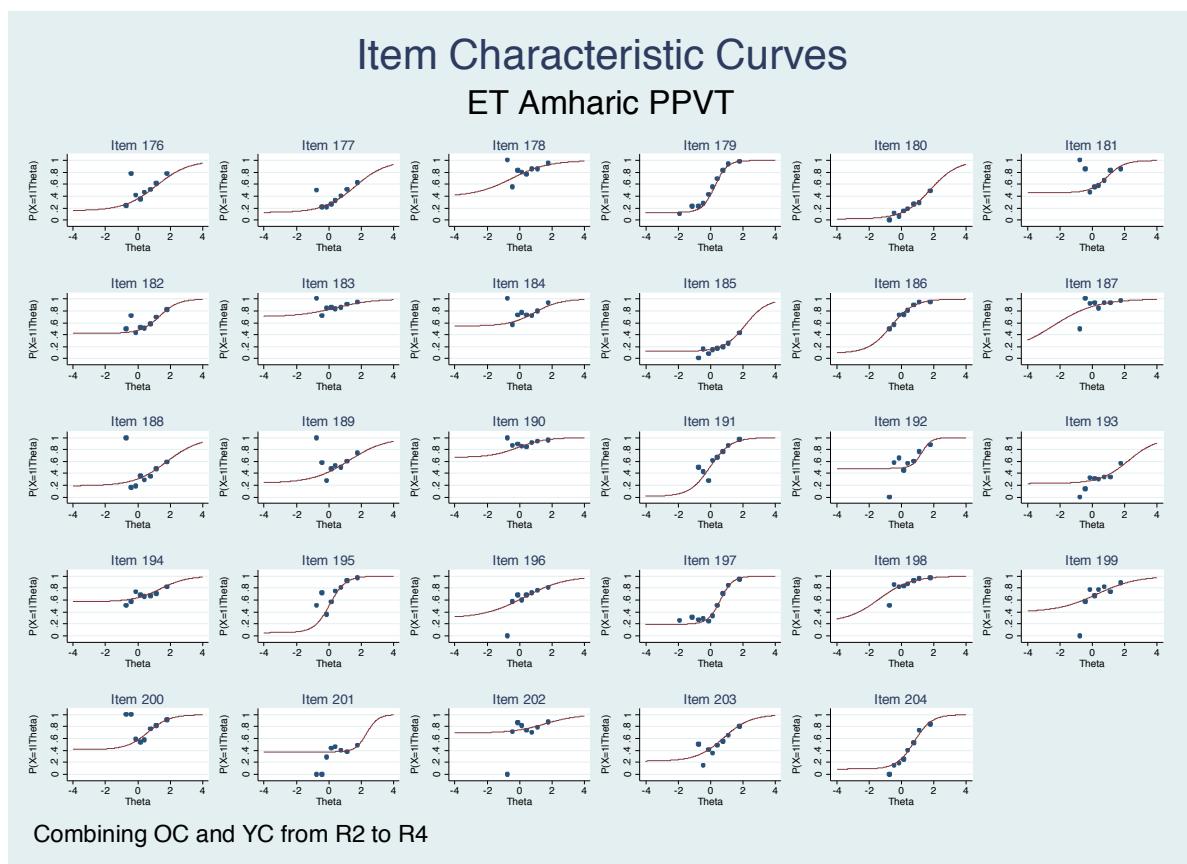
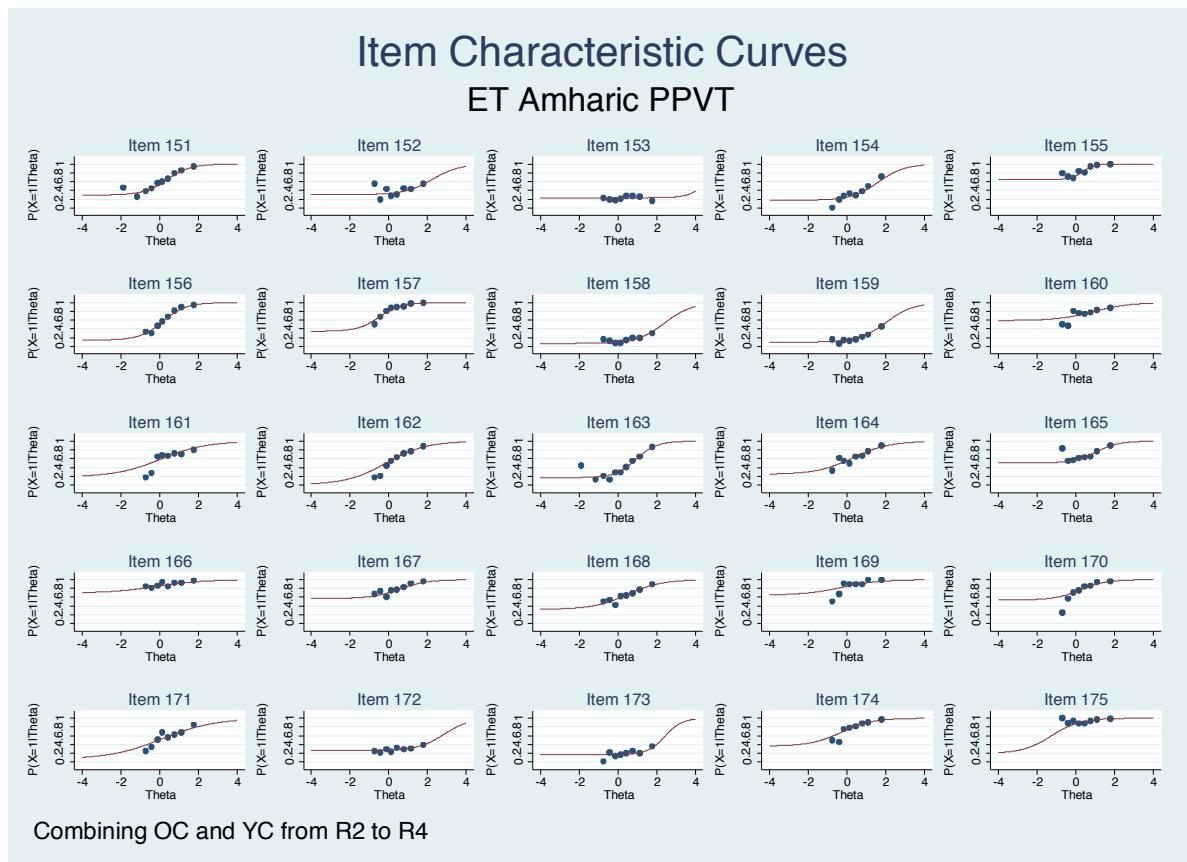
## Appendix B. ICC for each item by country and main language



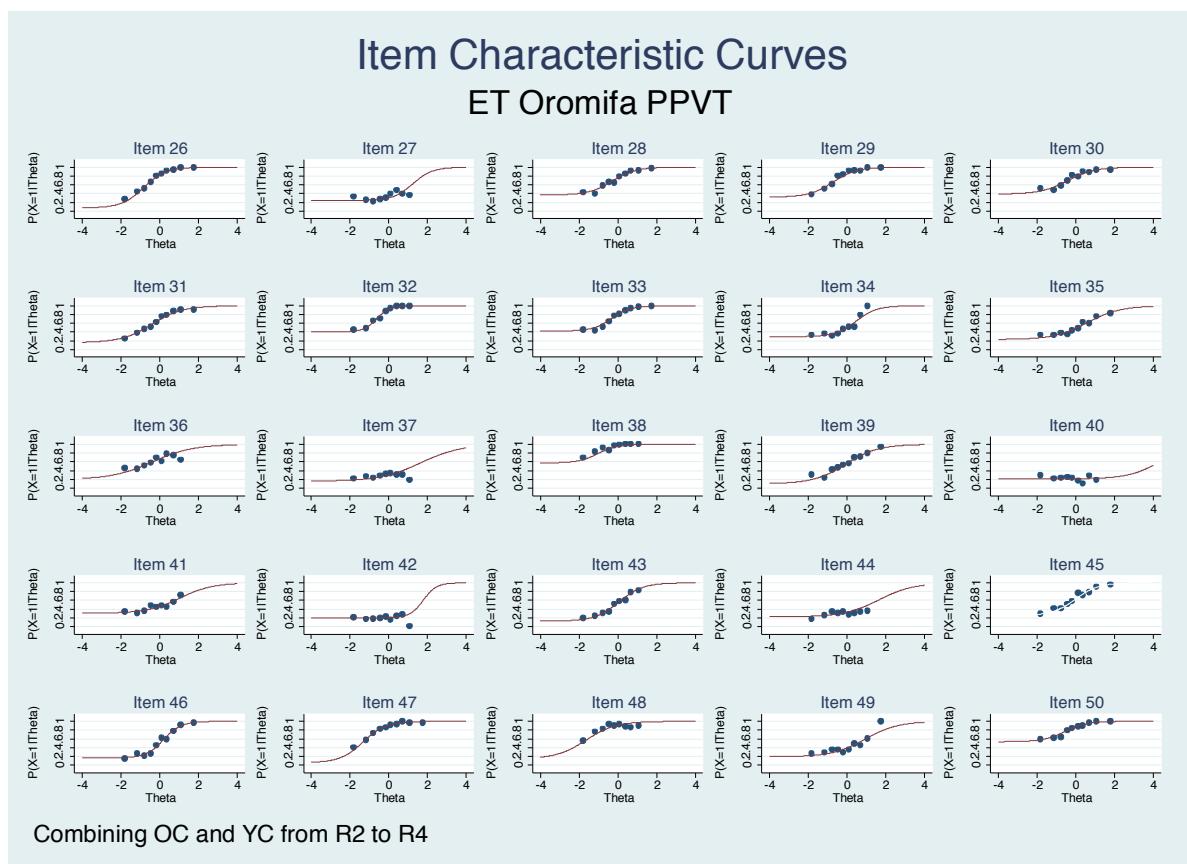
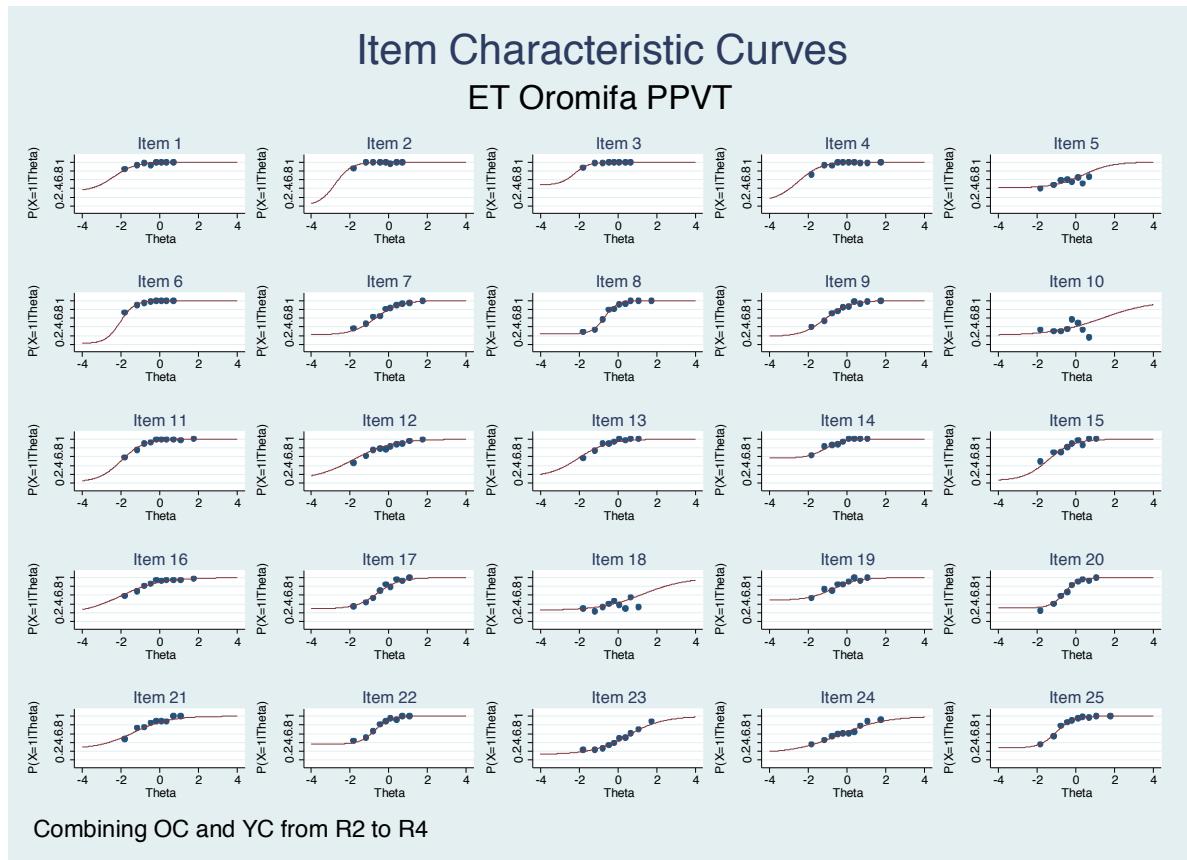


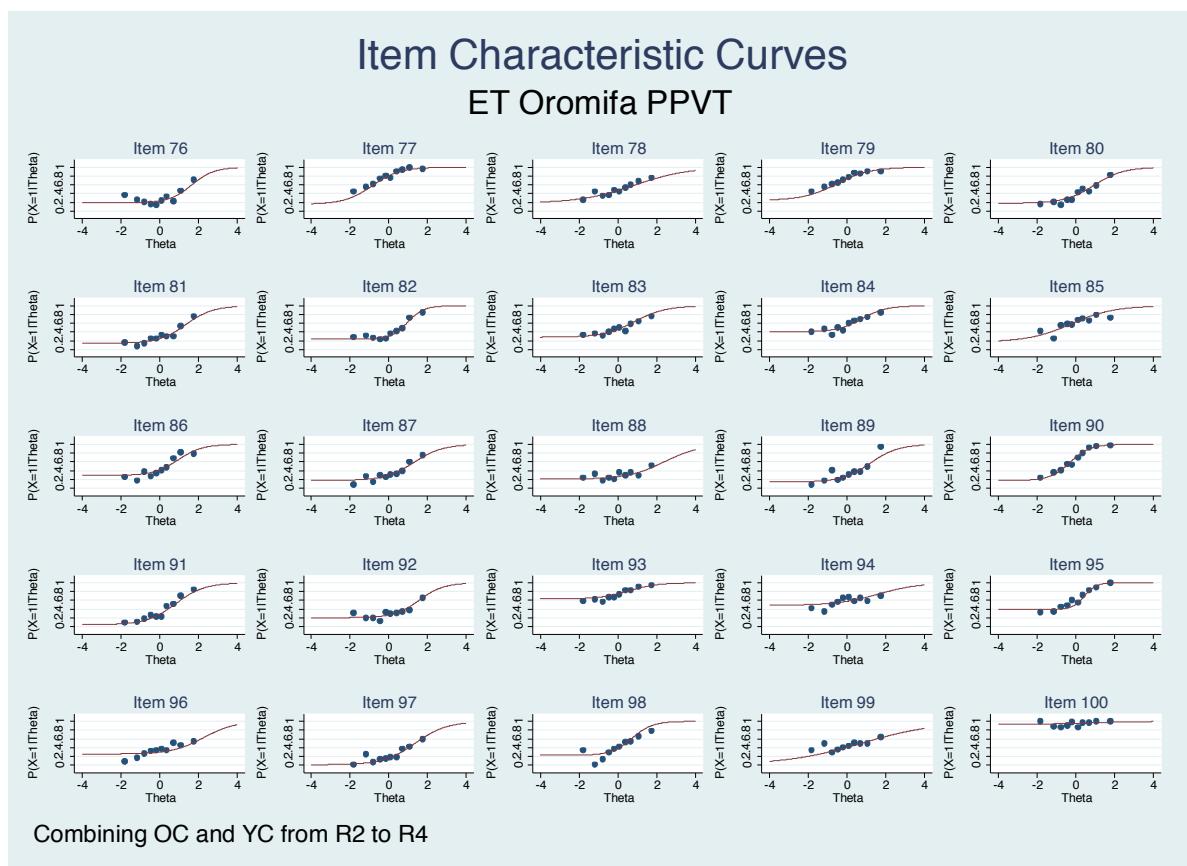
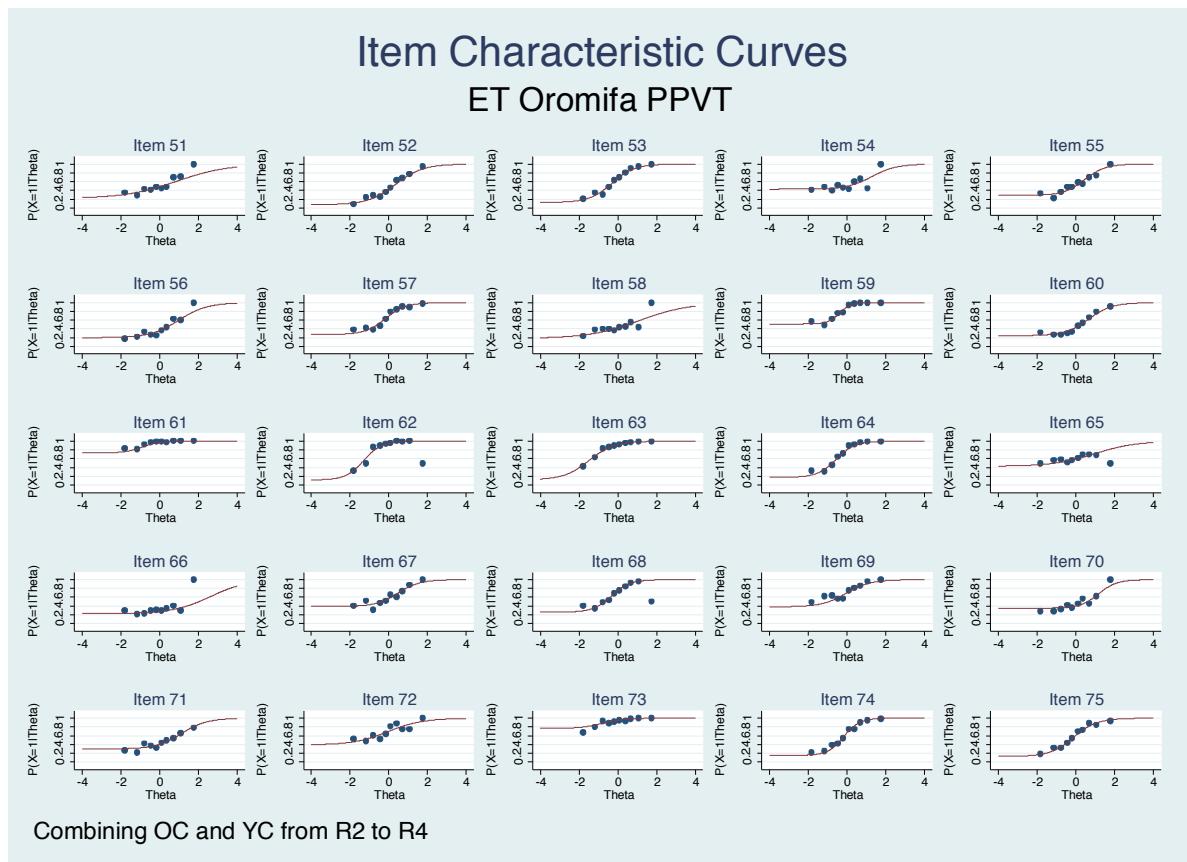
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



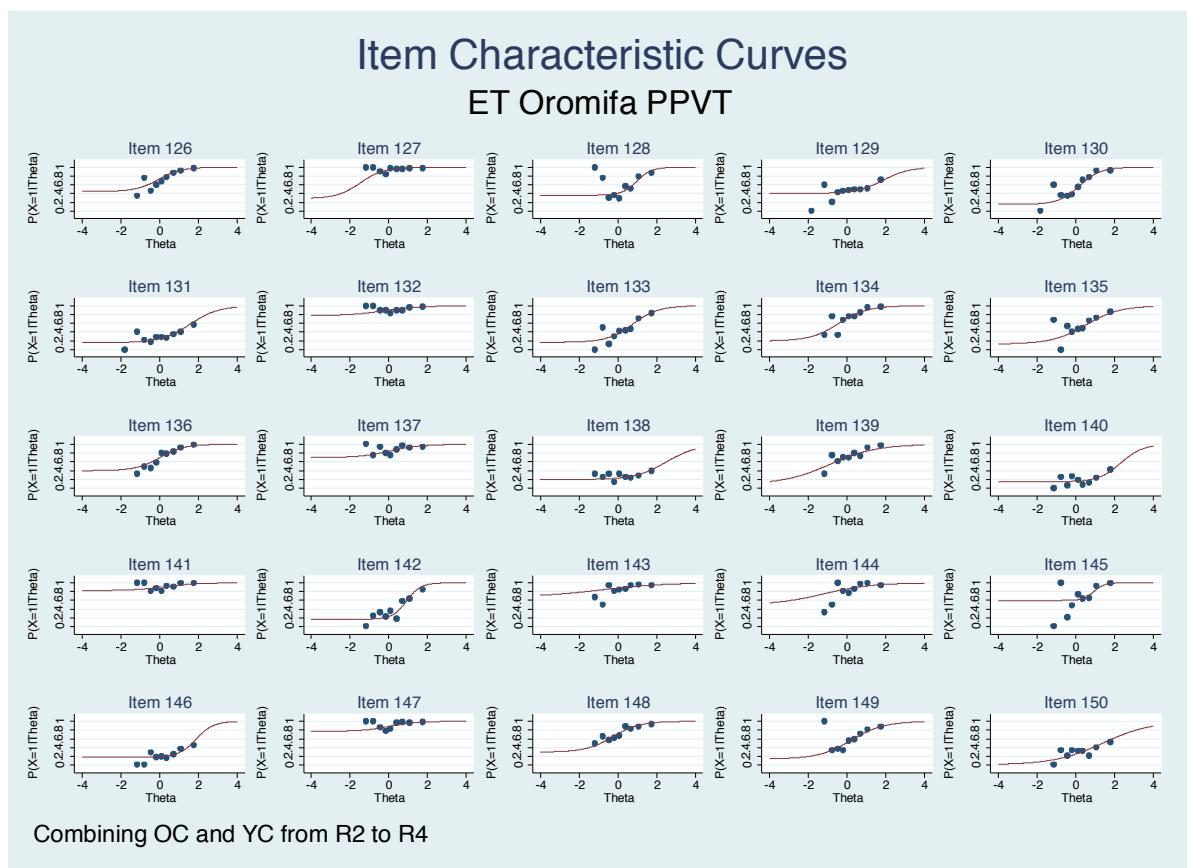
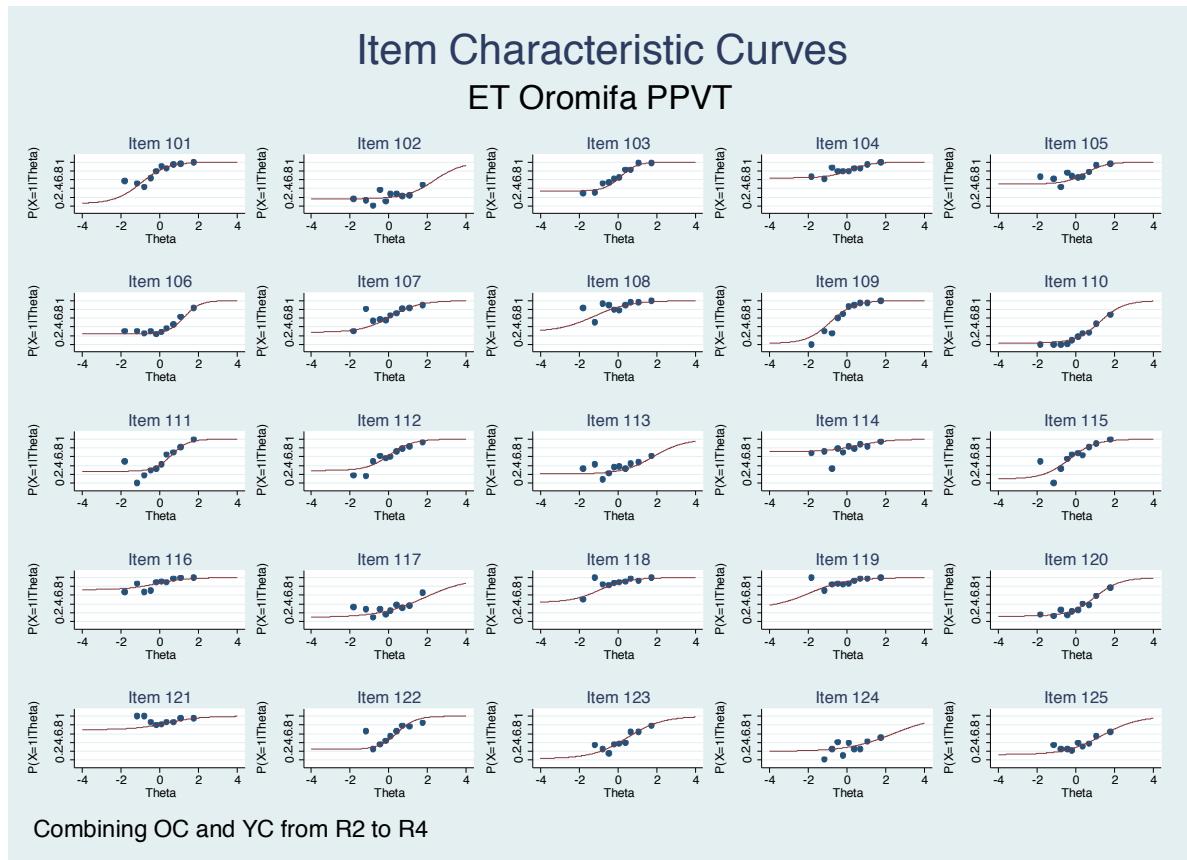


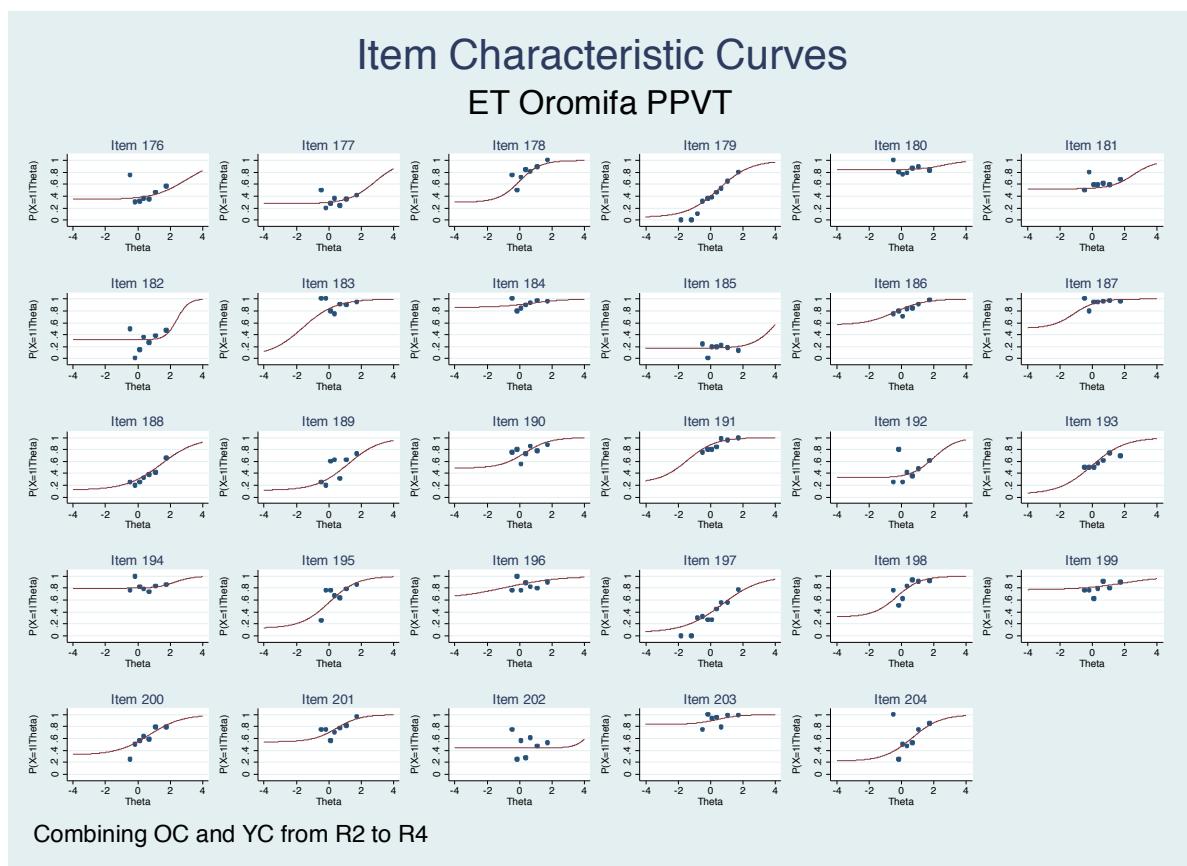
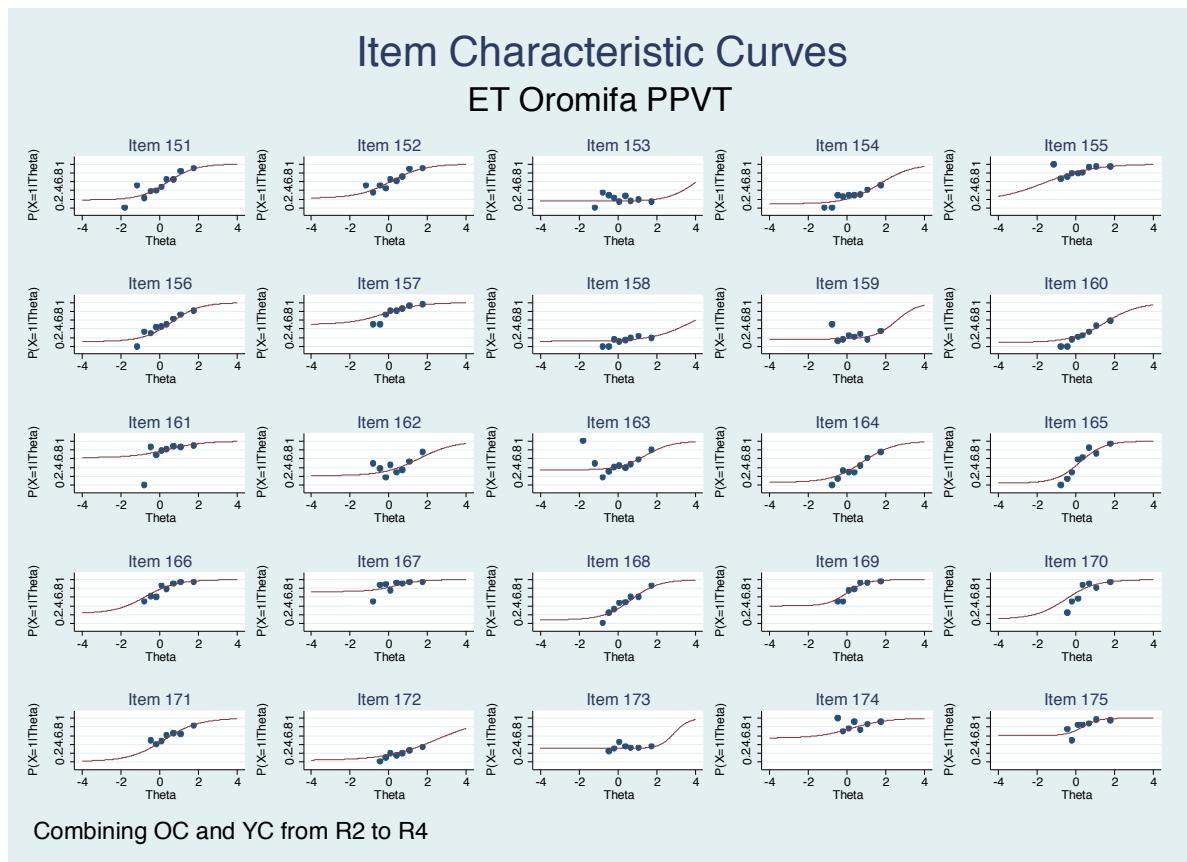
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



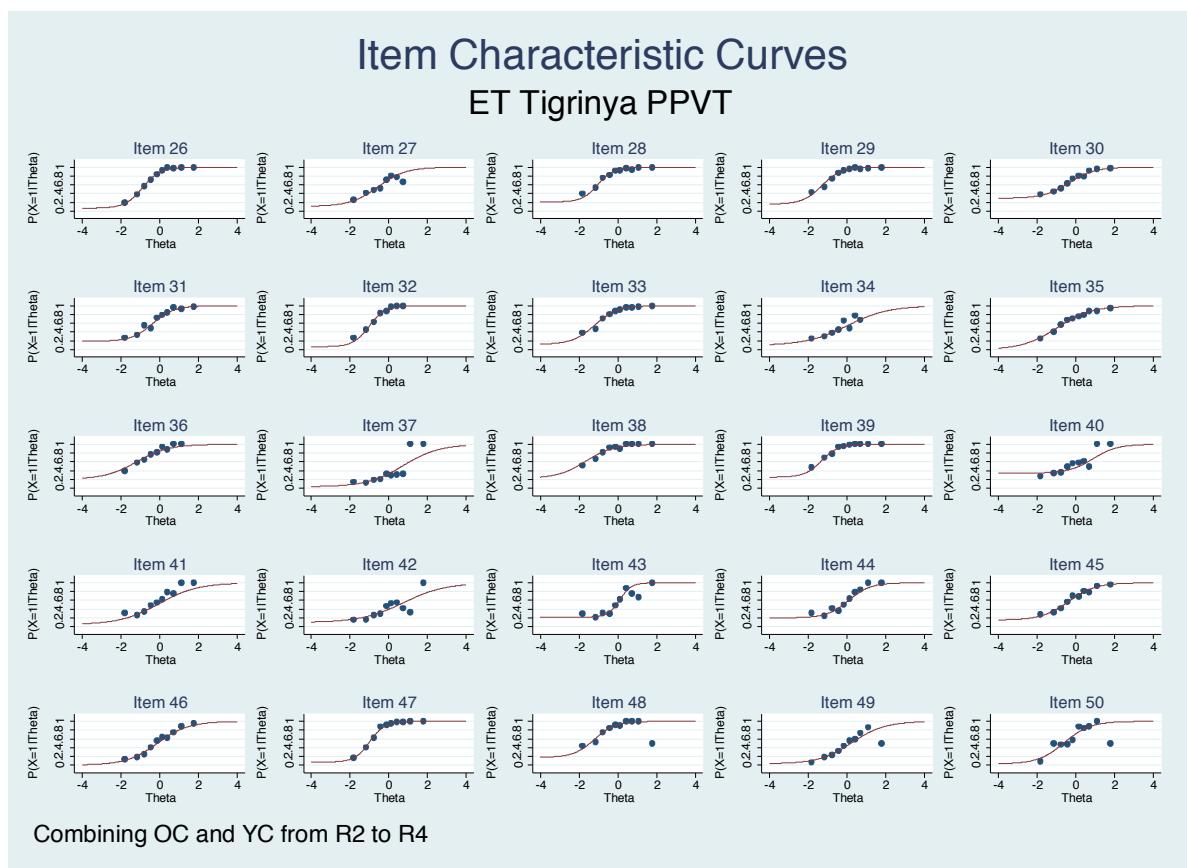
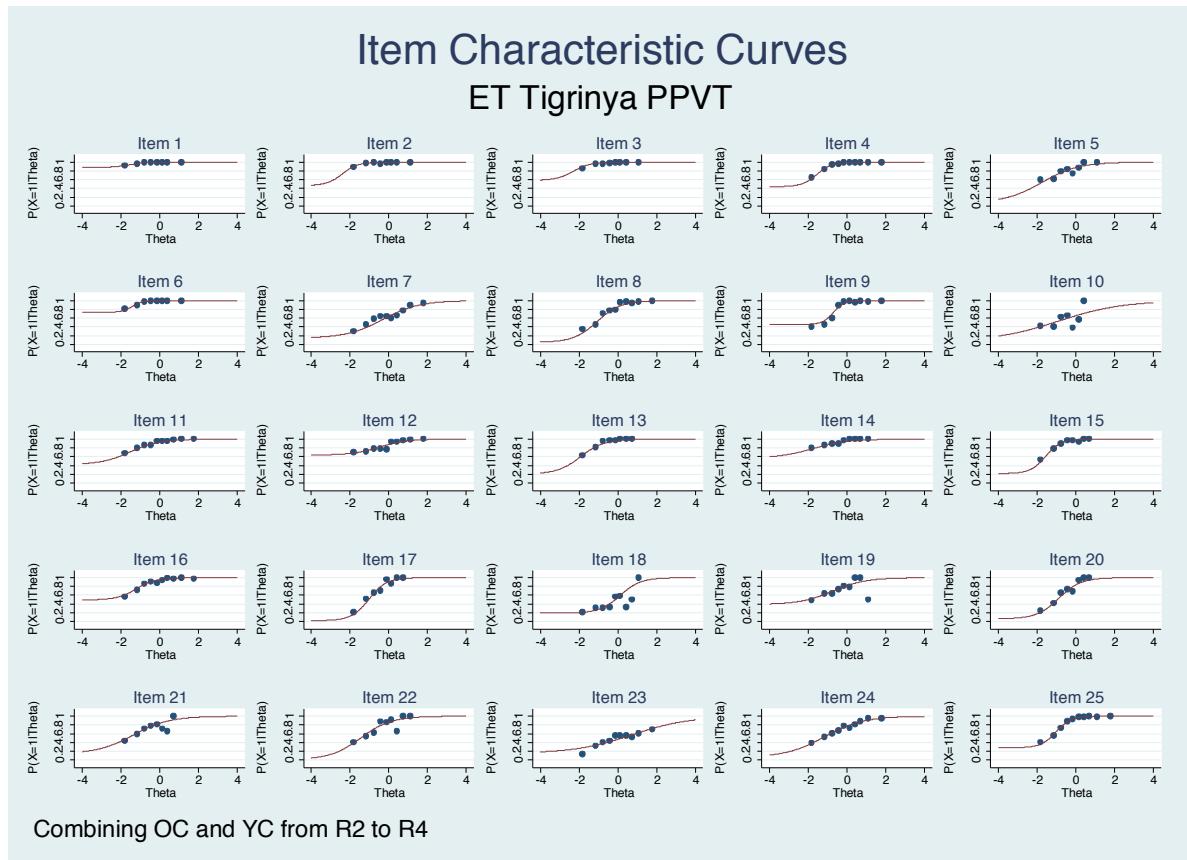


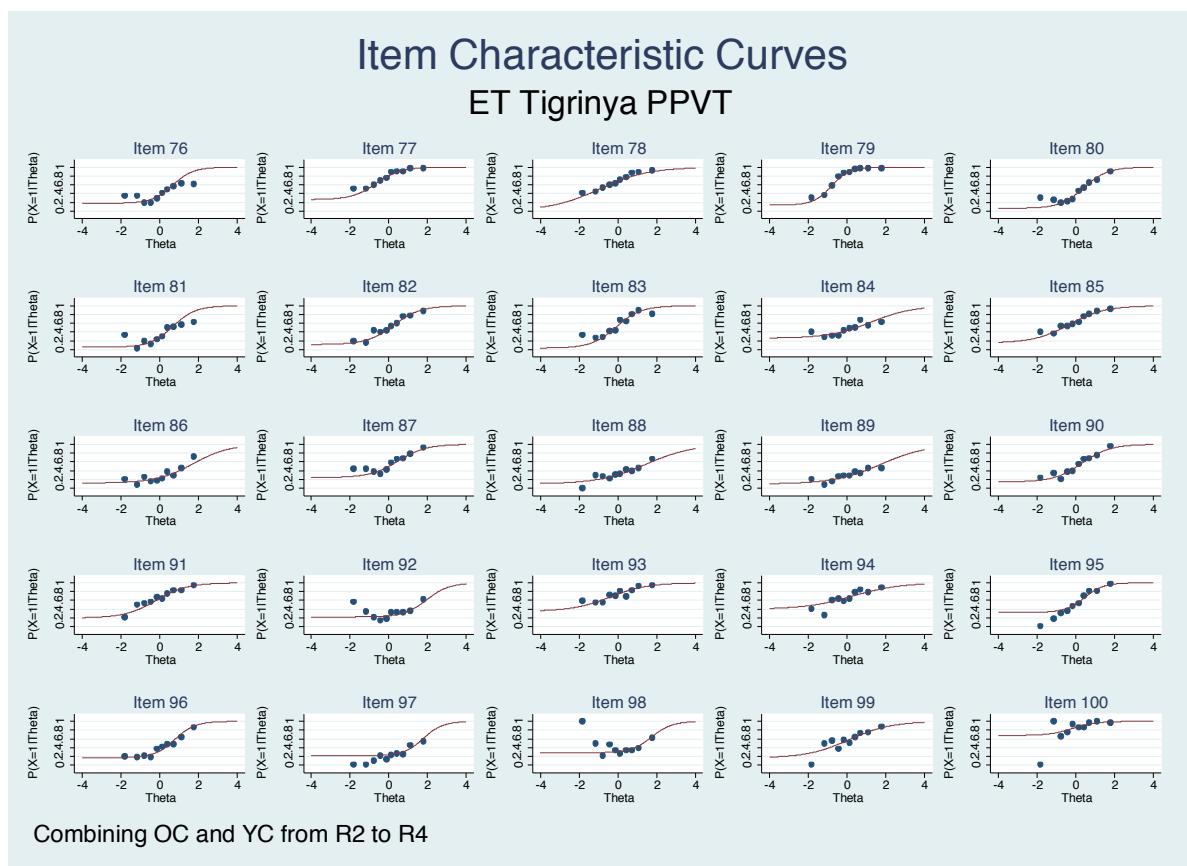
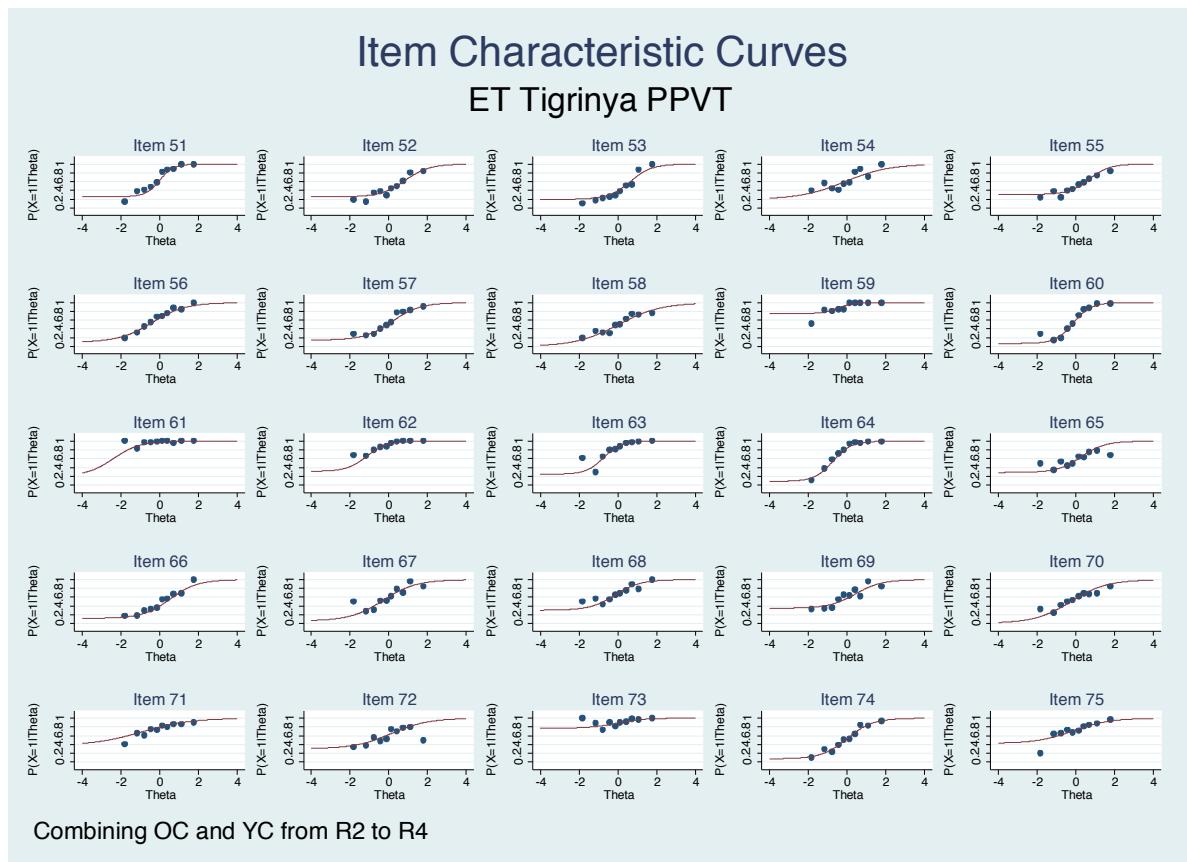
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



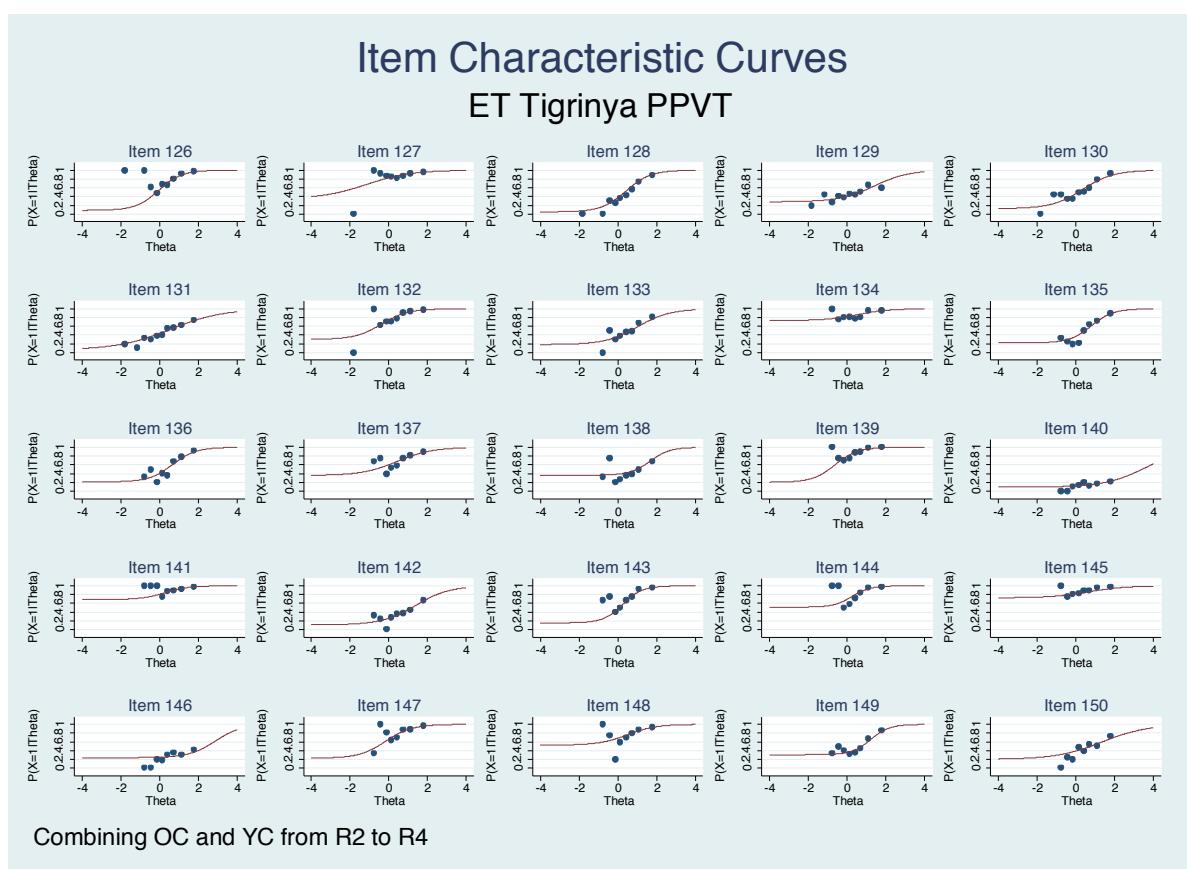
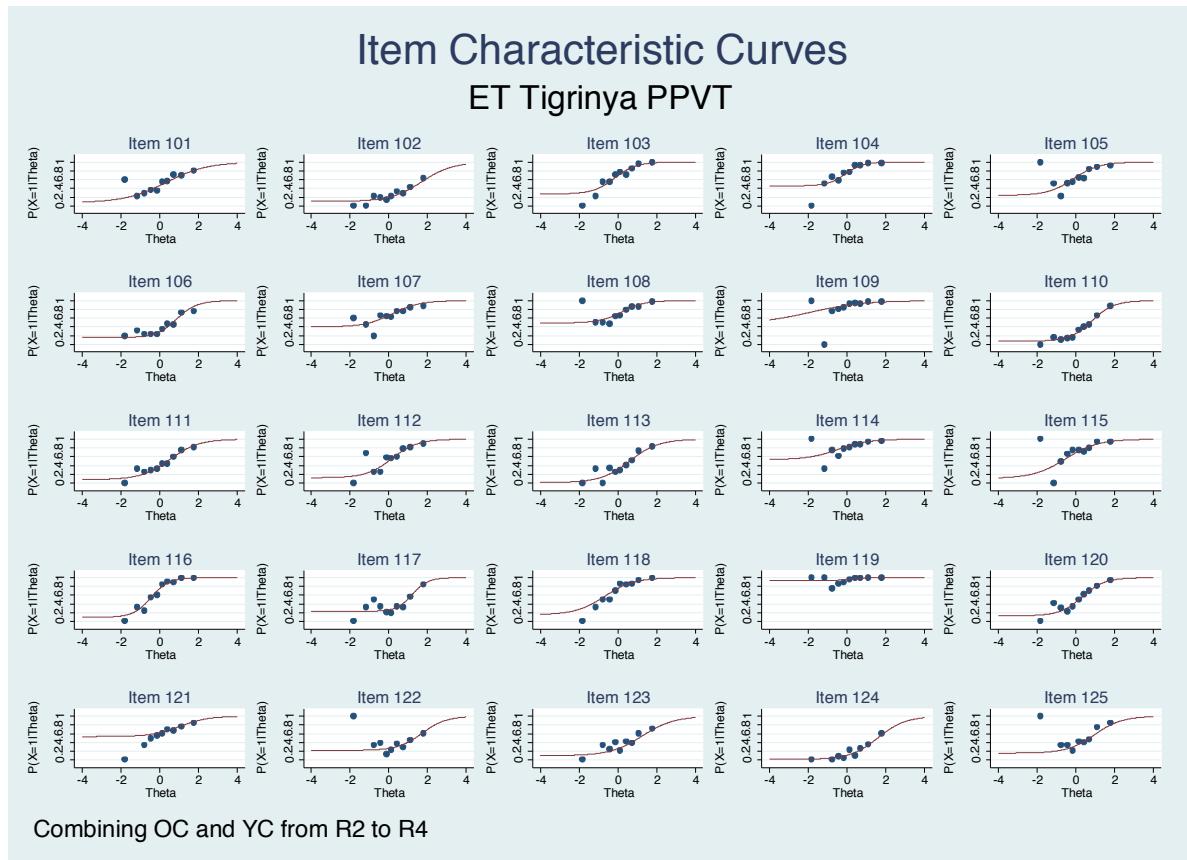


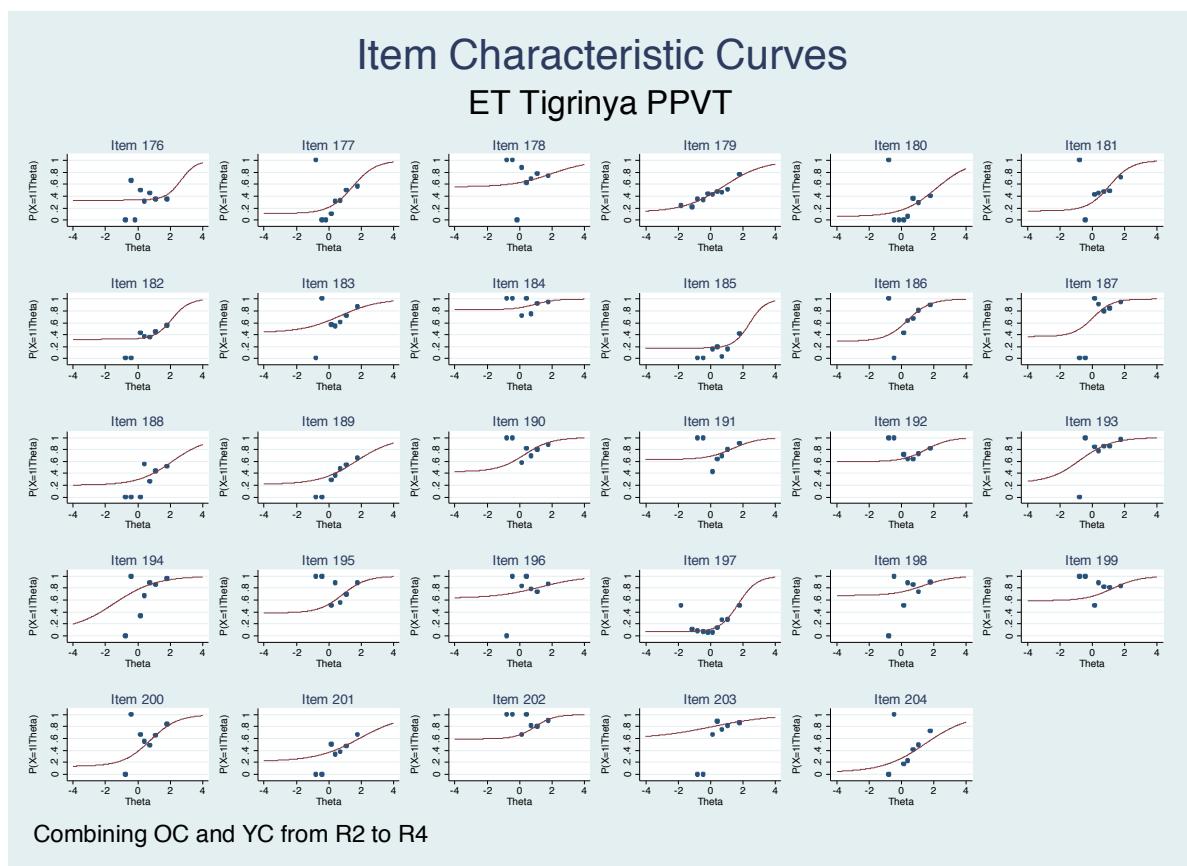
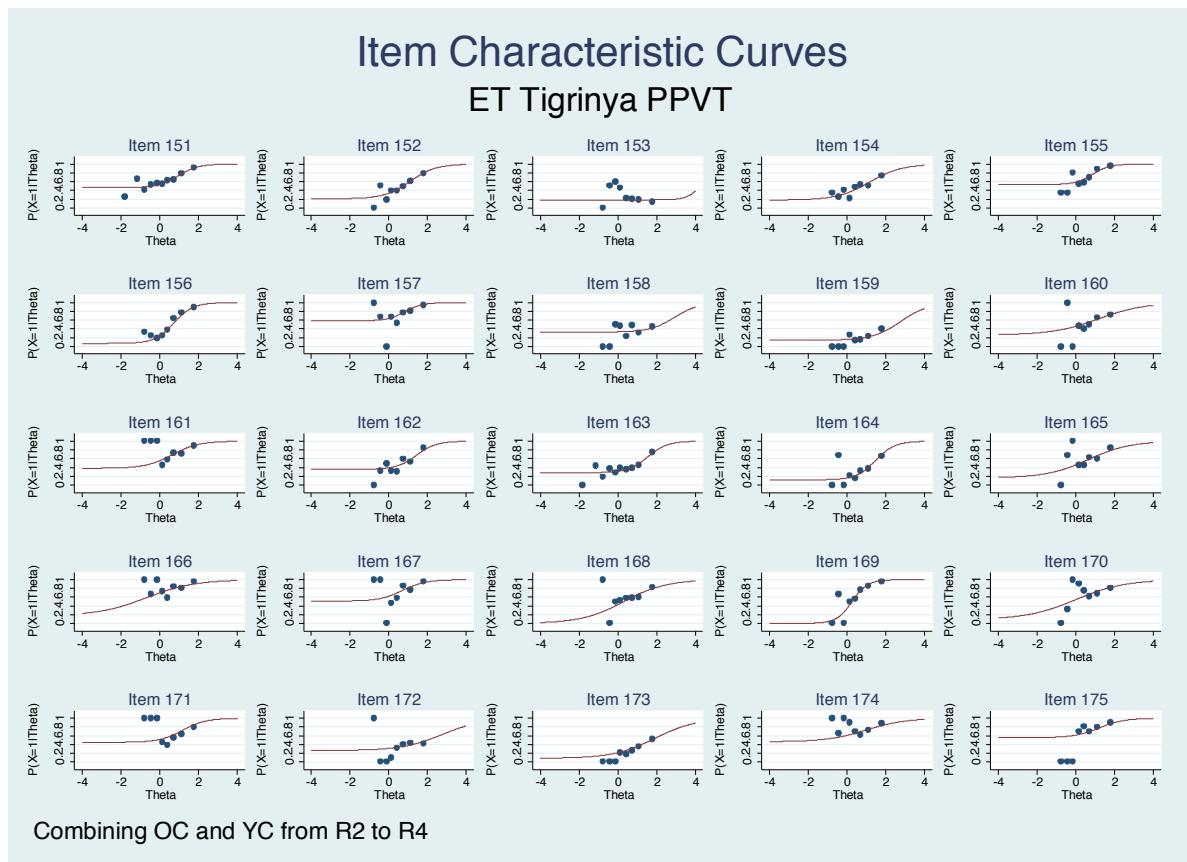
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



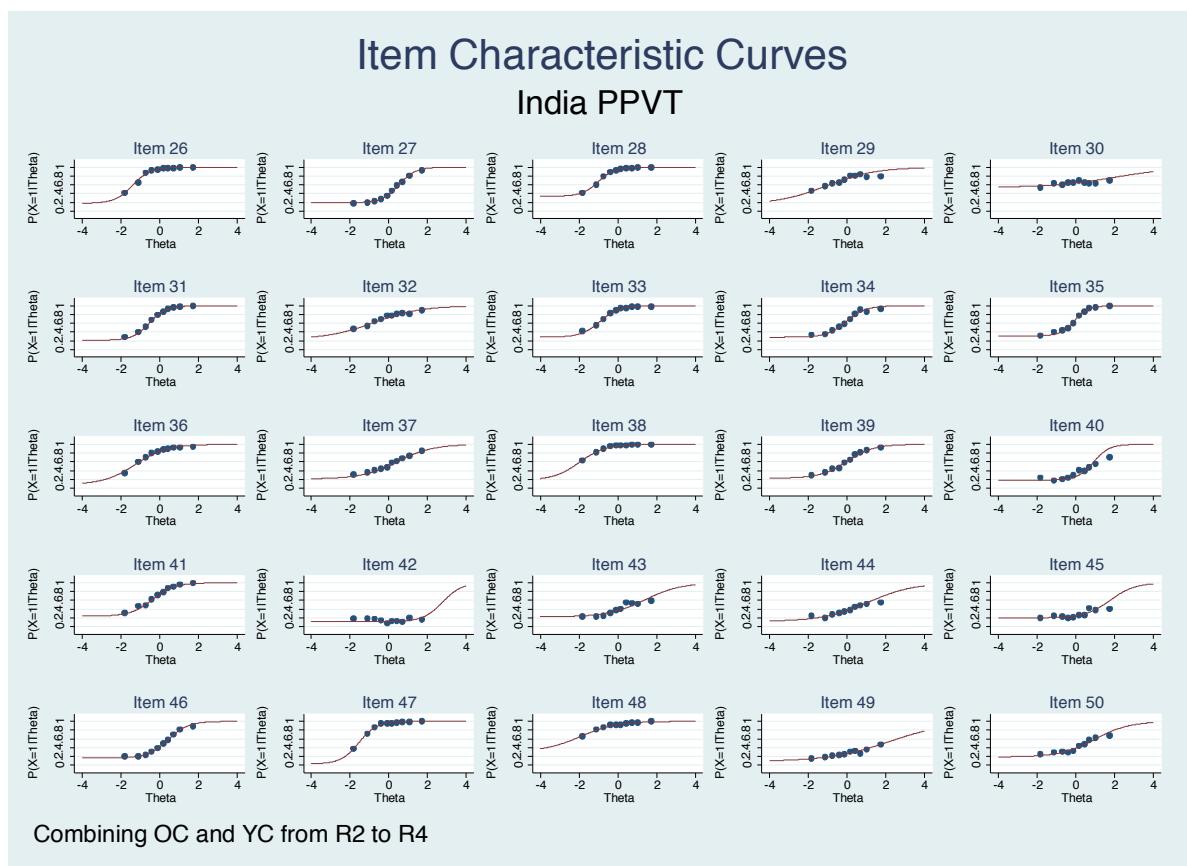
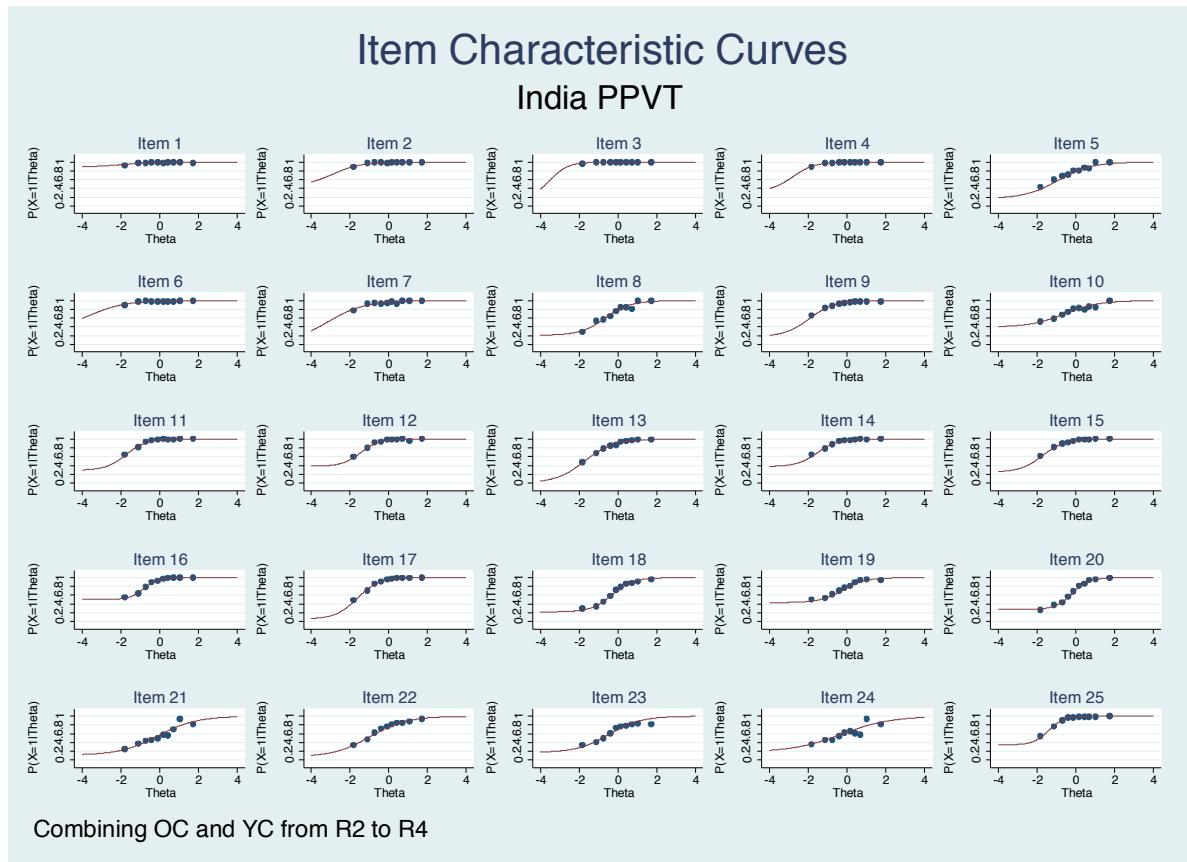


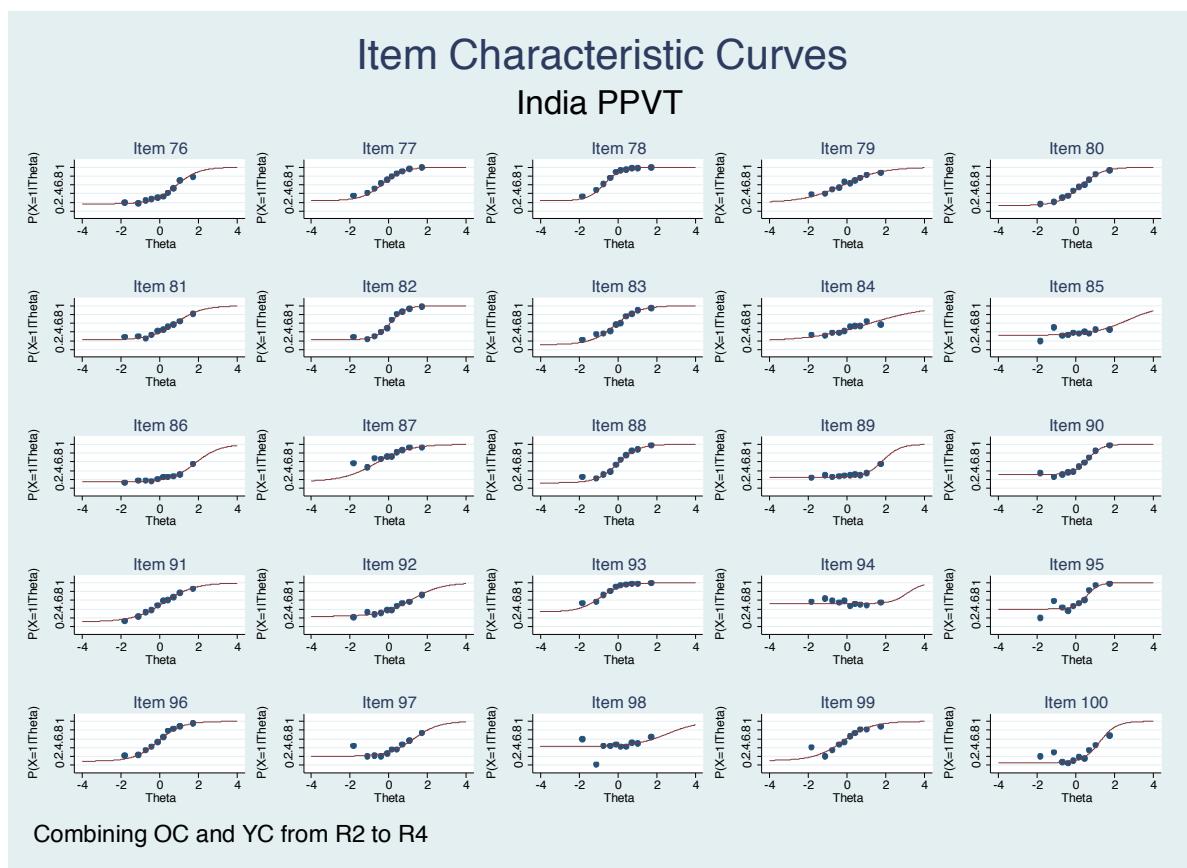
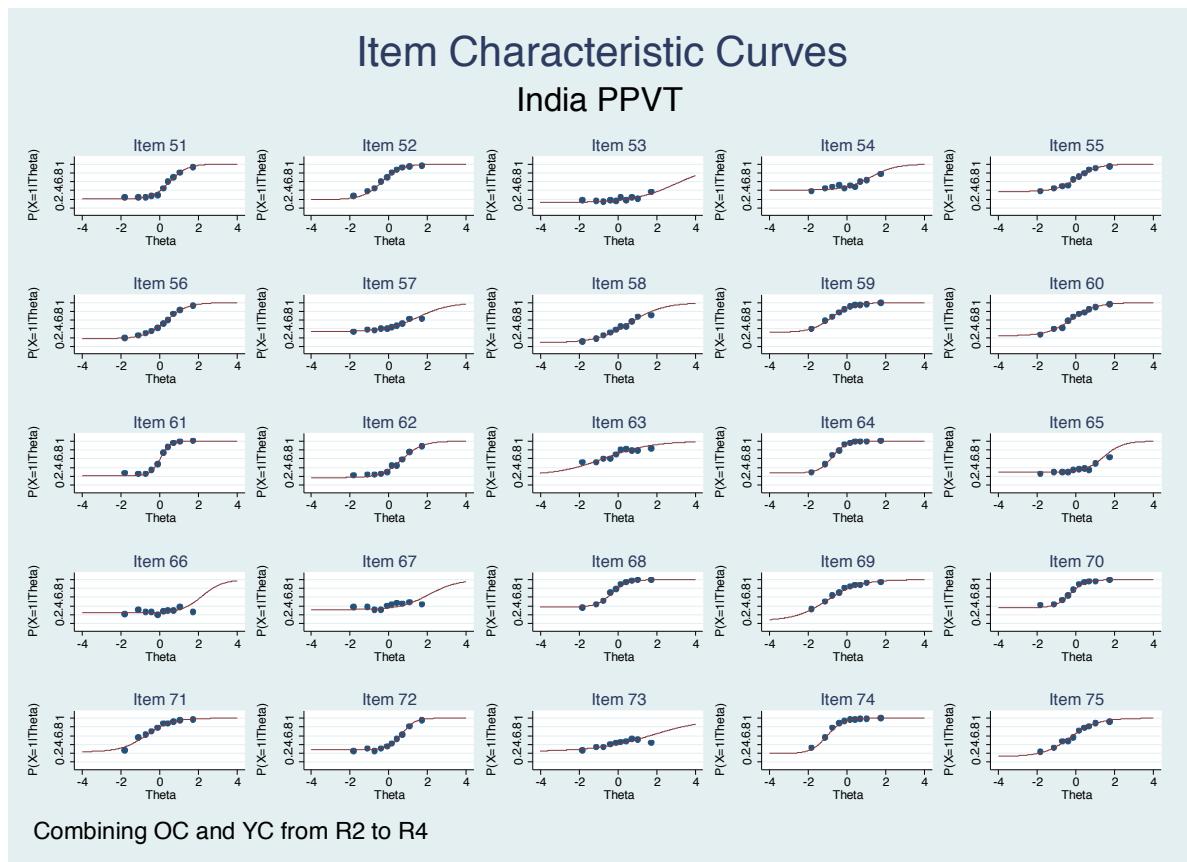
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

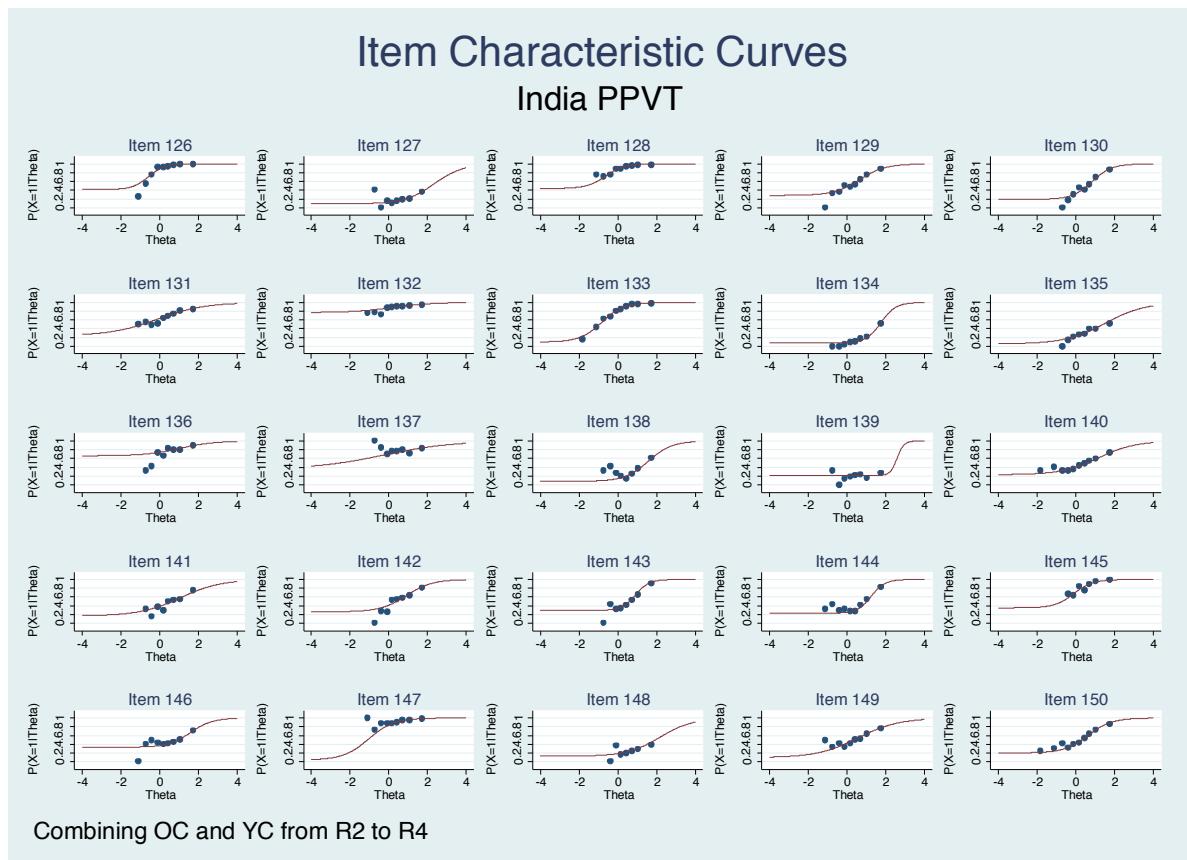
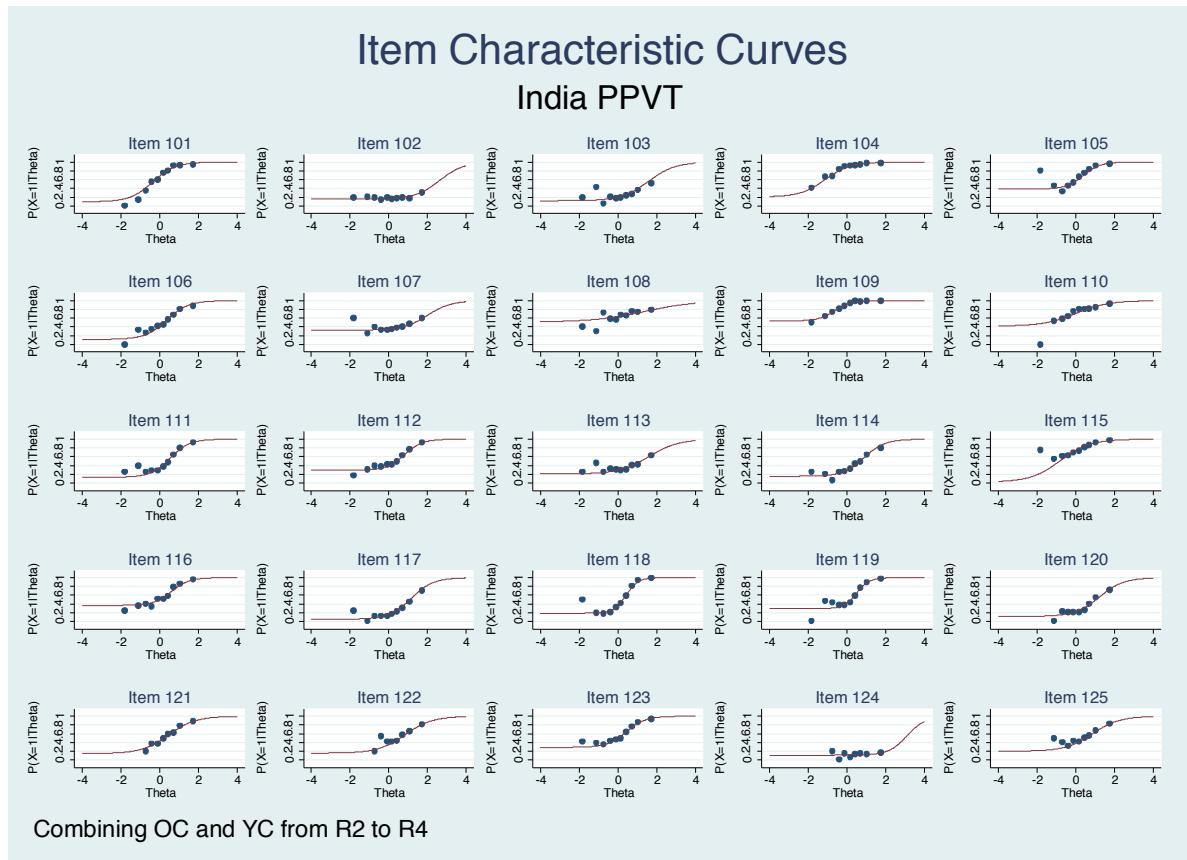


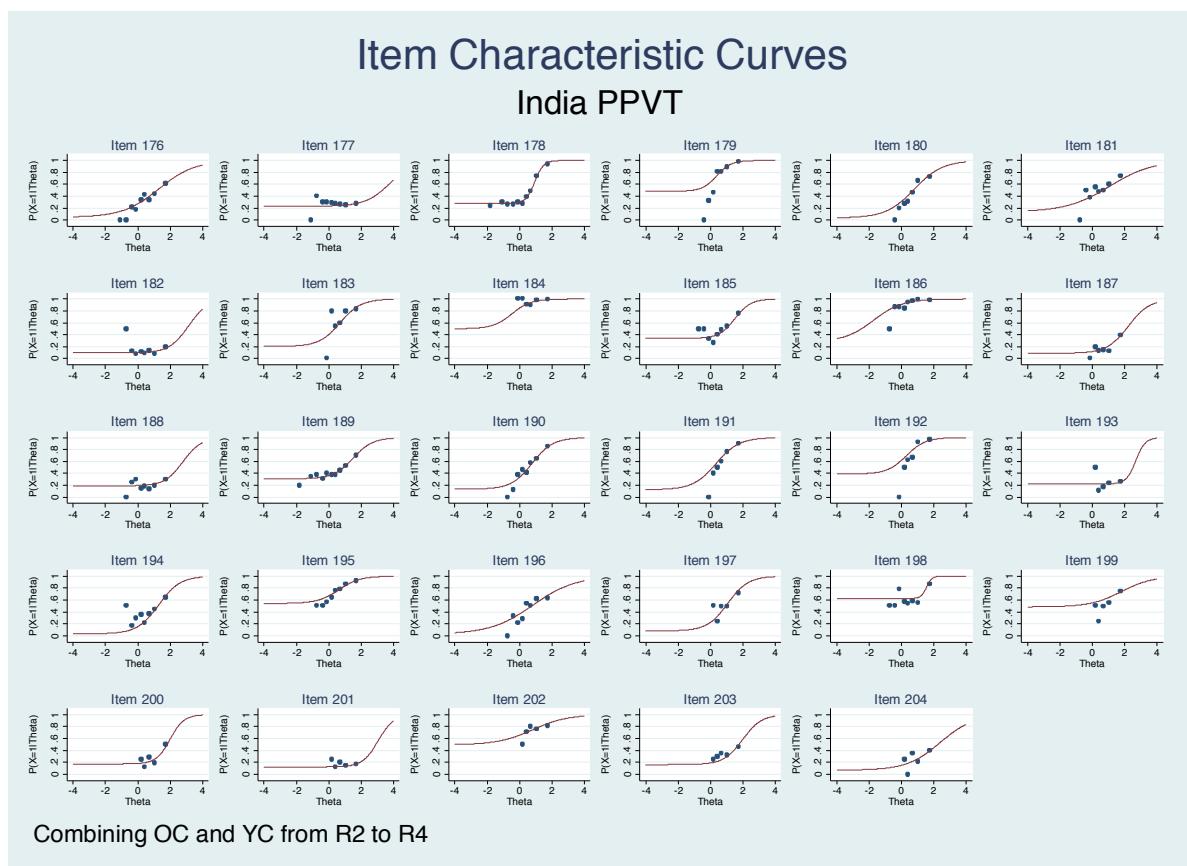
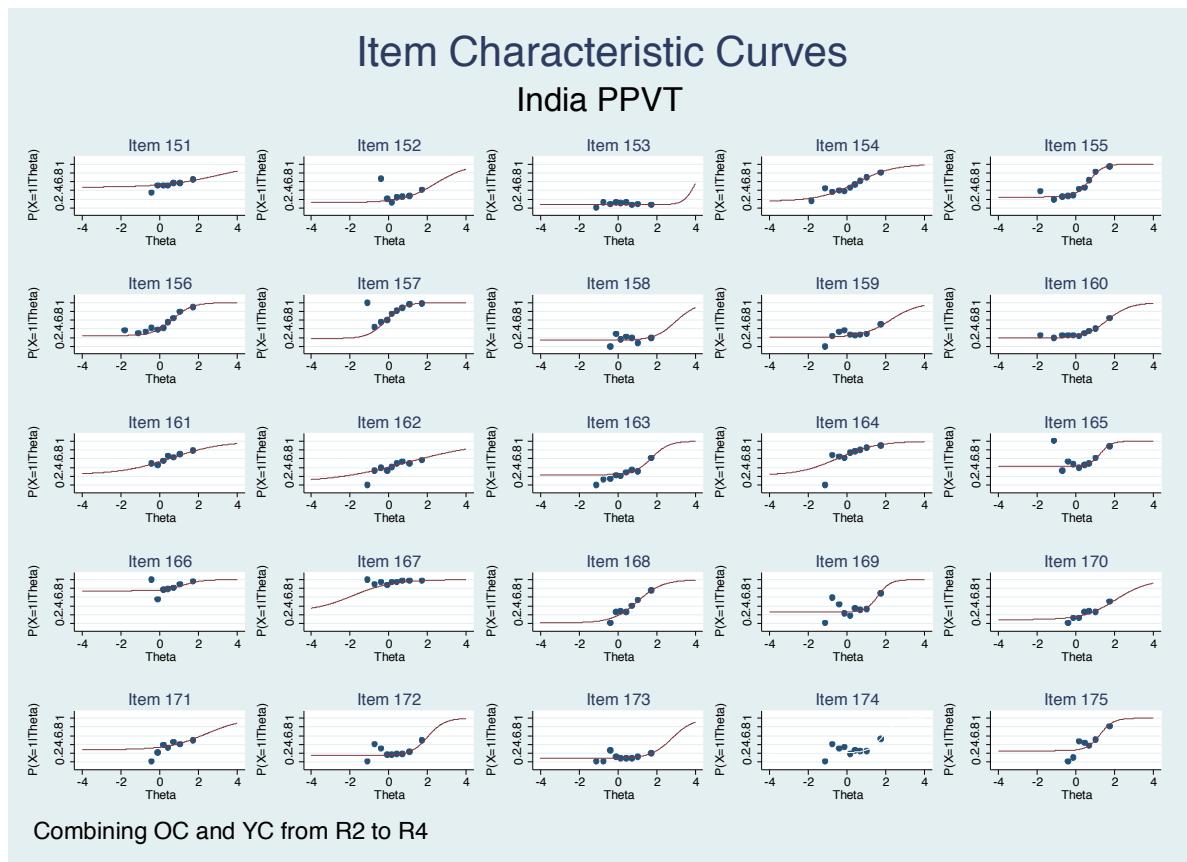


EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

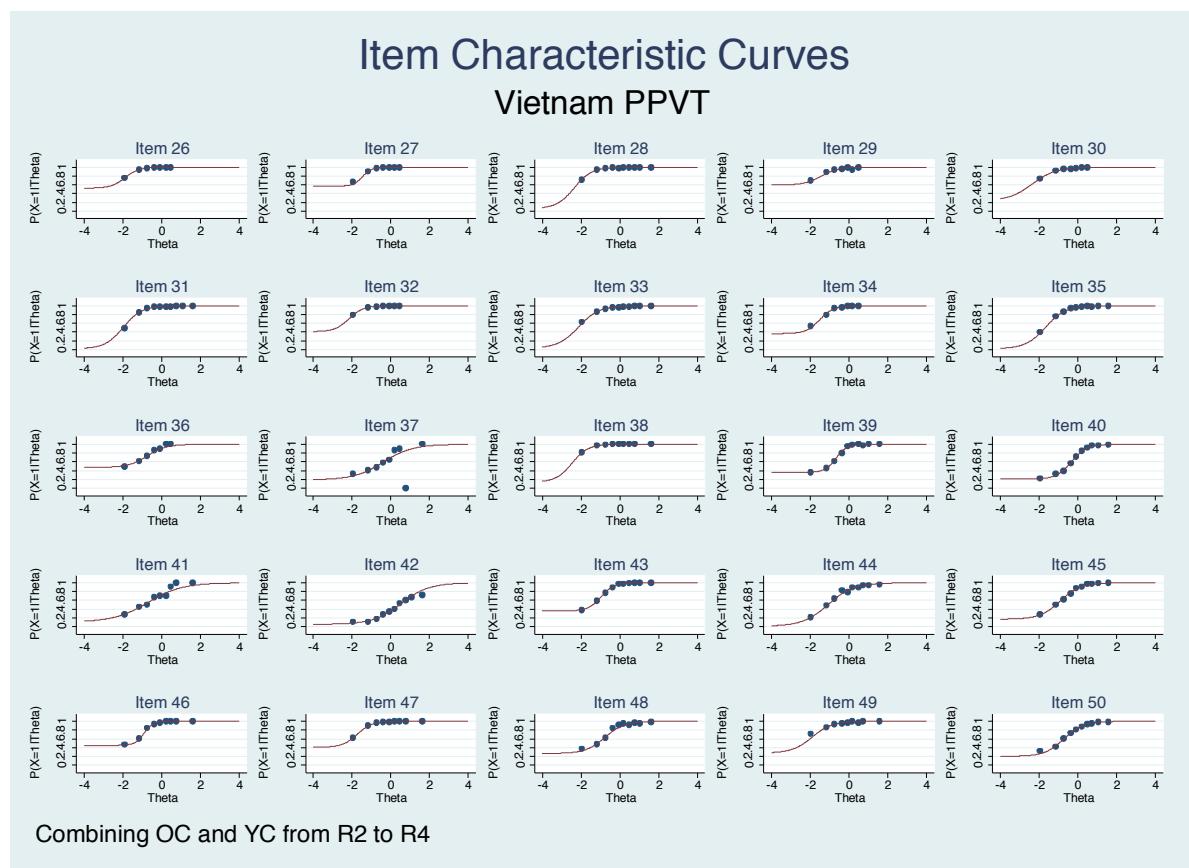
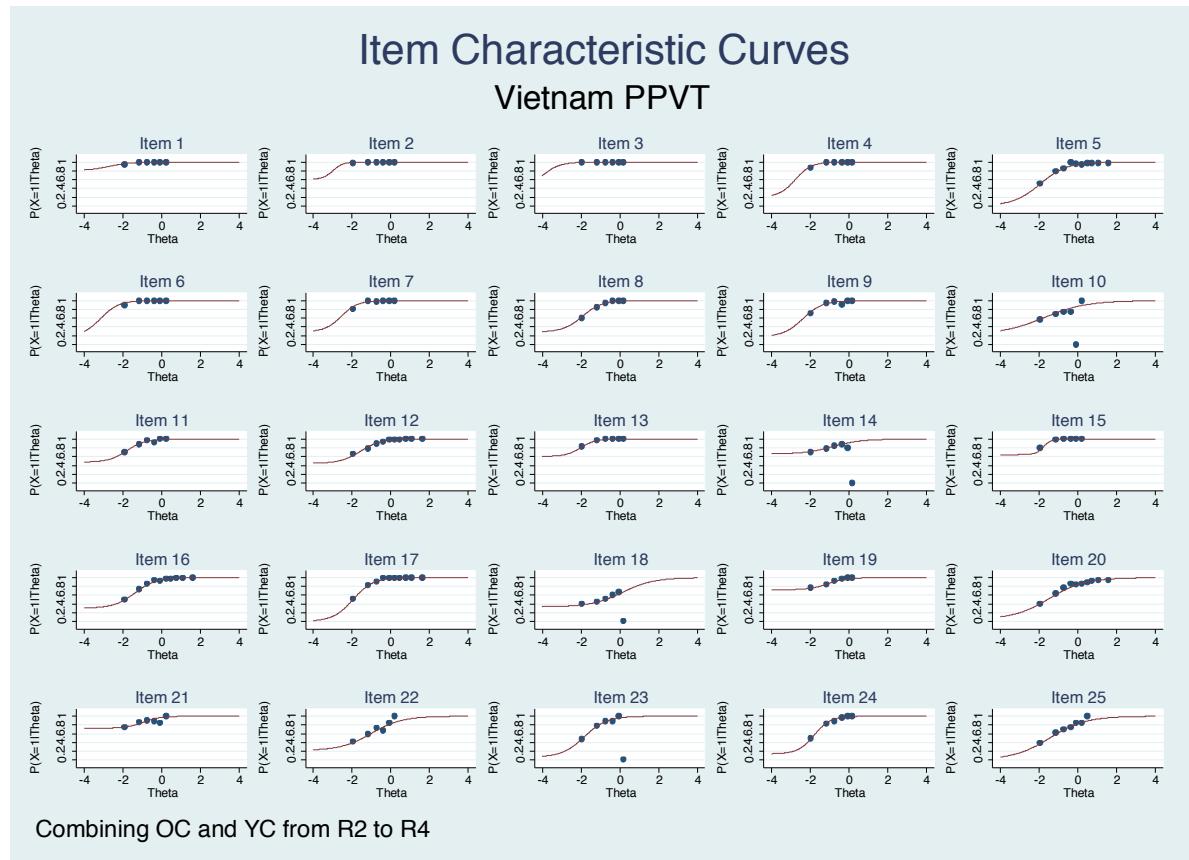


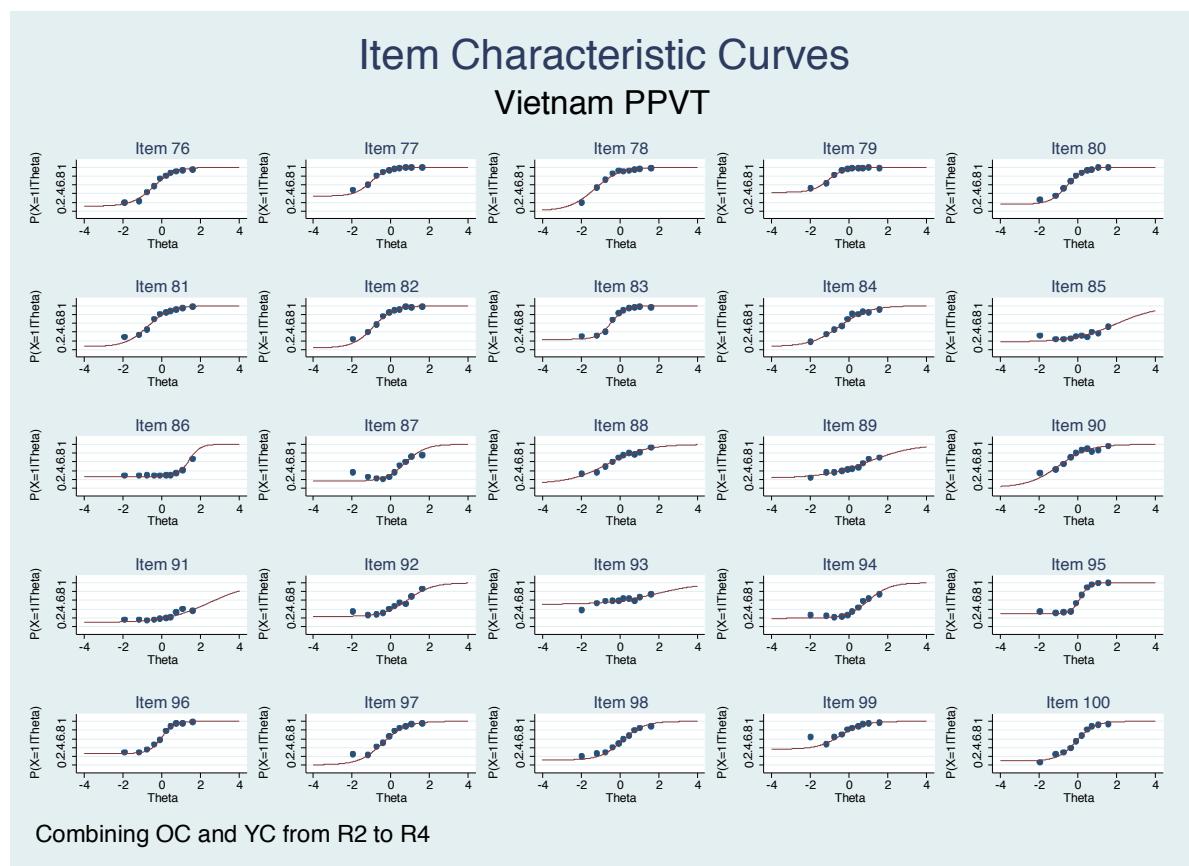
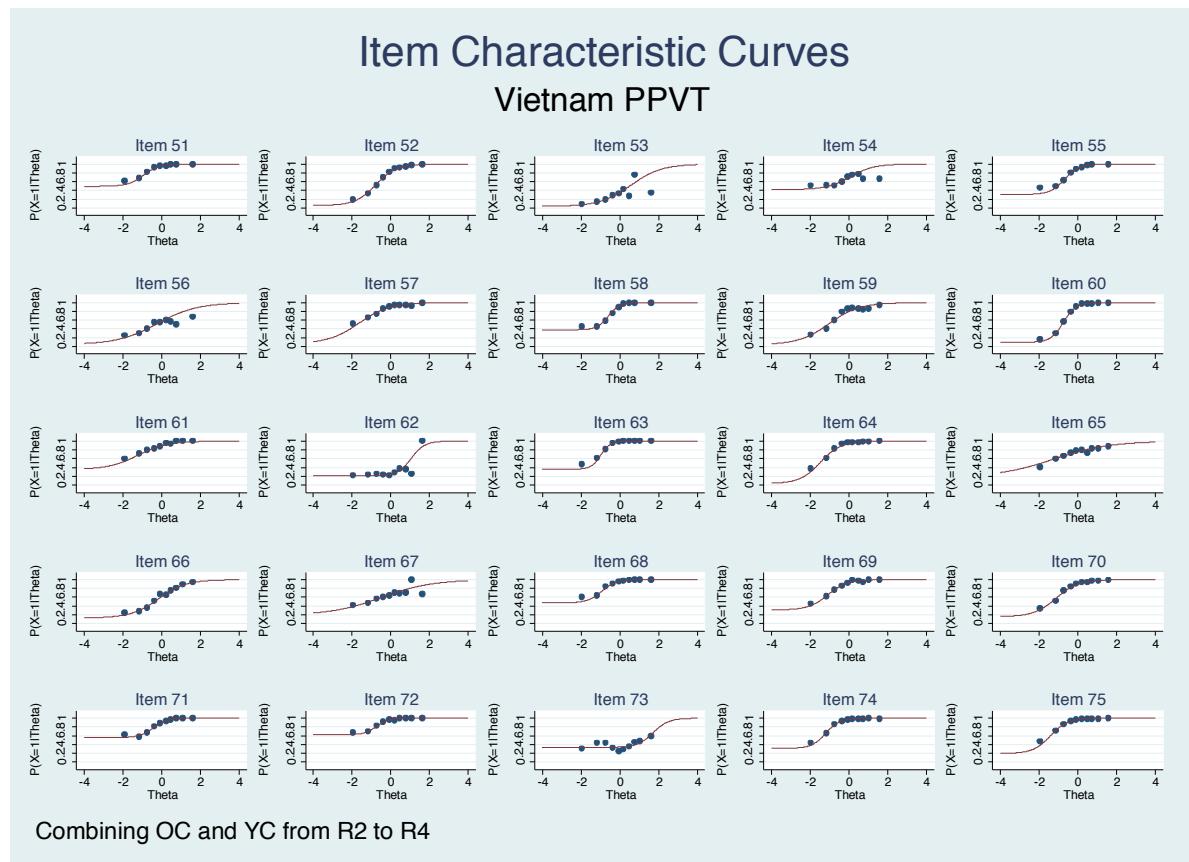




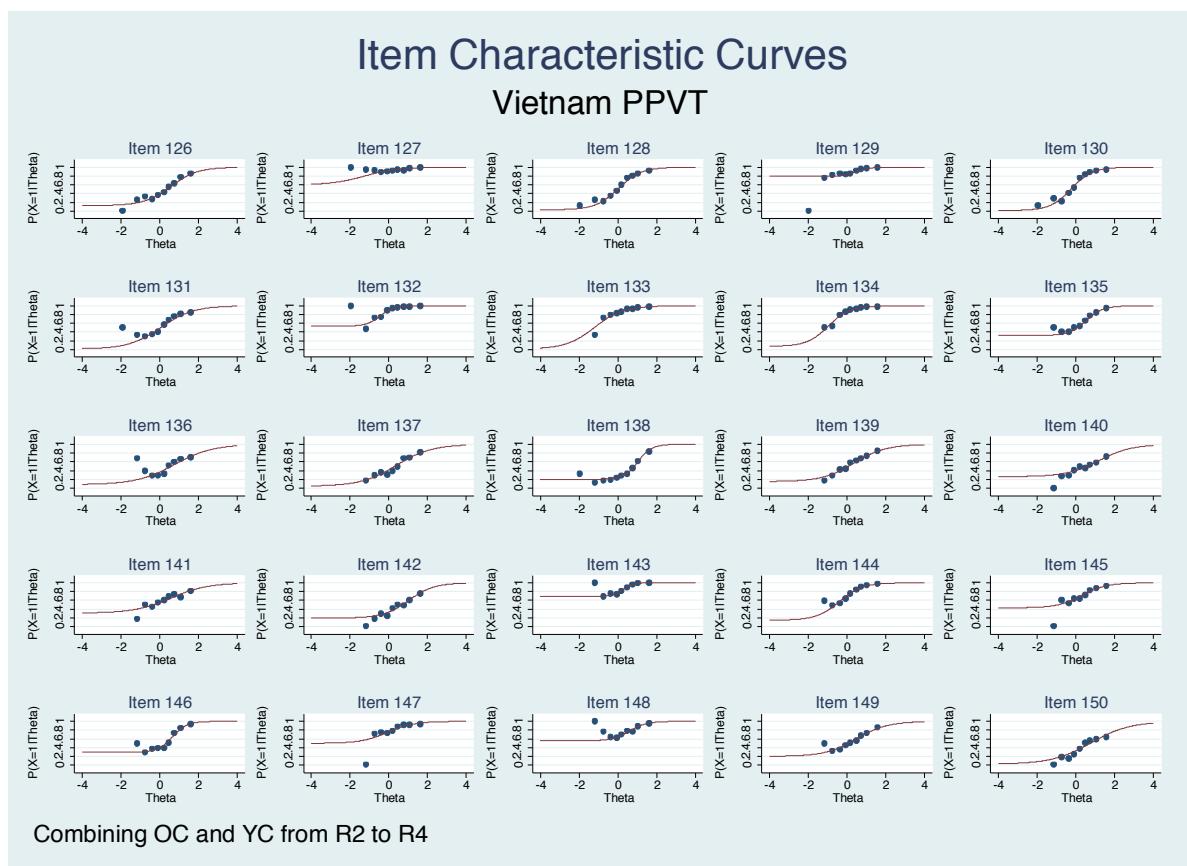
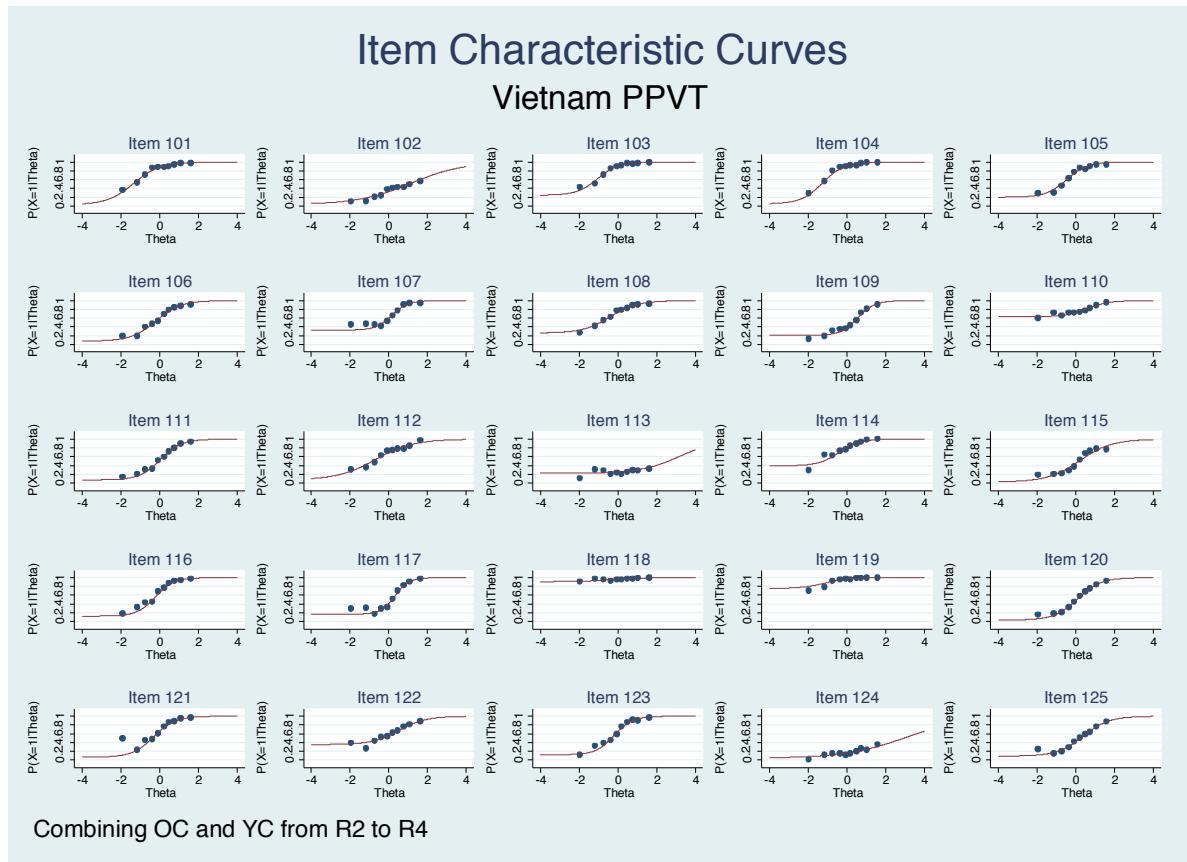


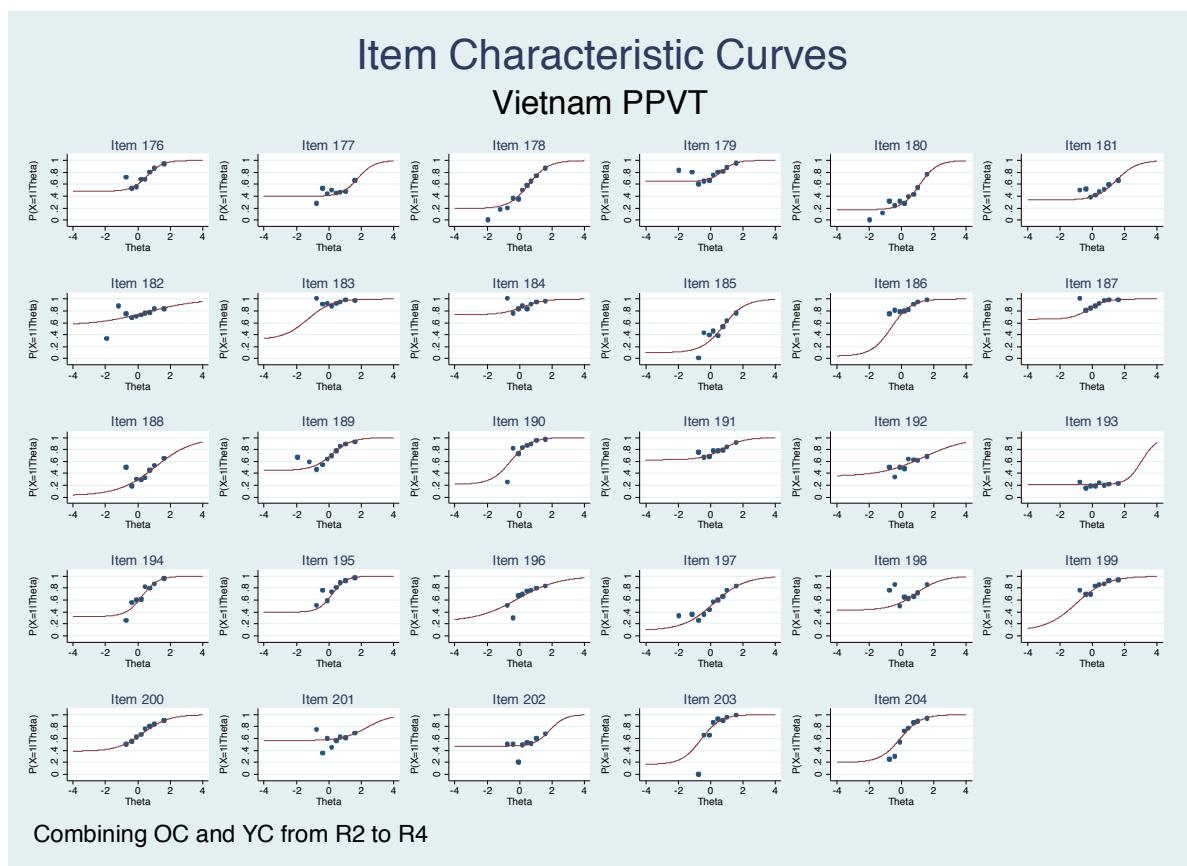
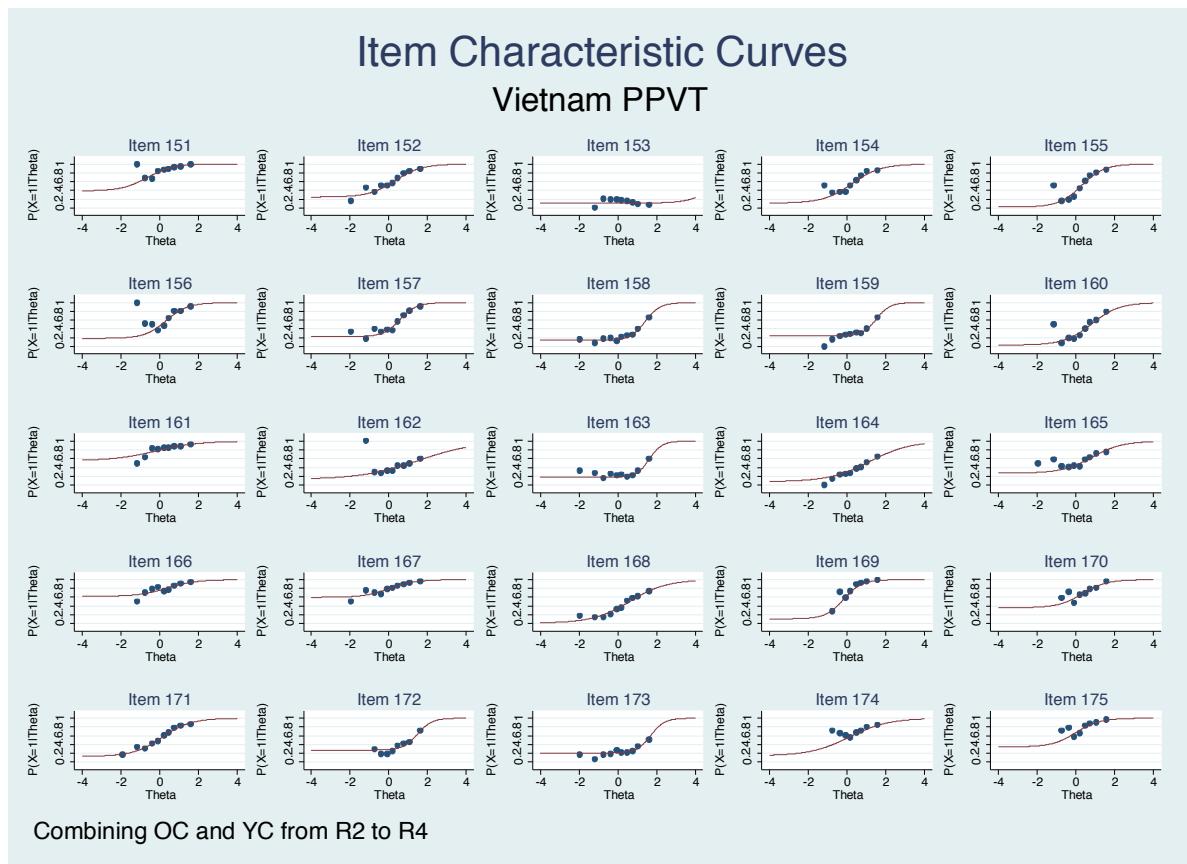
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



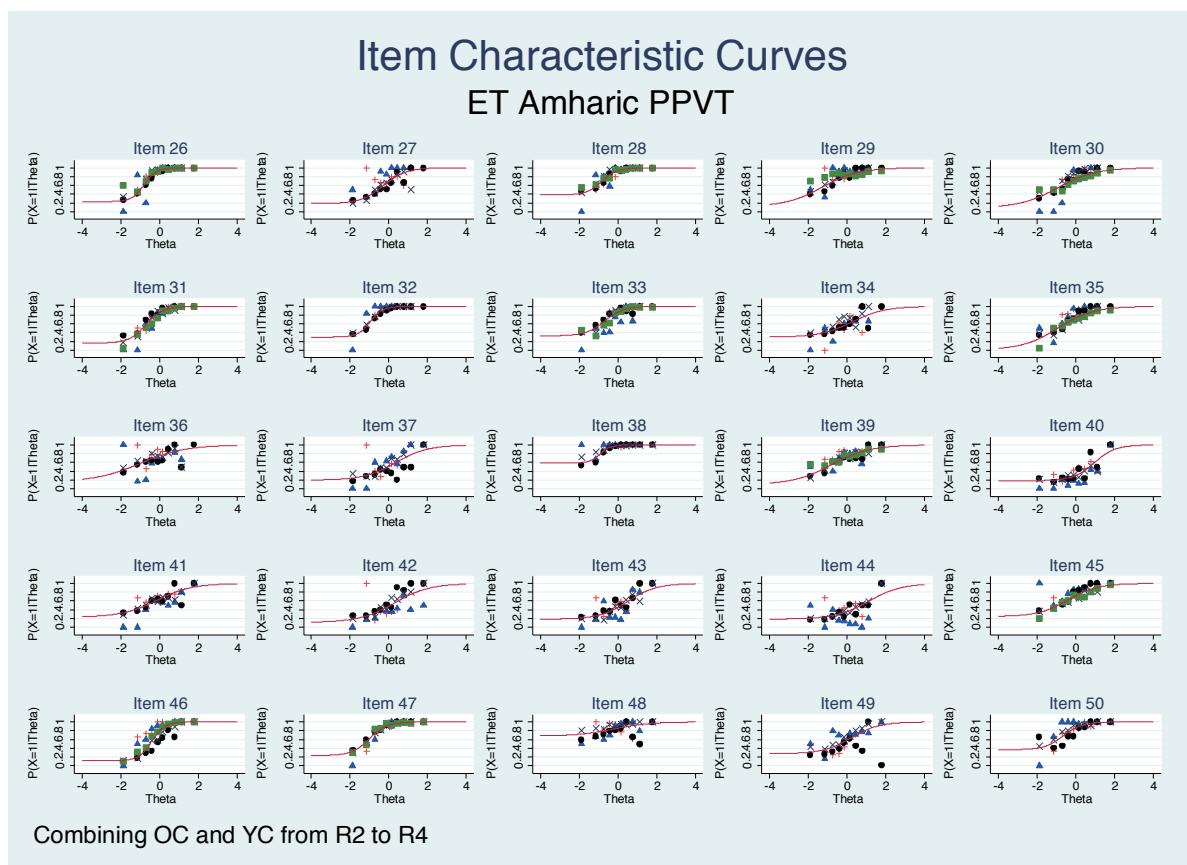
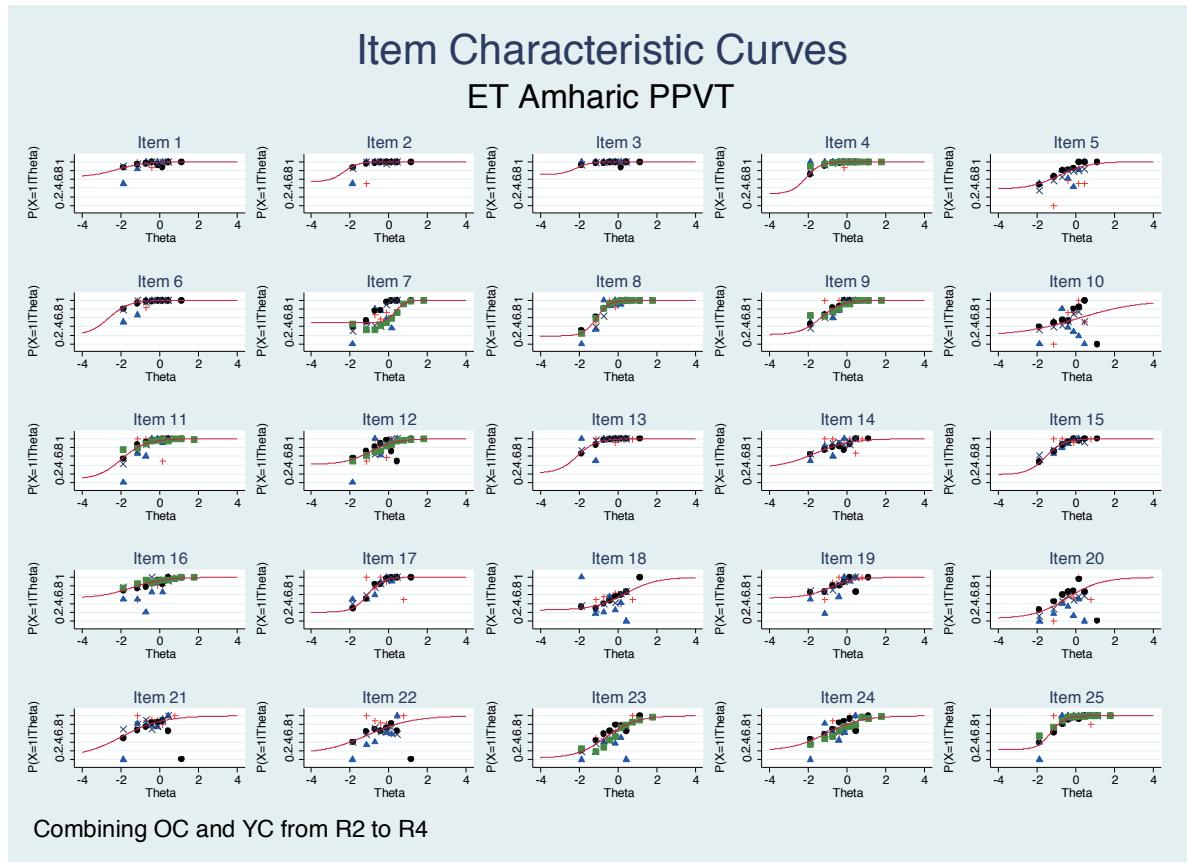


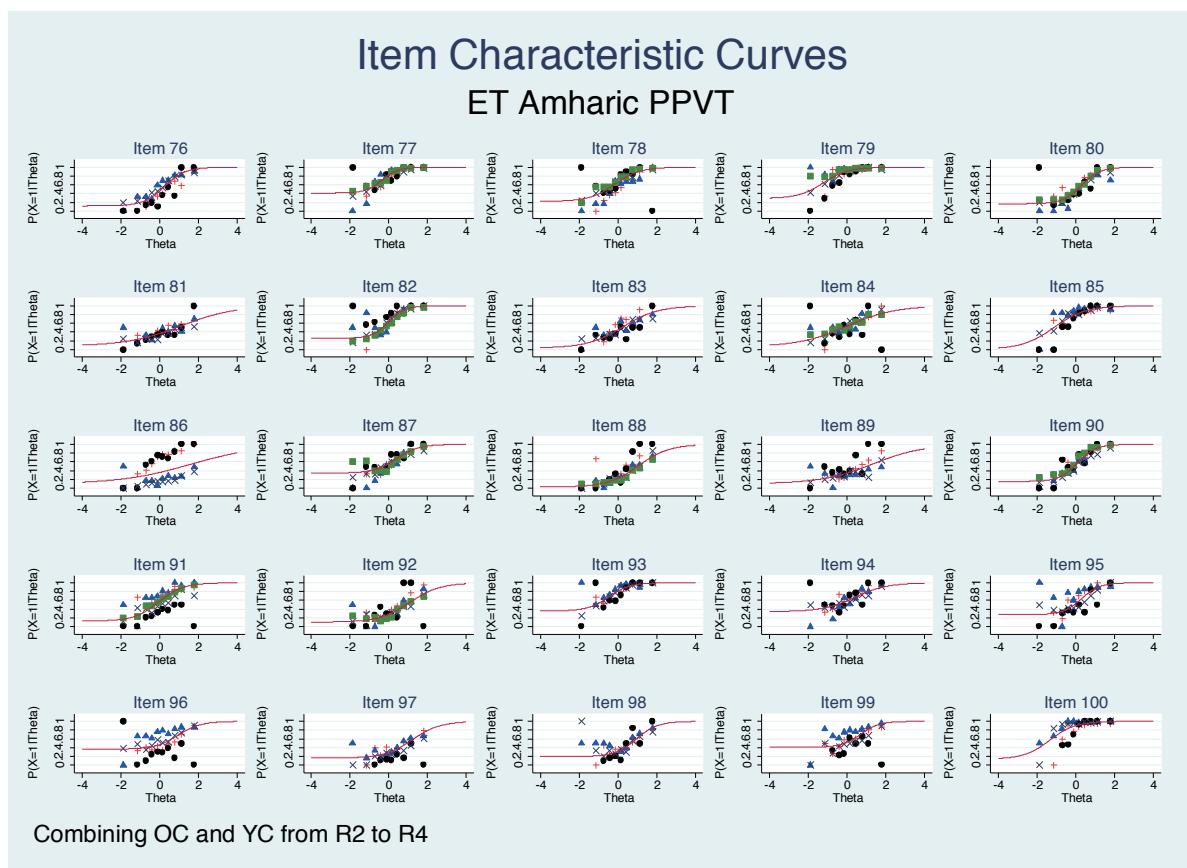
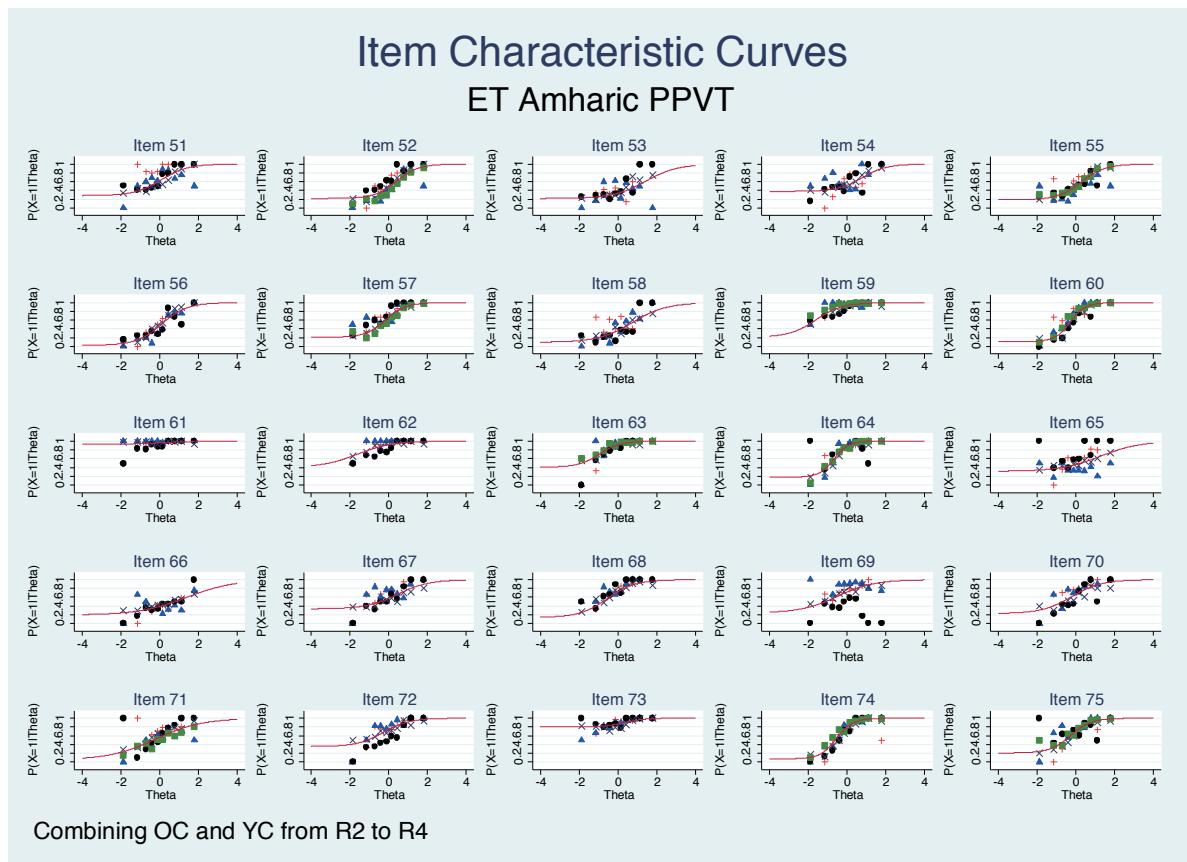
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



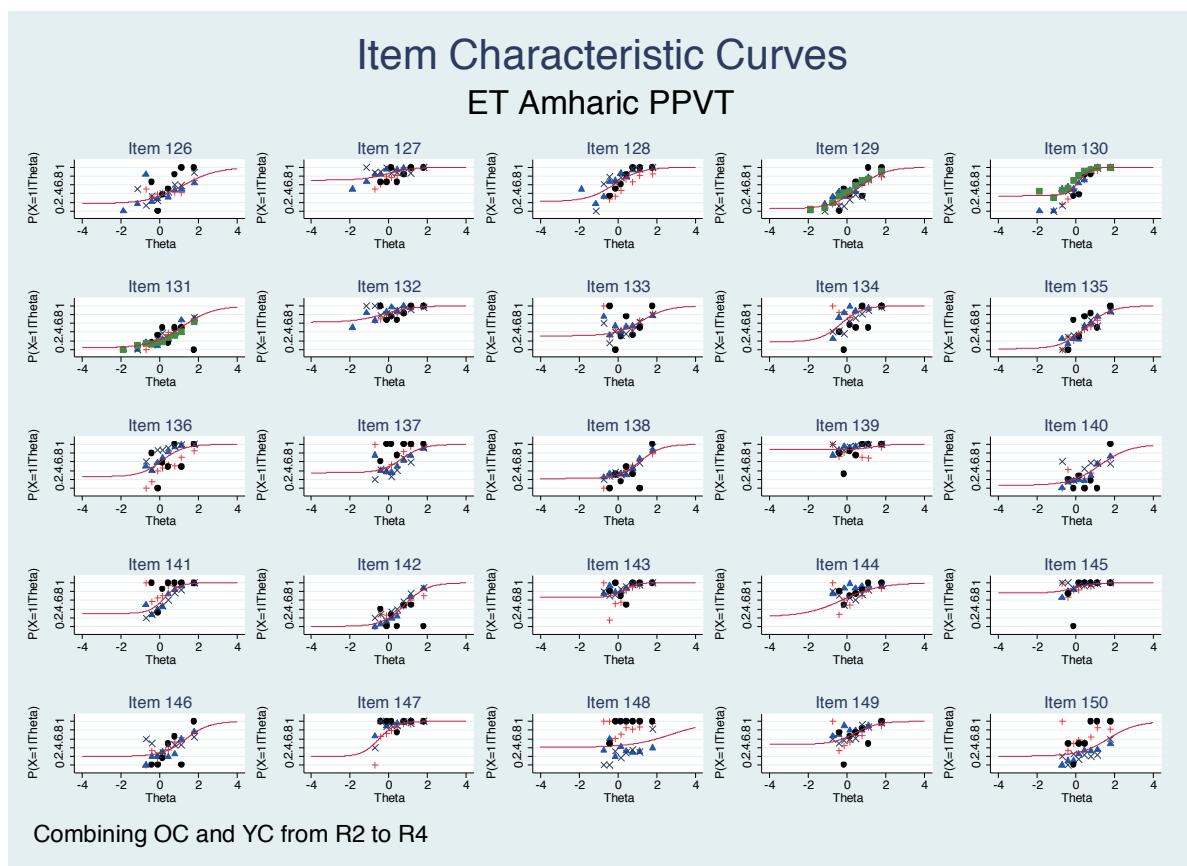
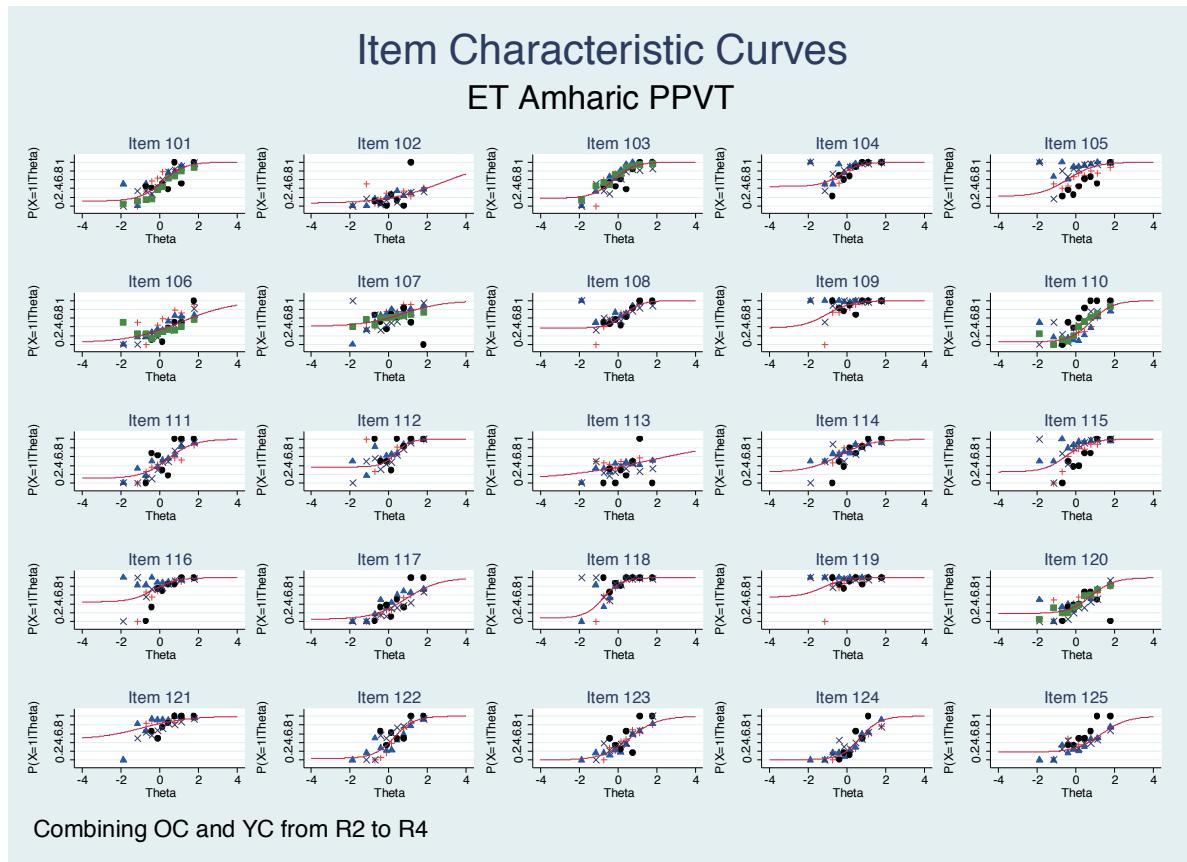


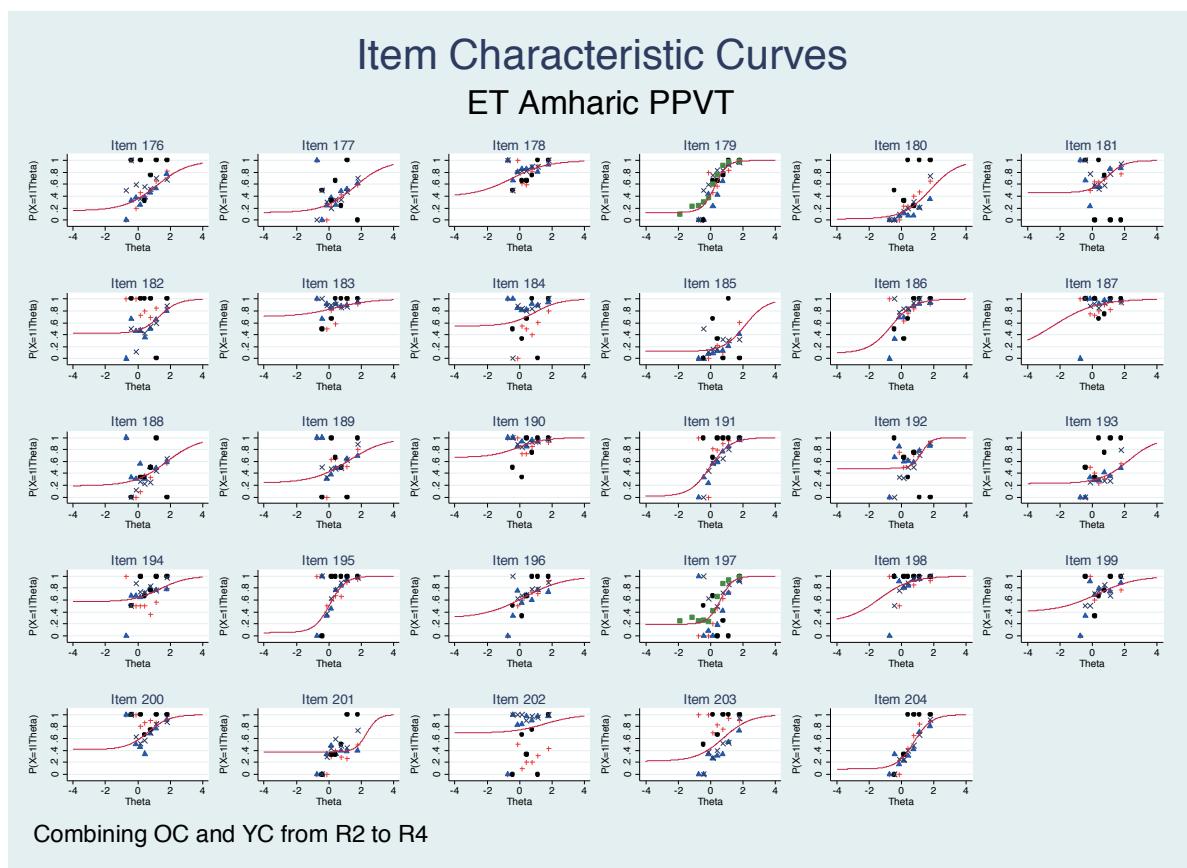
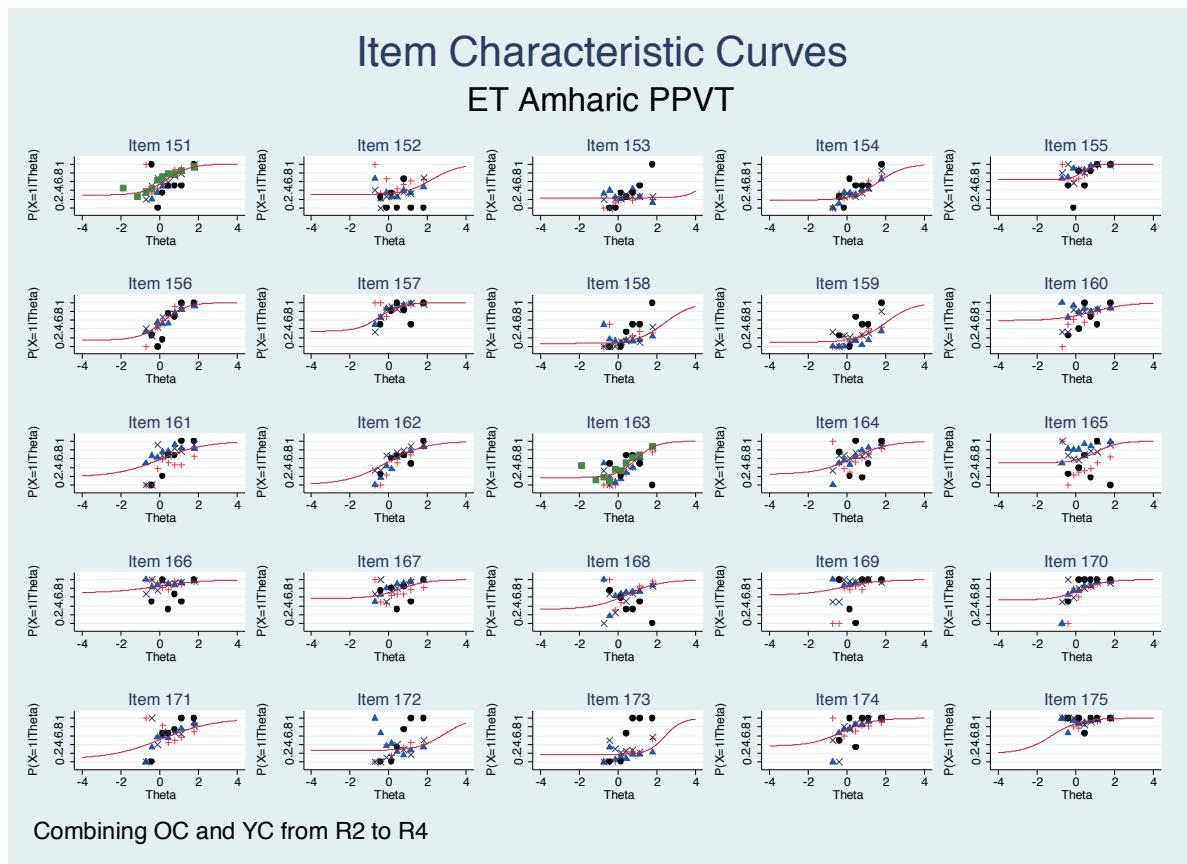
### Appendix C. DIF analysis for each item by country and main language



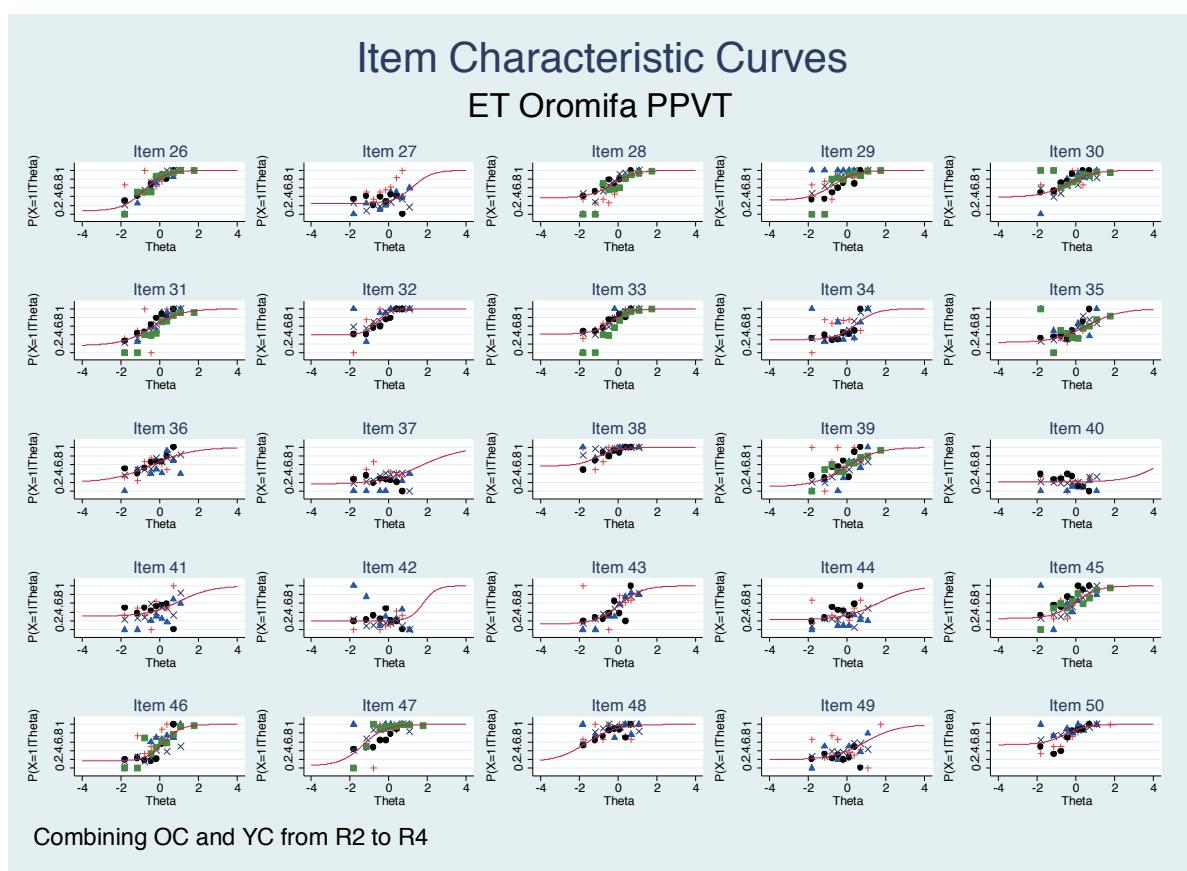
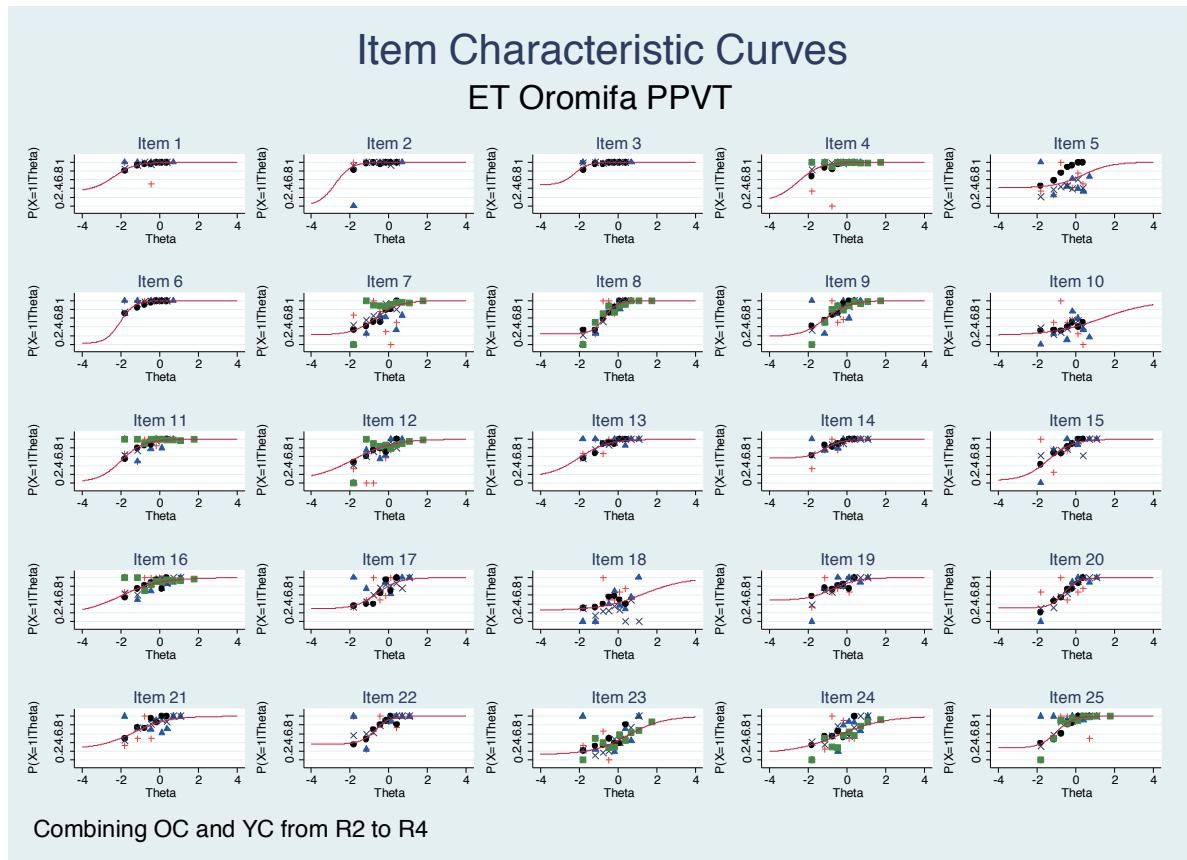


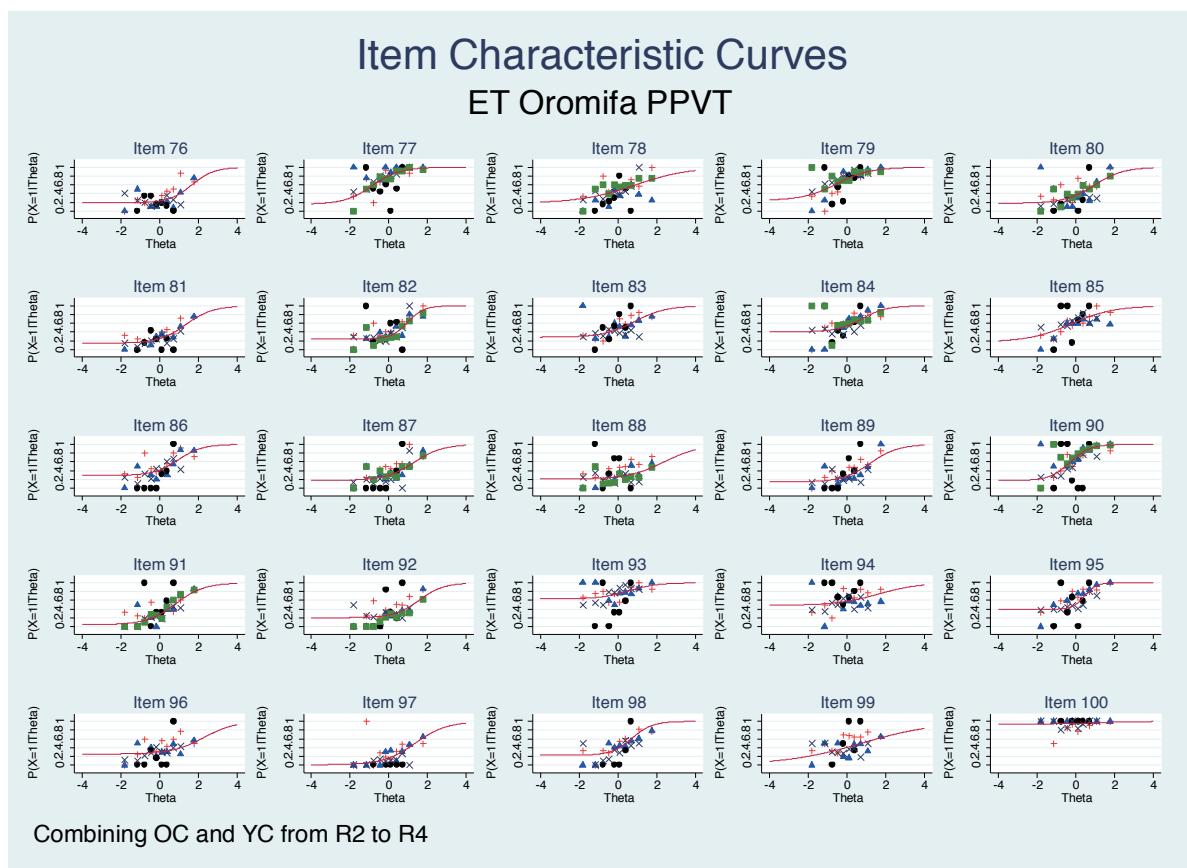
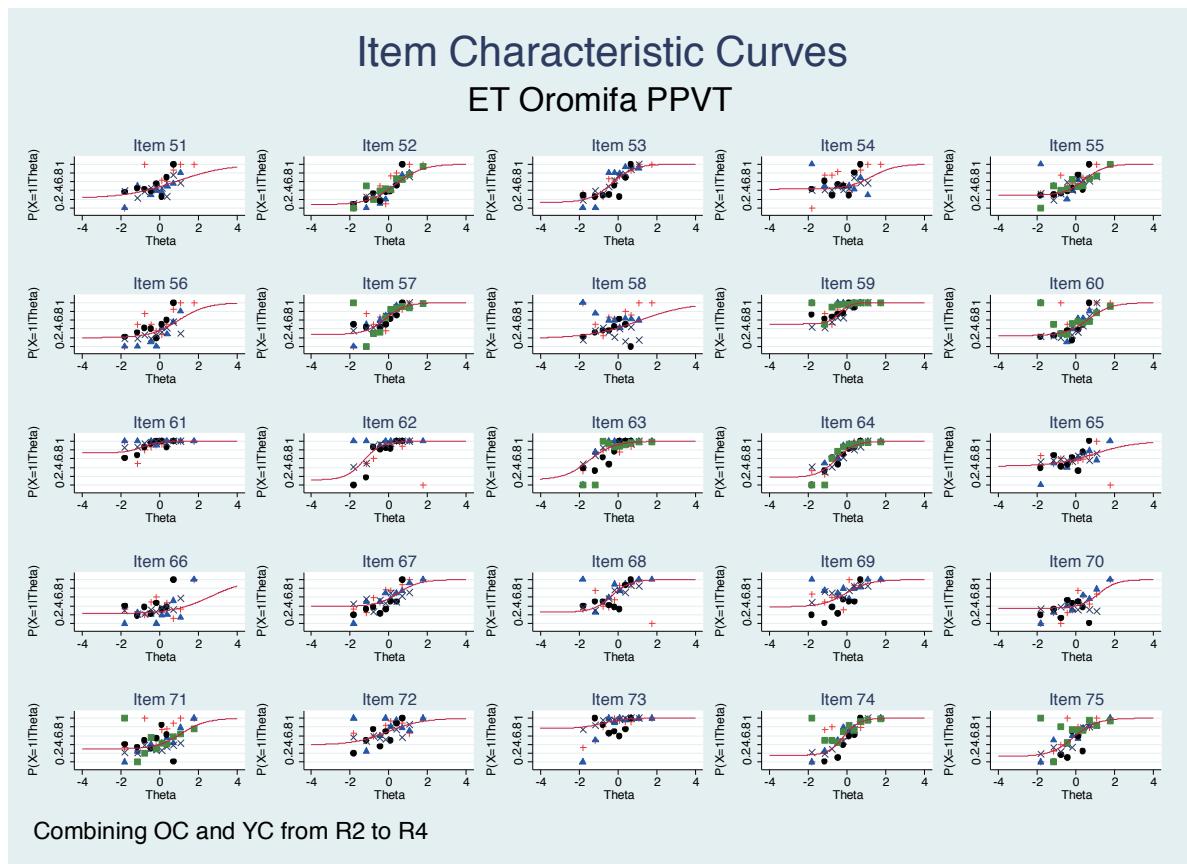
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



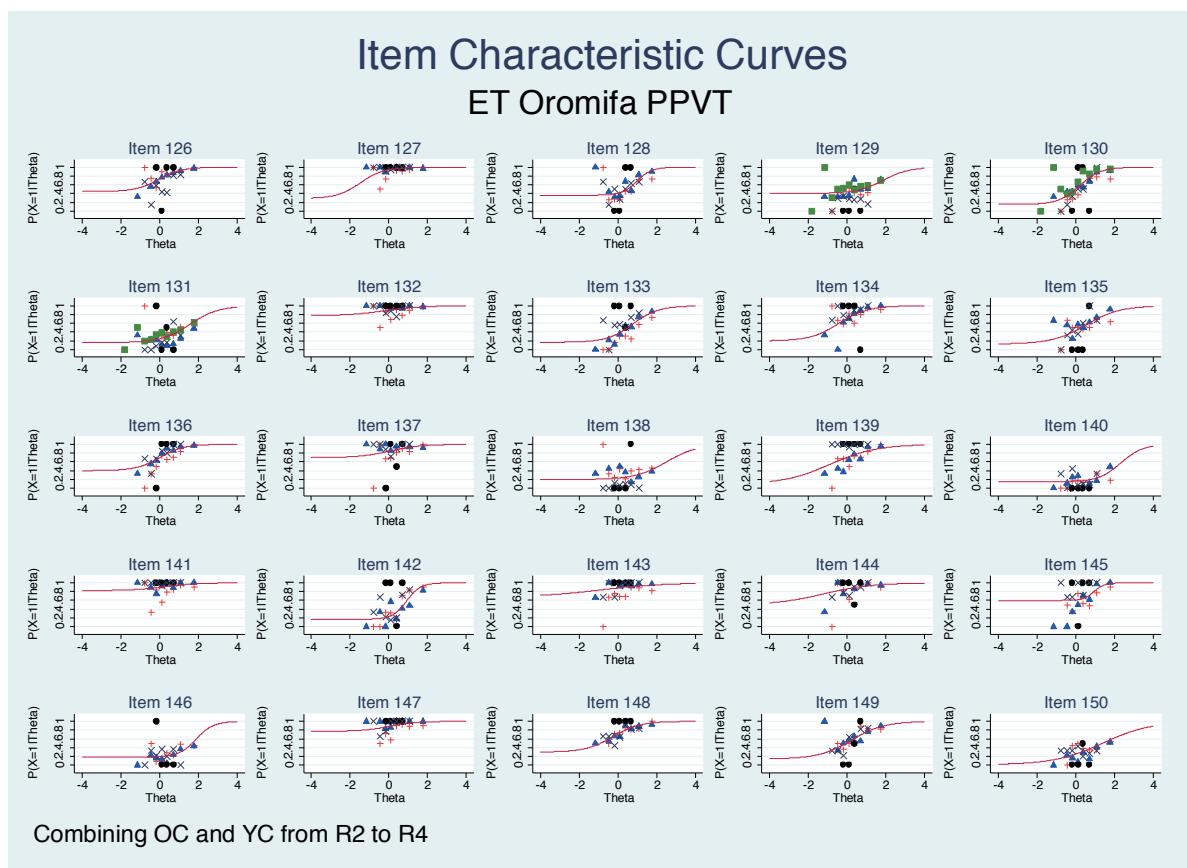
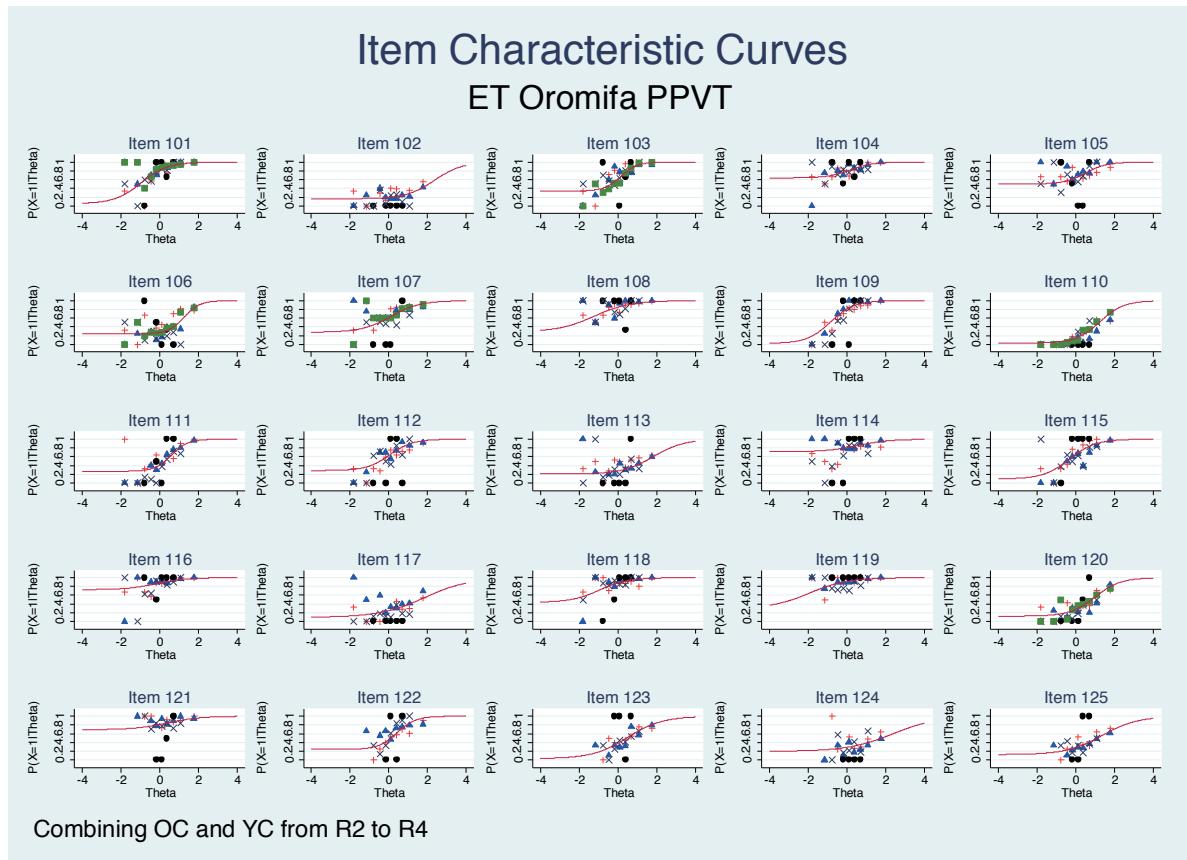


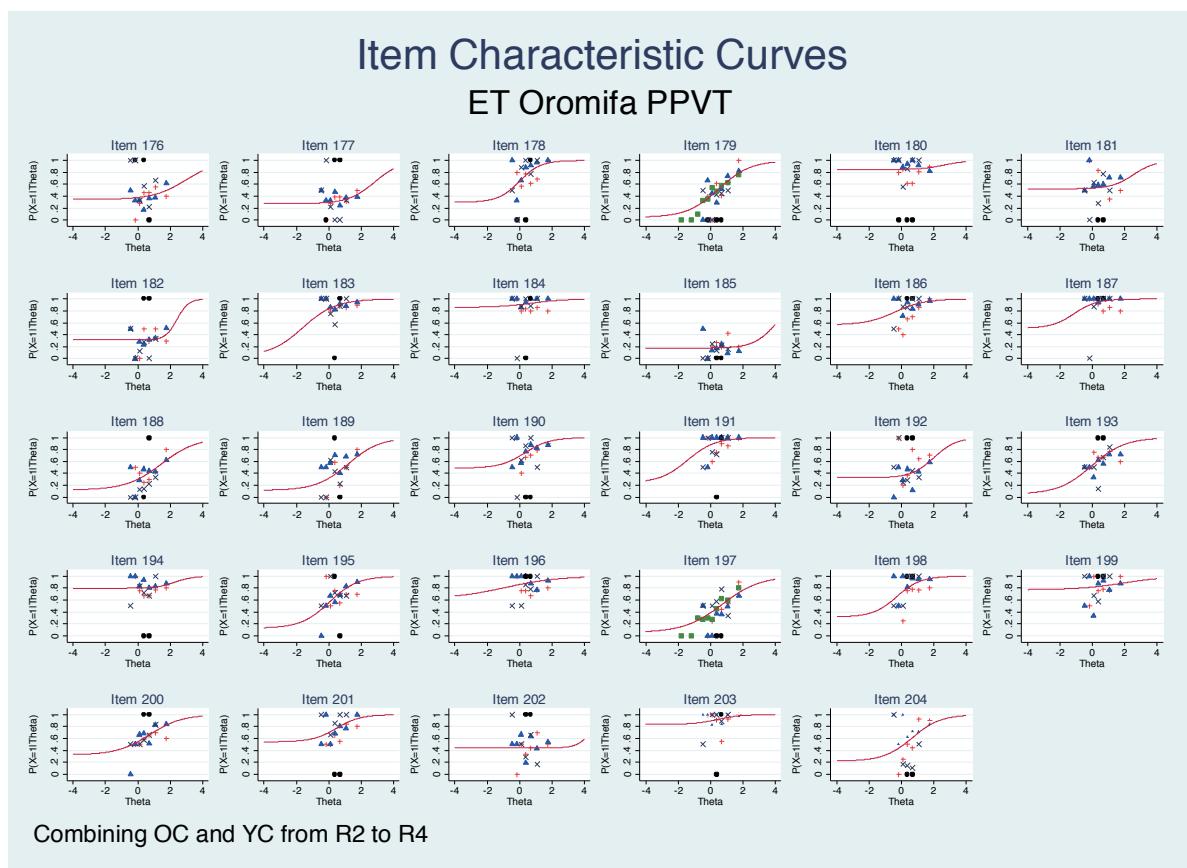
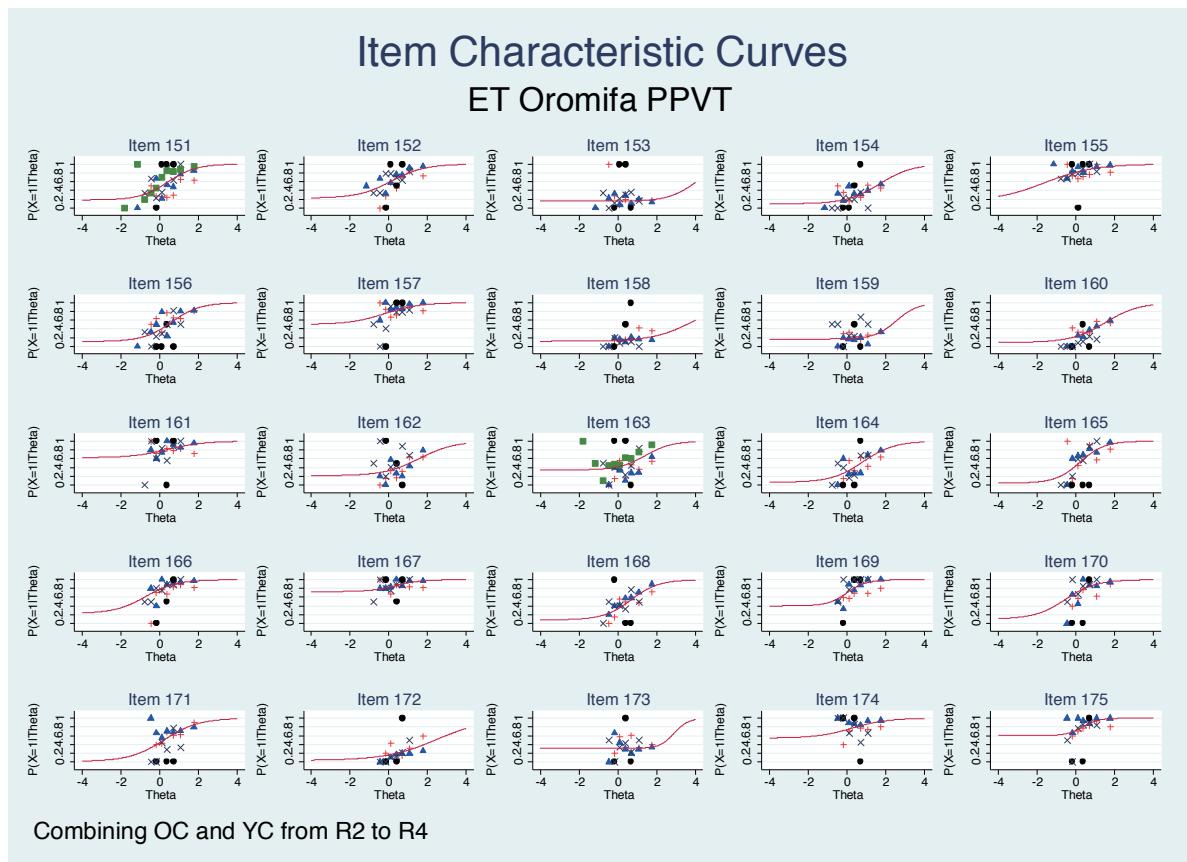
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



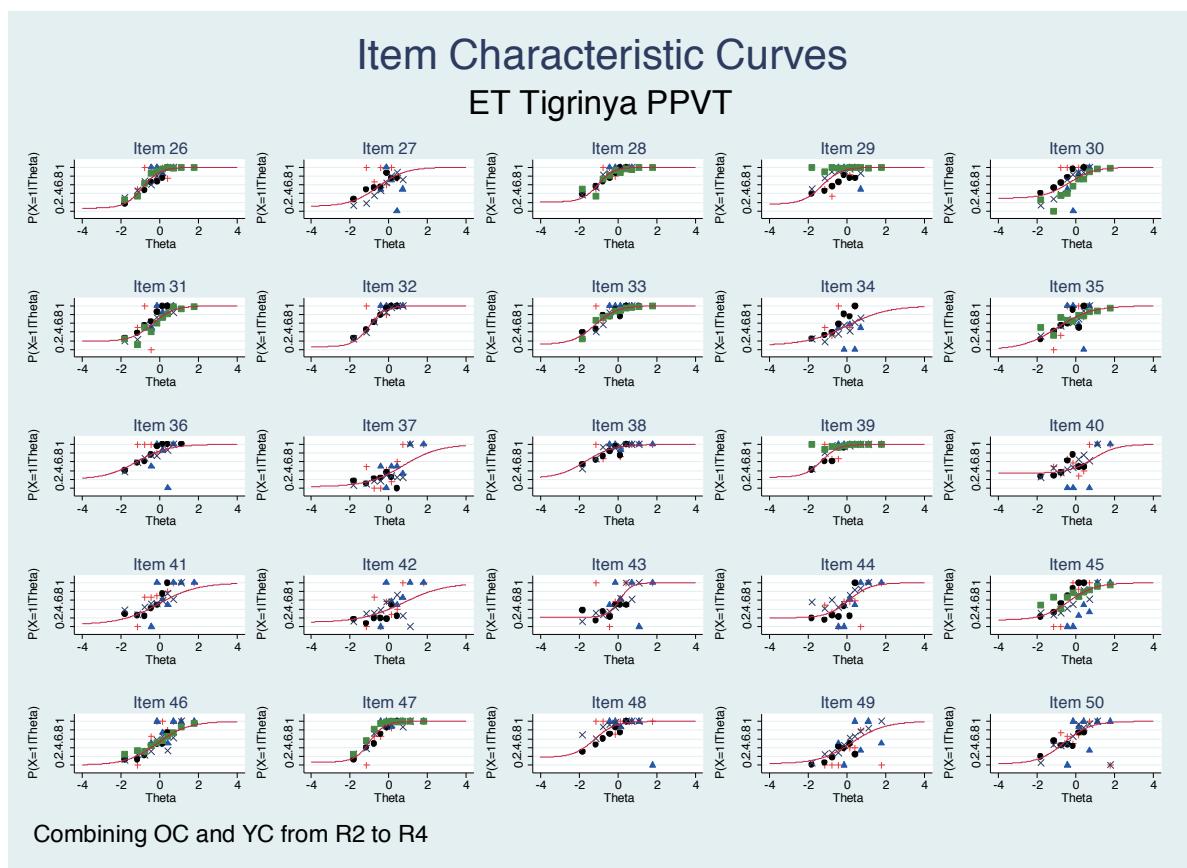
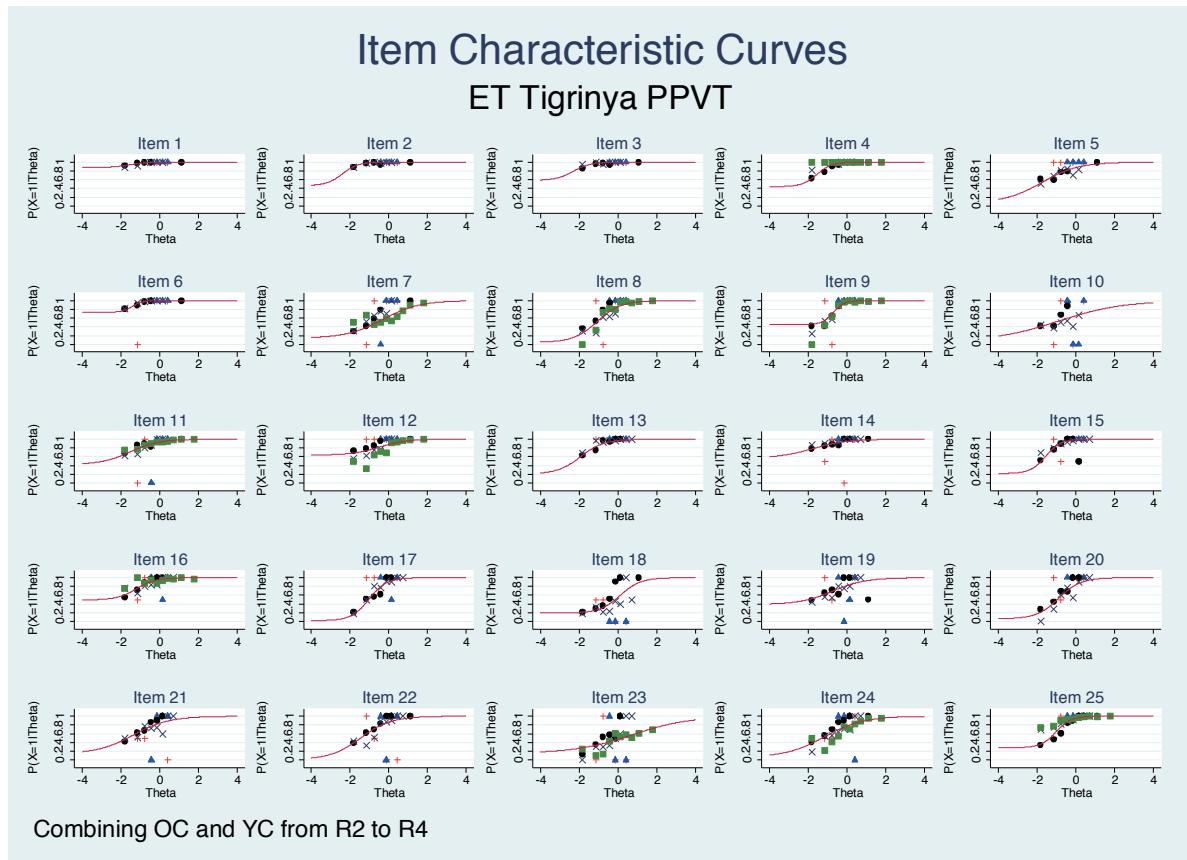


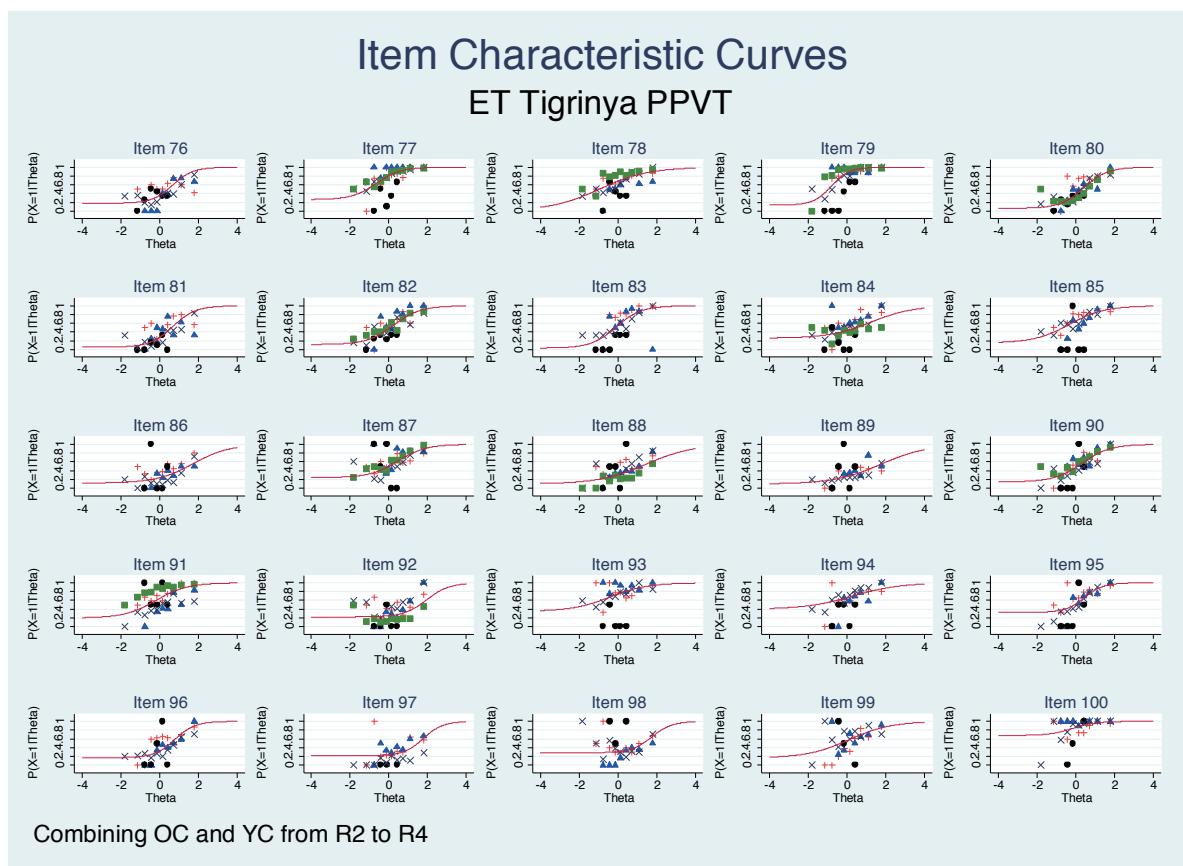
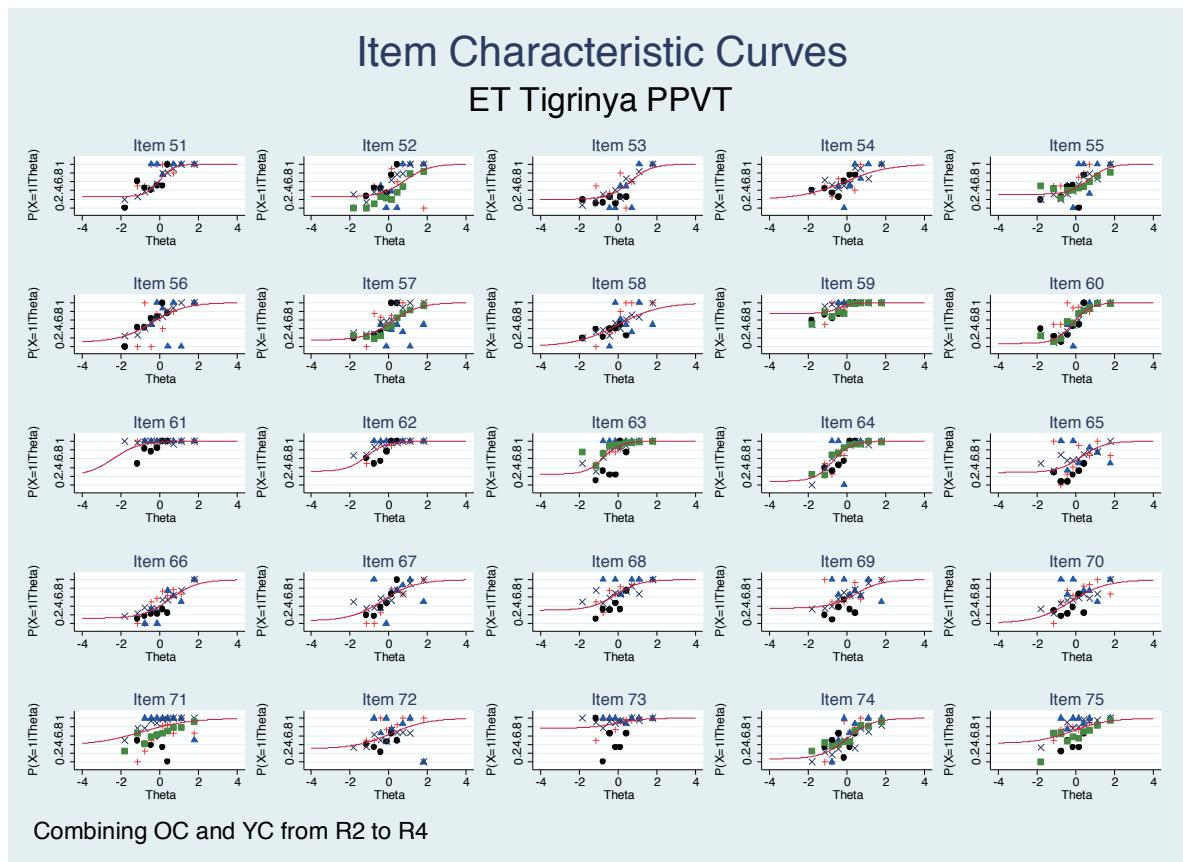
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



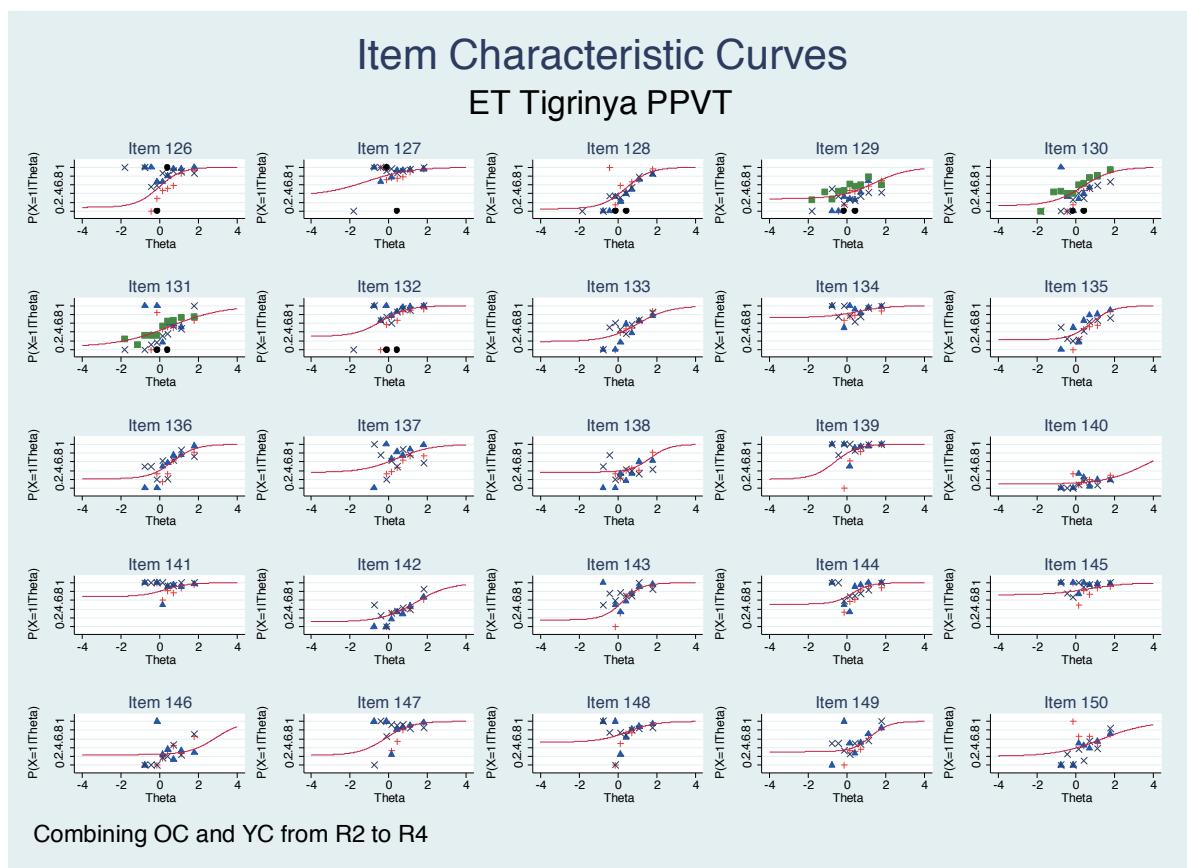
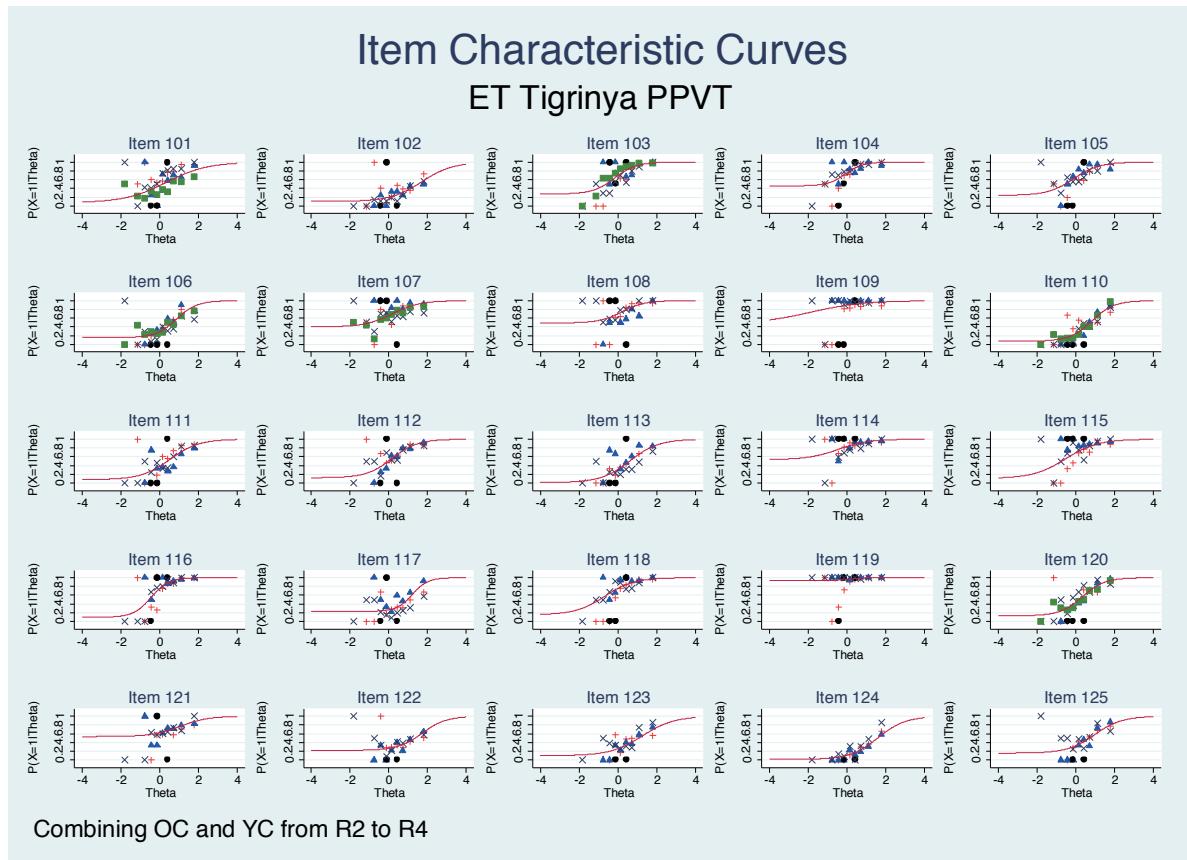


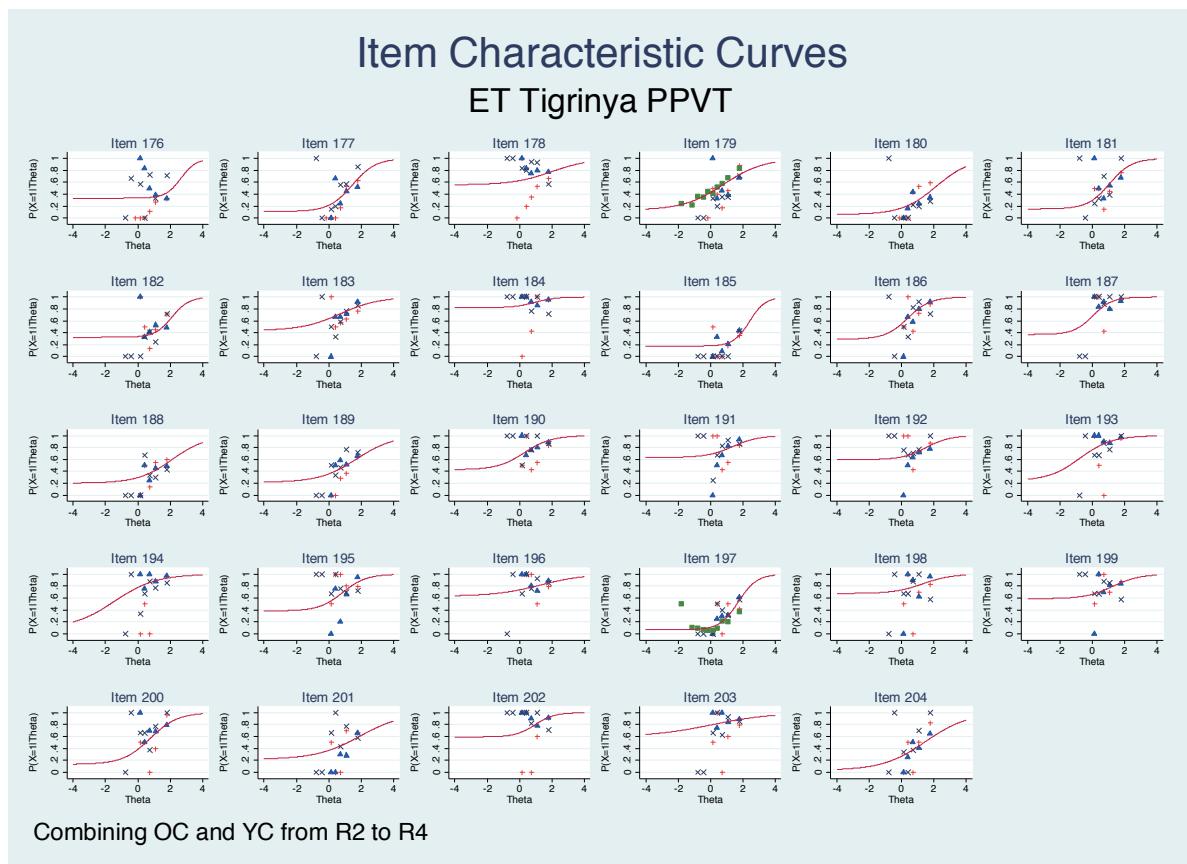
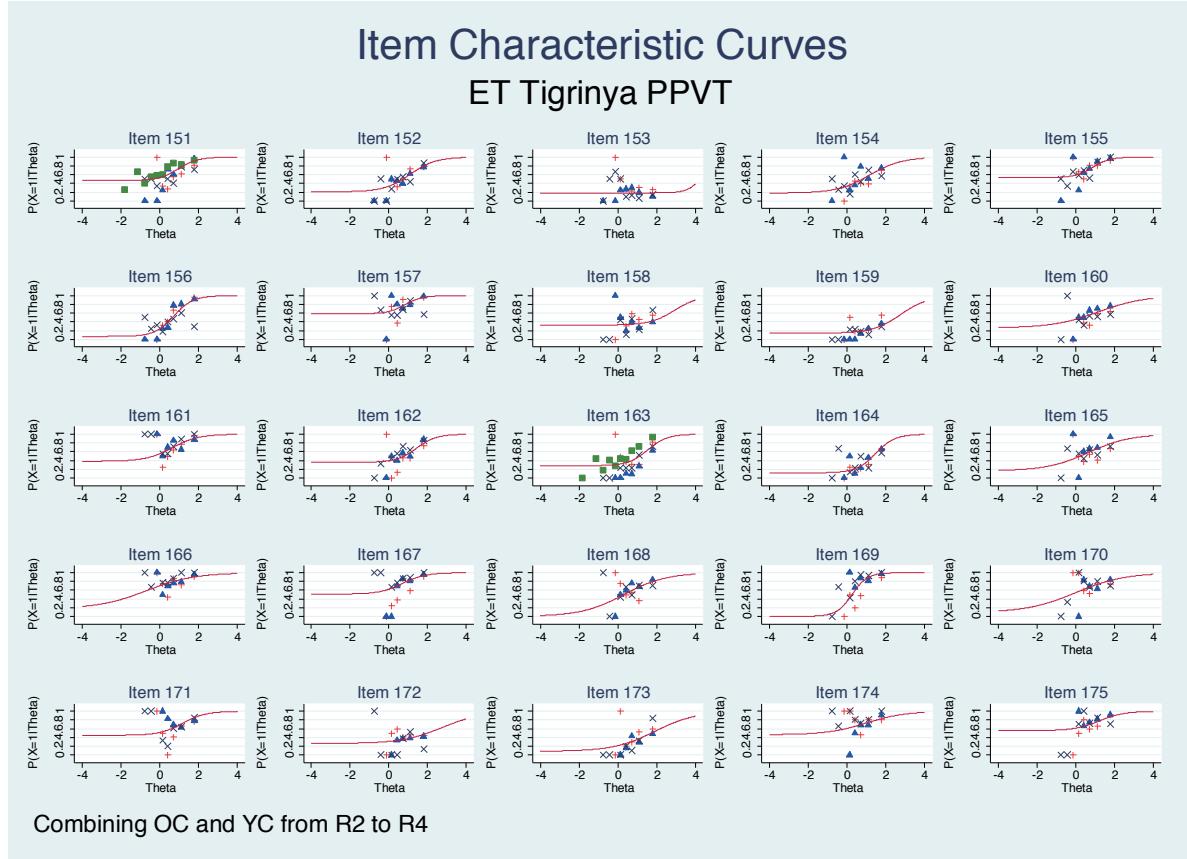
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



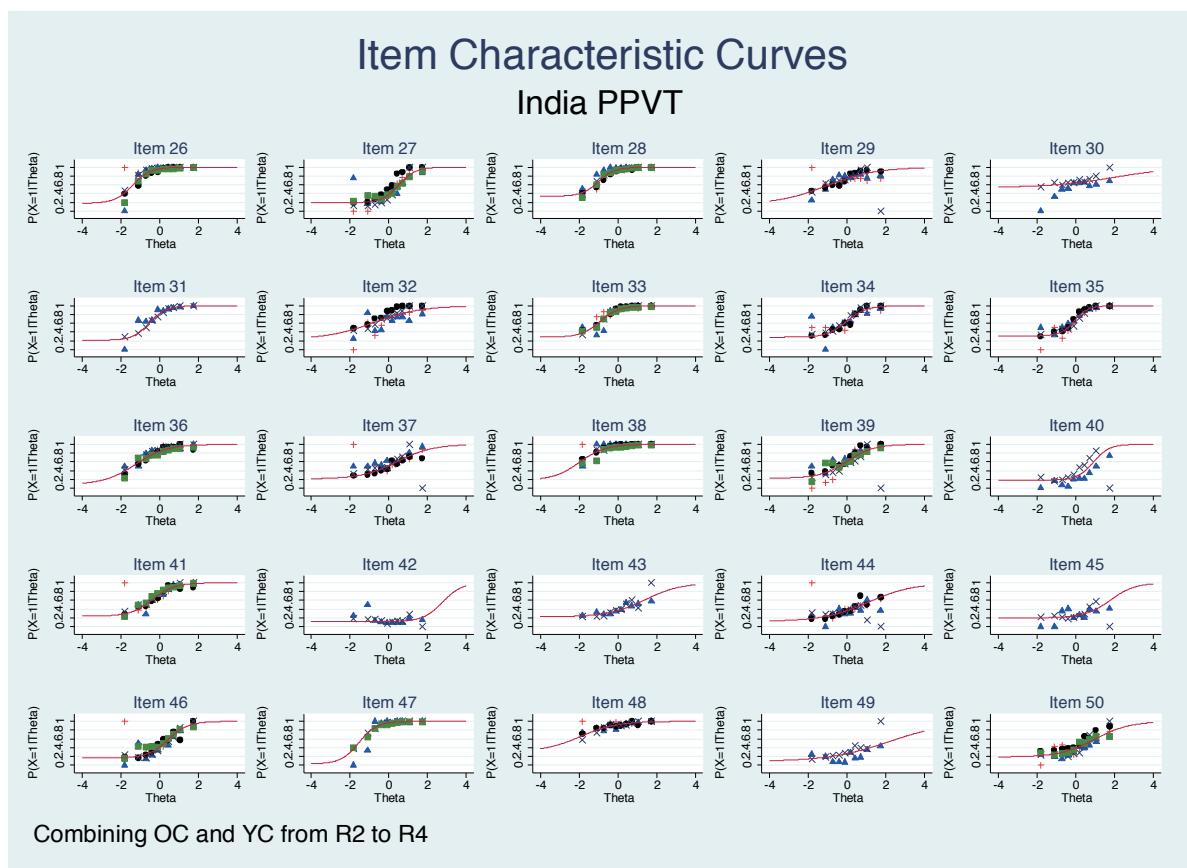
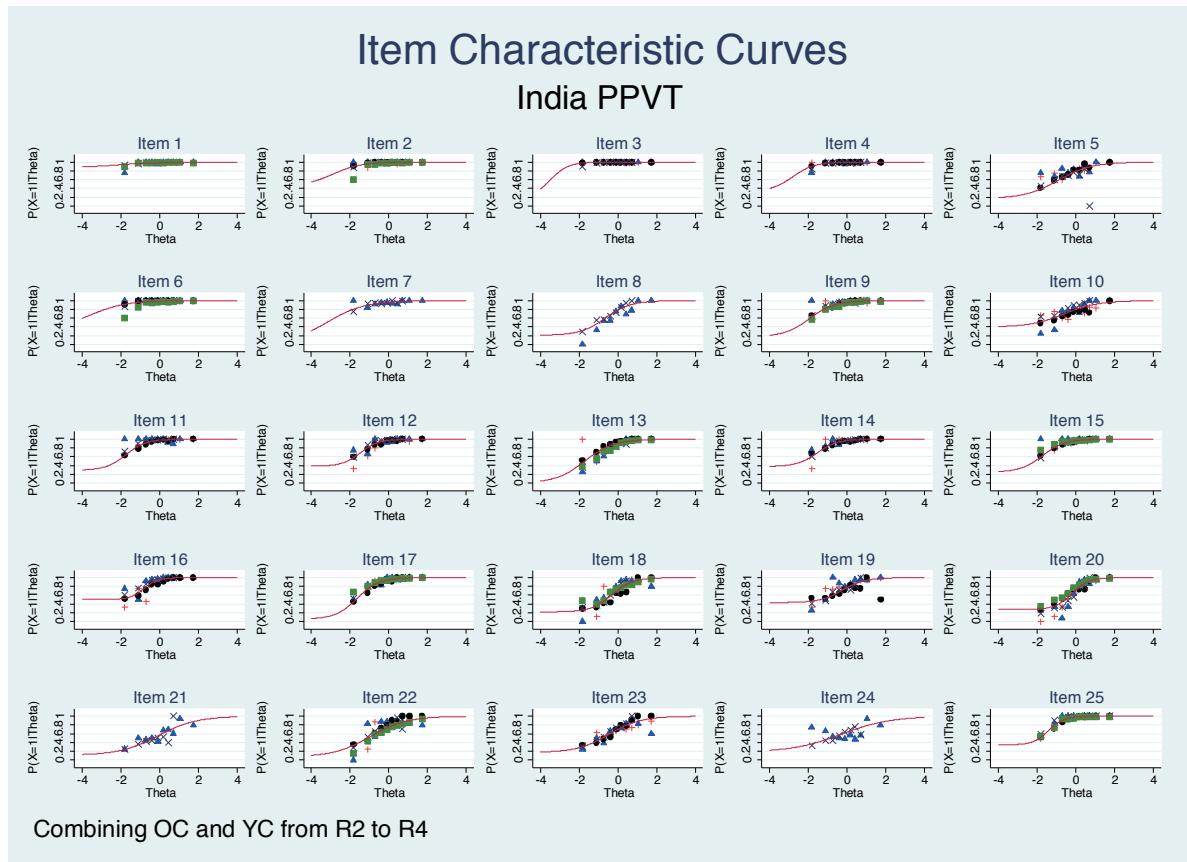


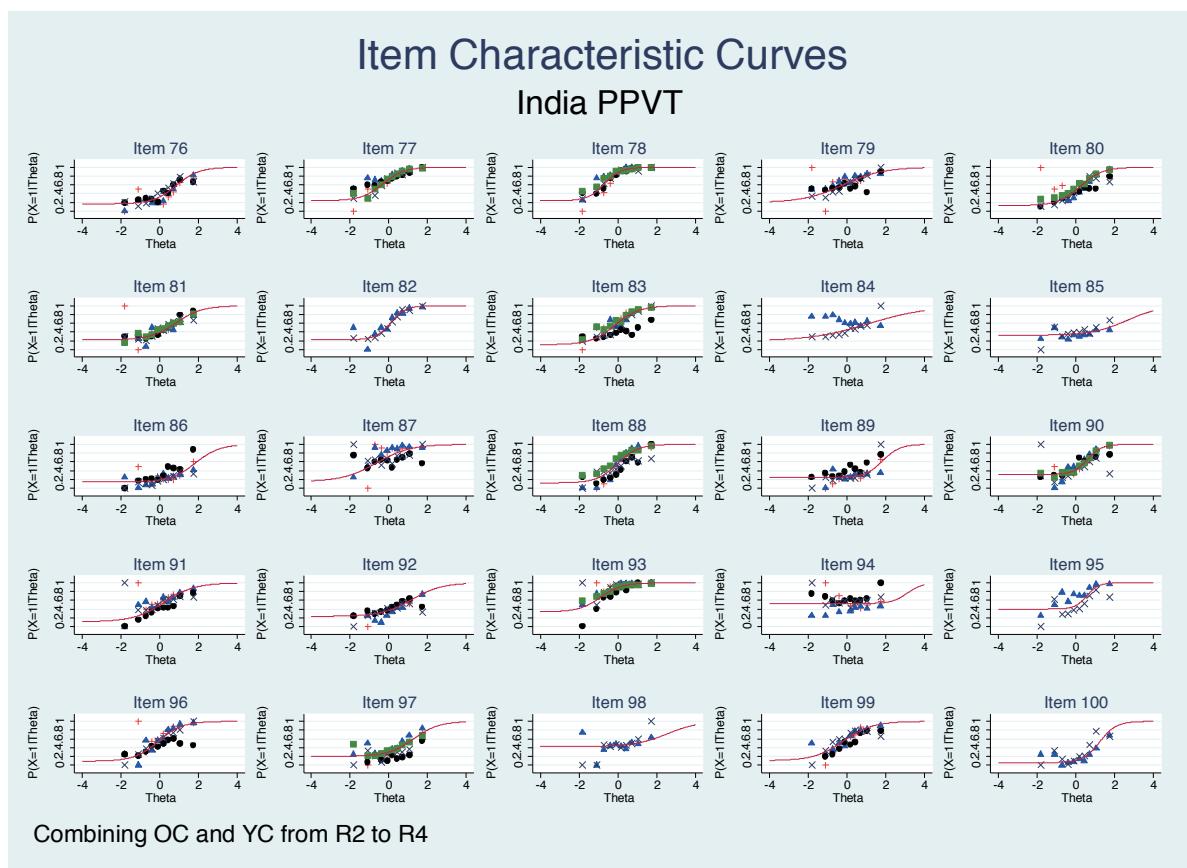
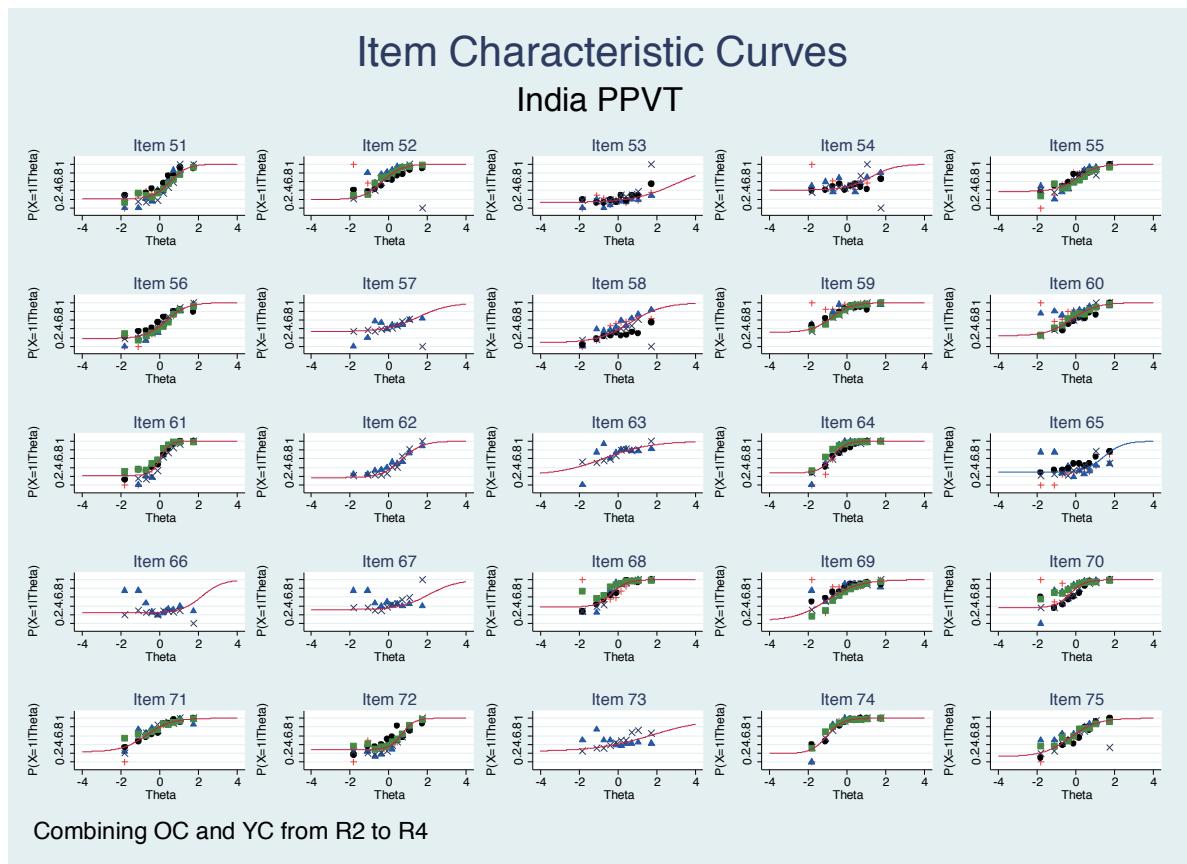
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



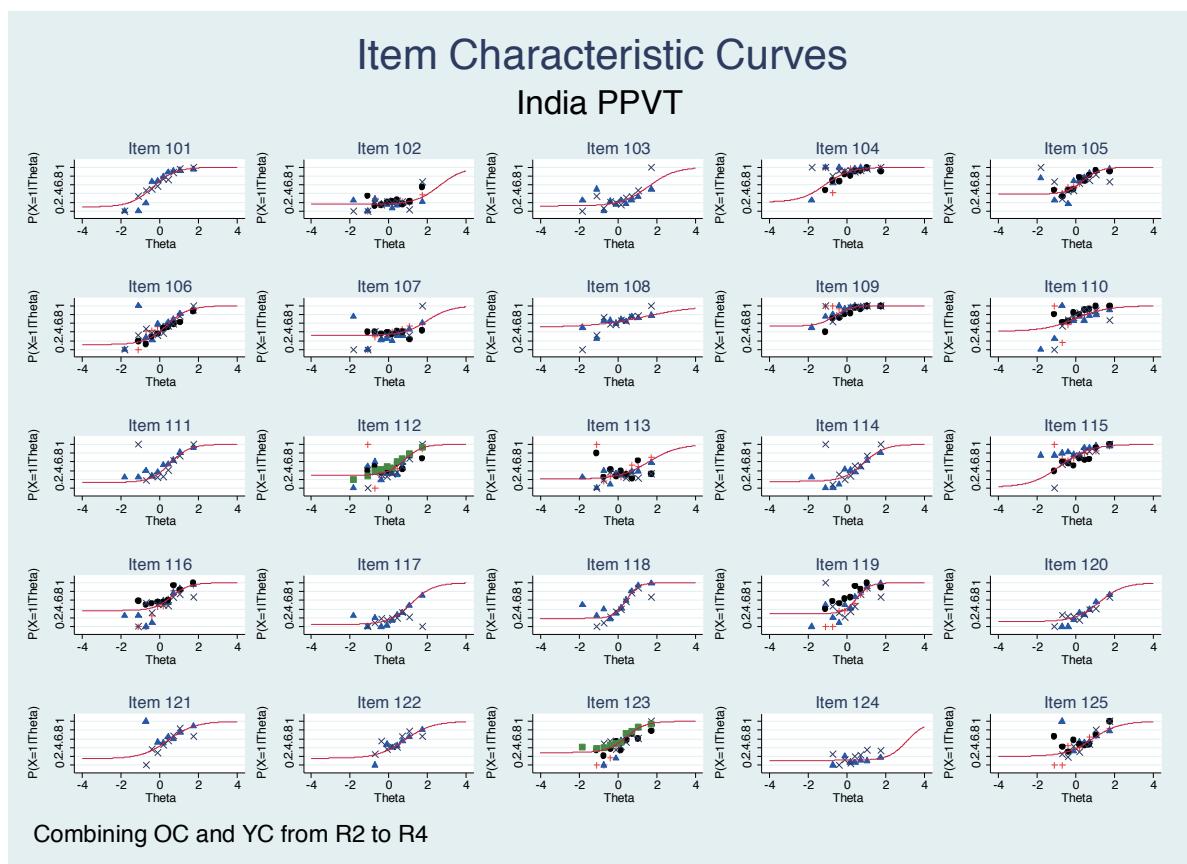
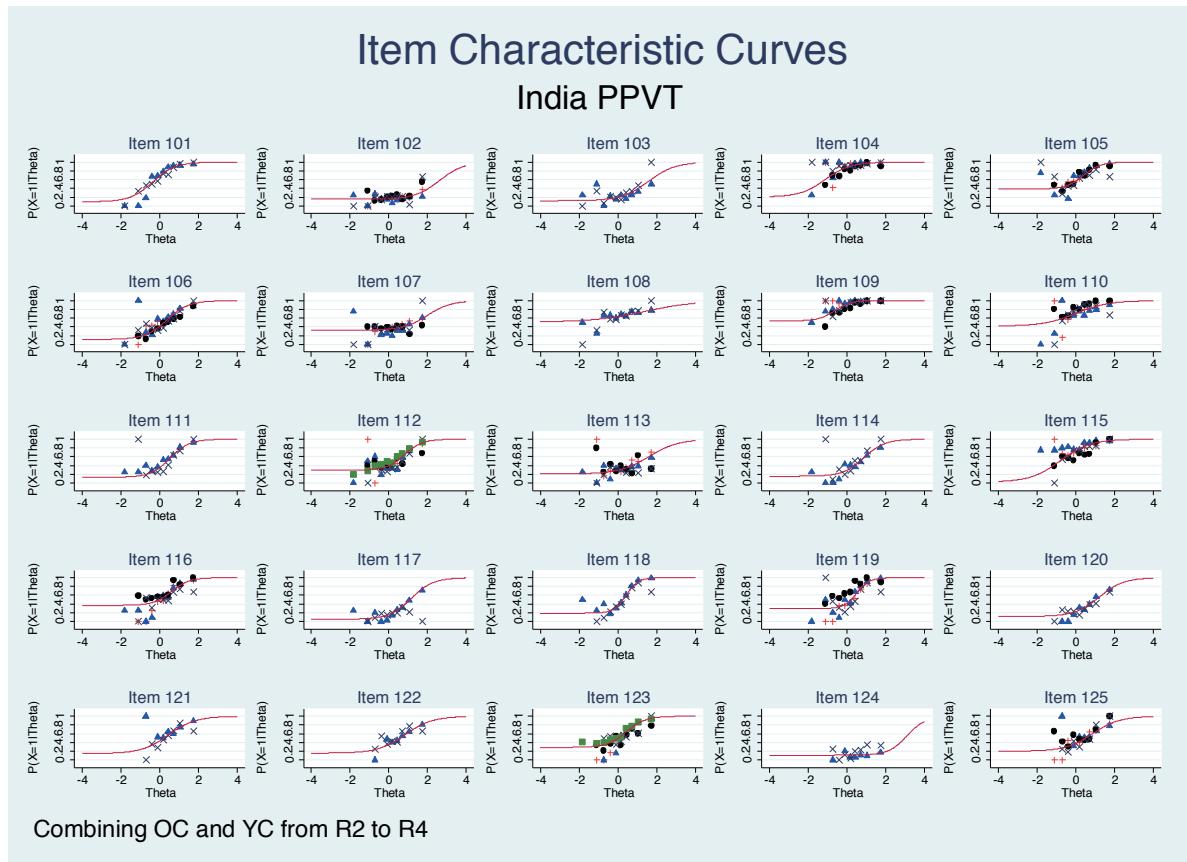


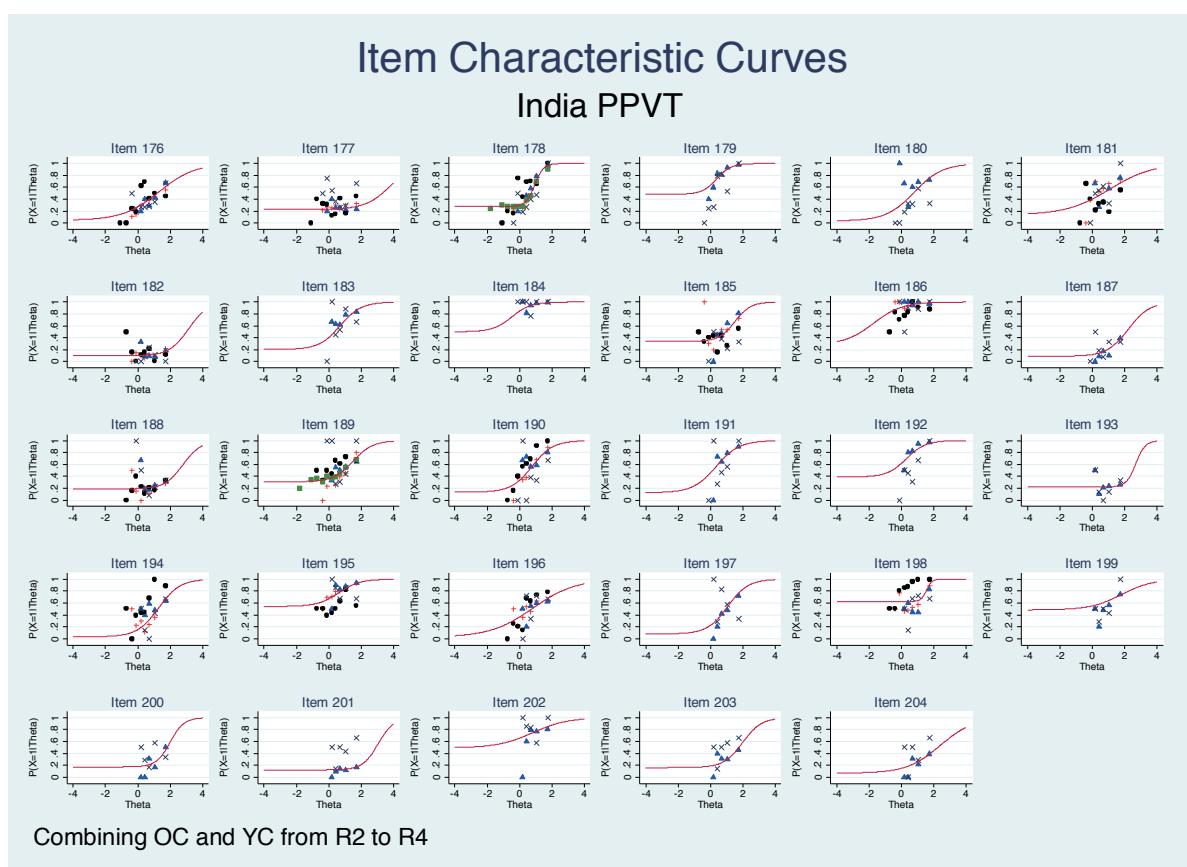
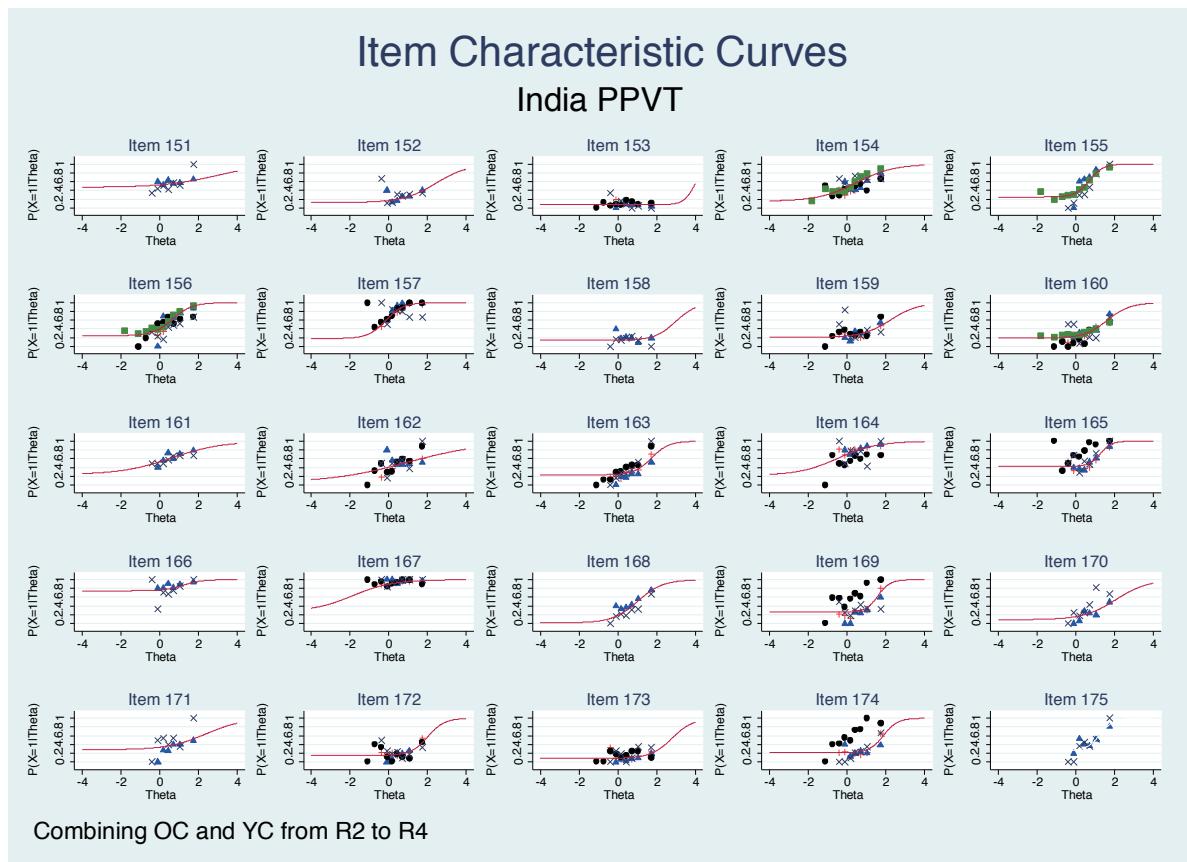
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



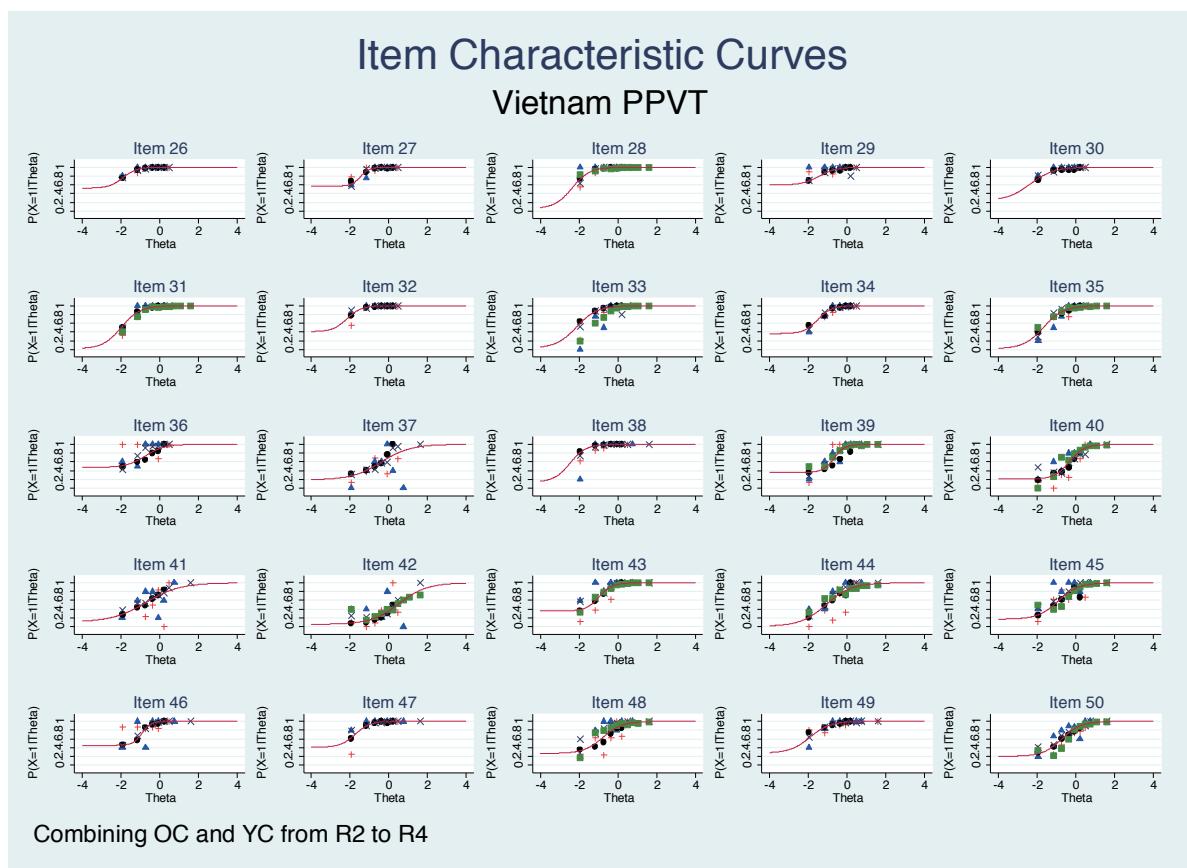
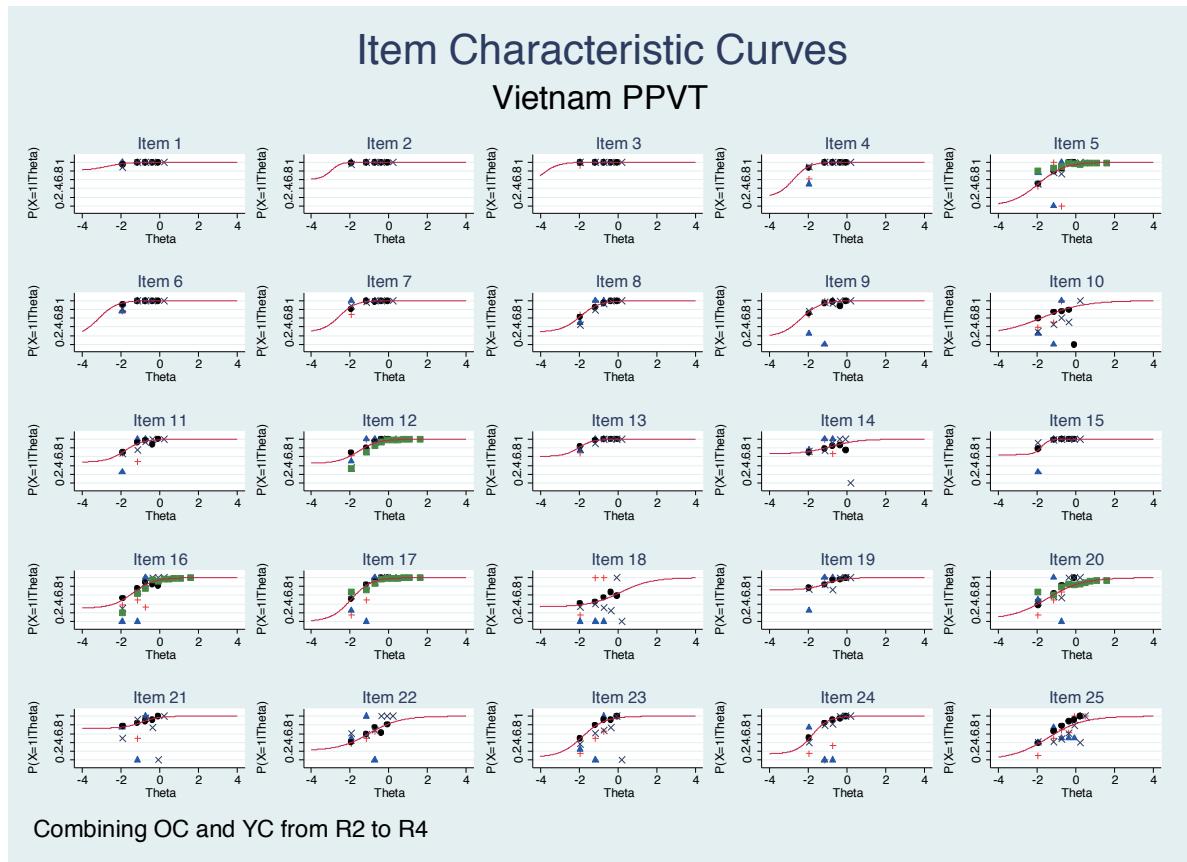


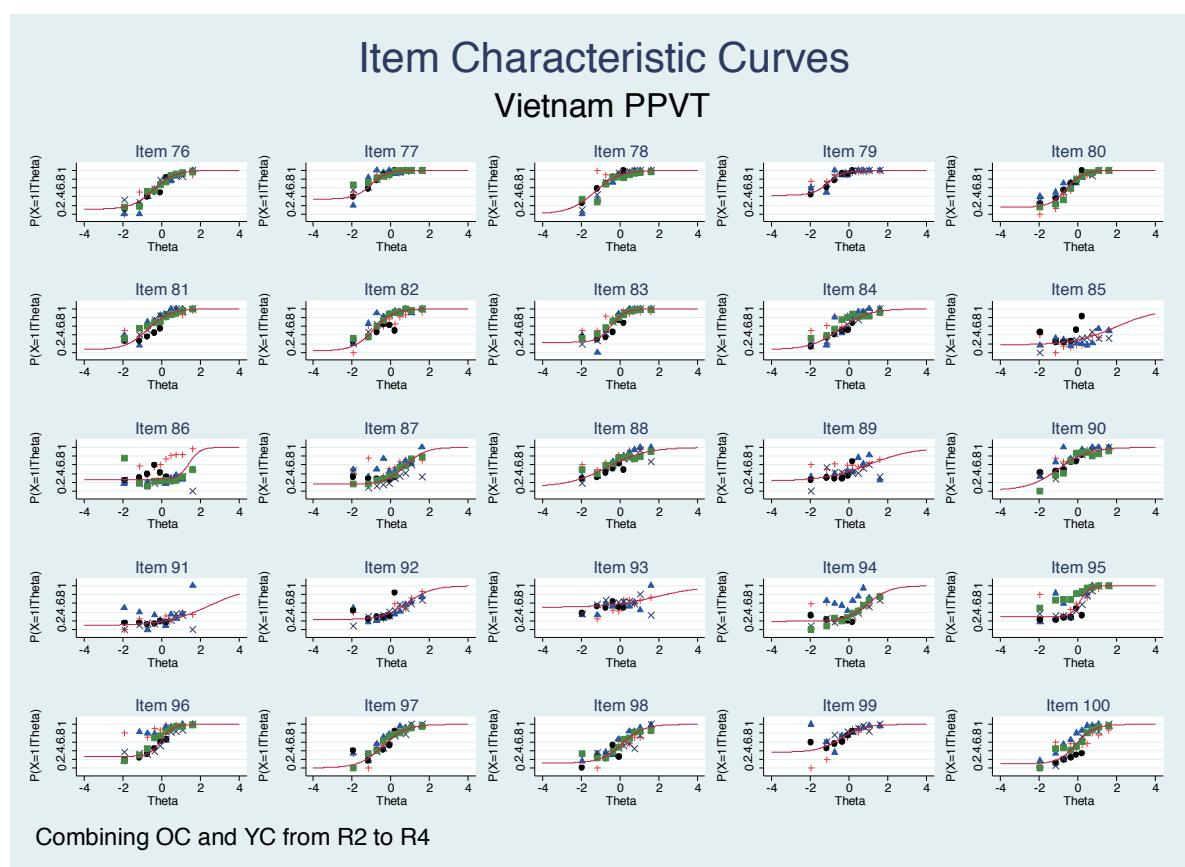
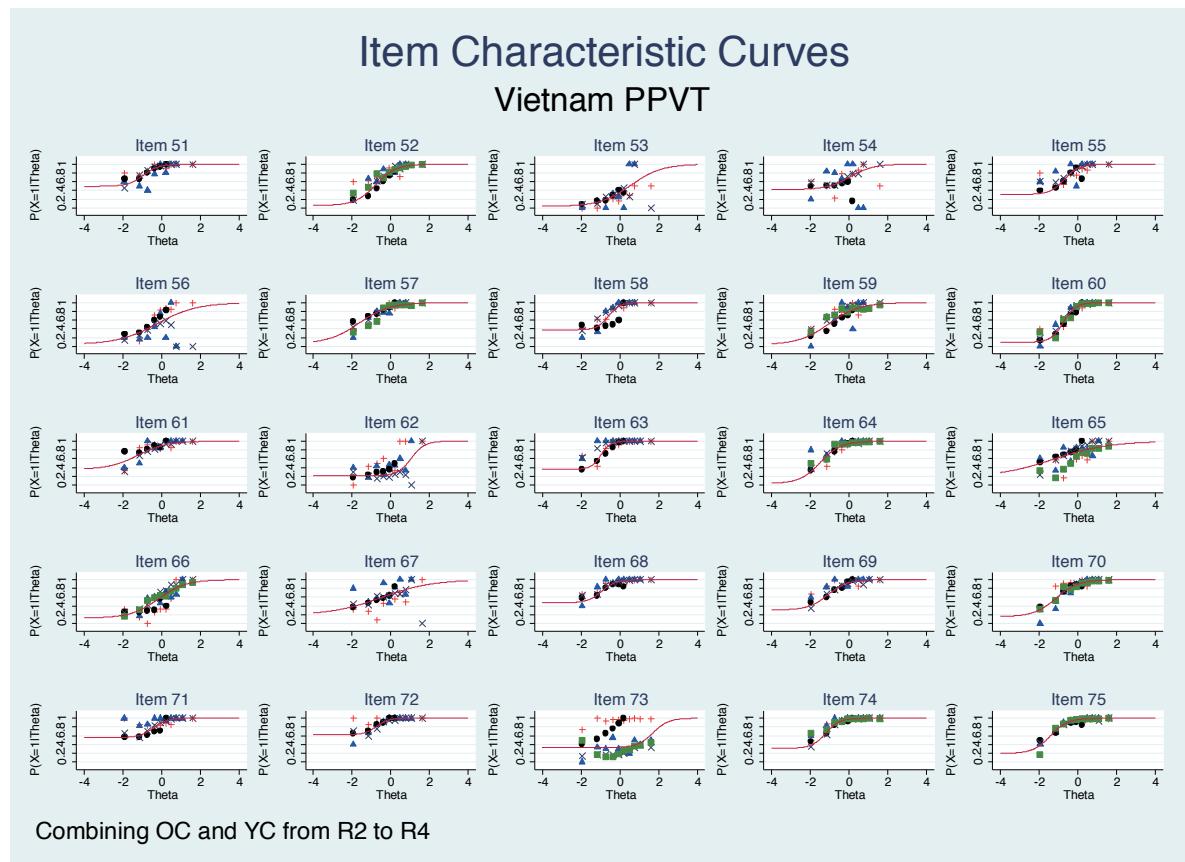
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



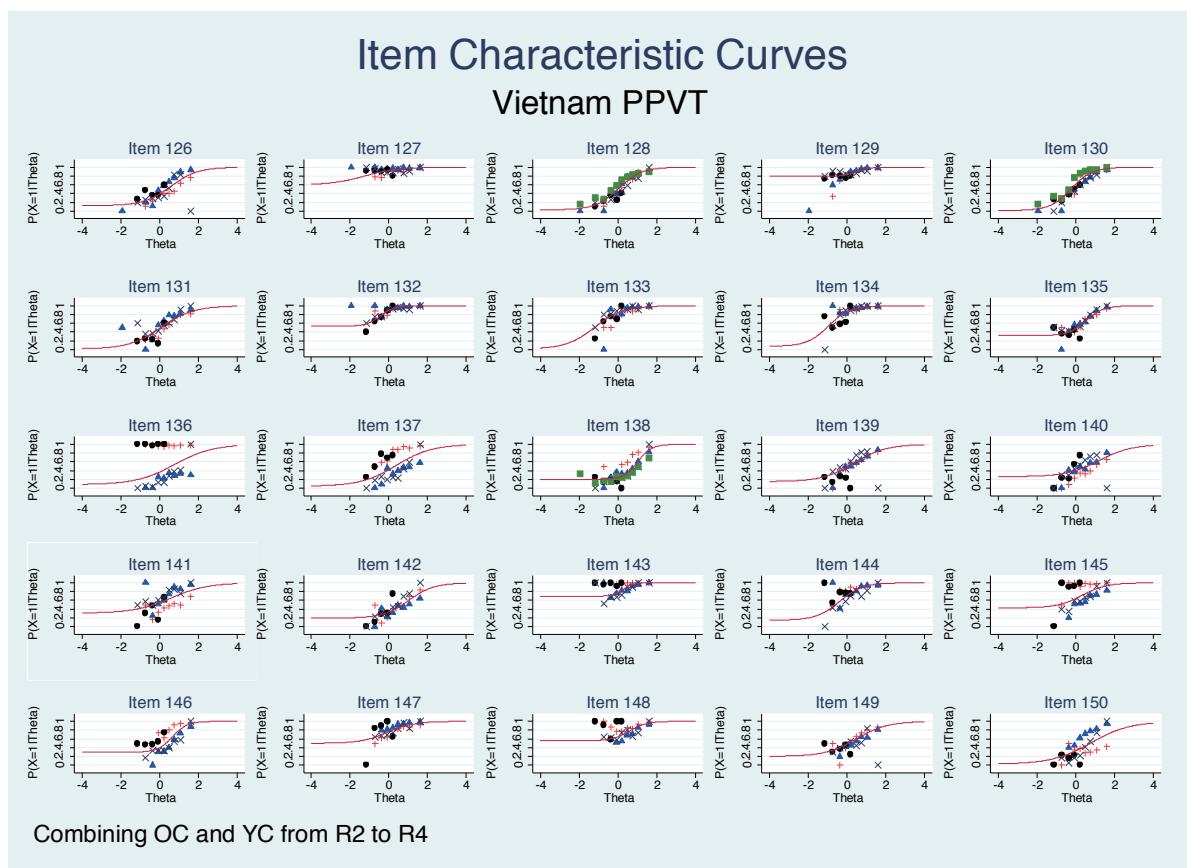
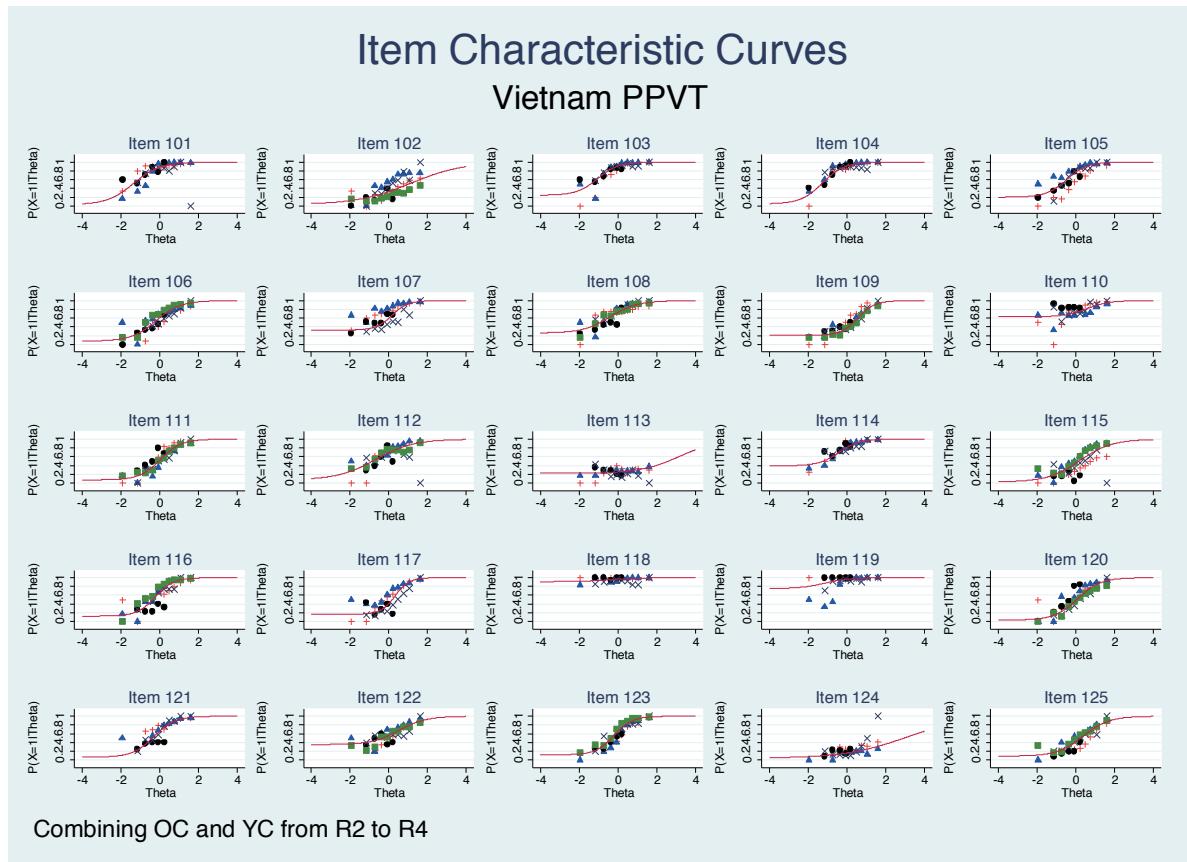


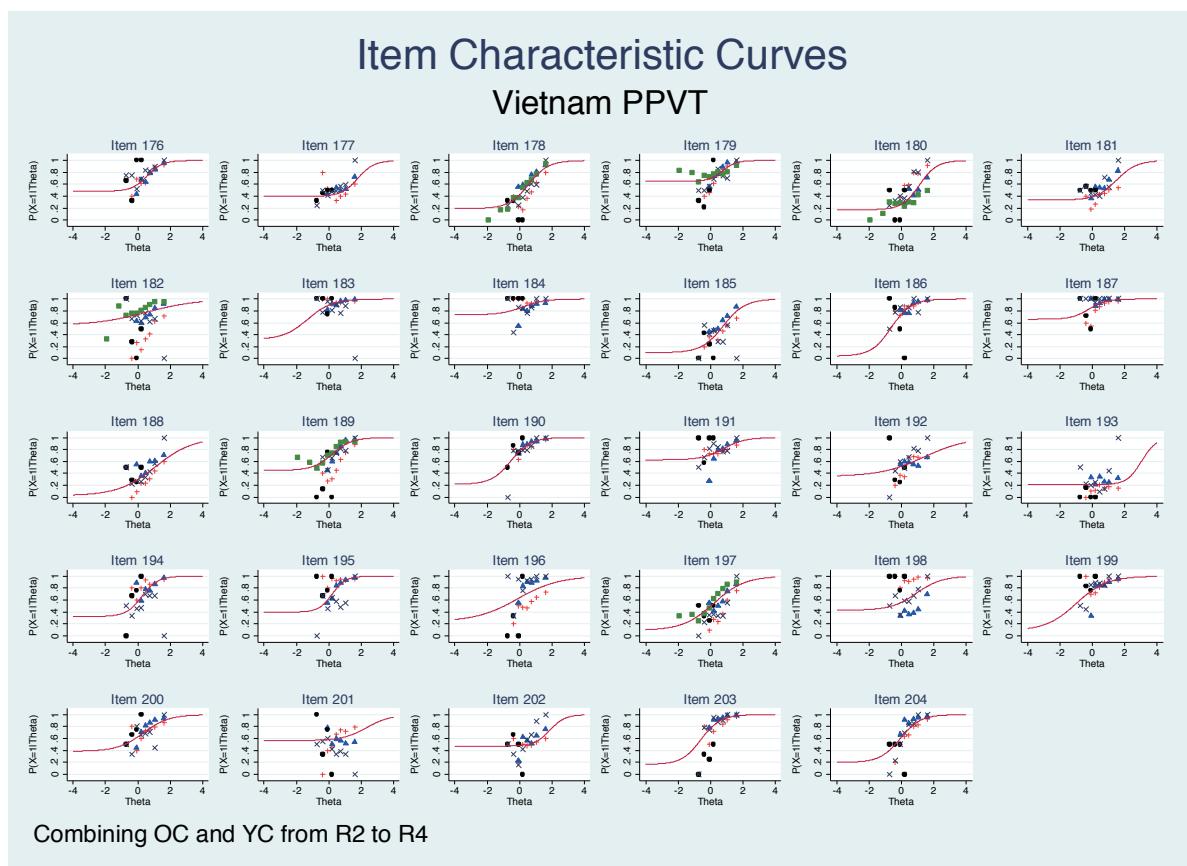
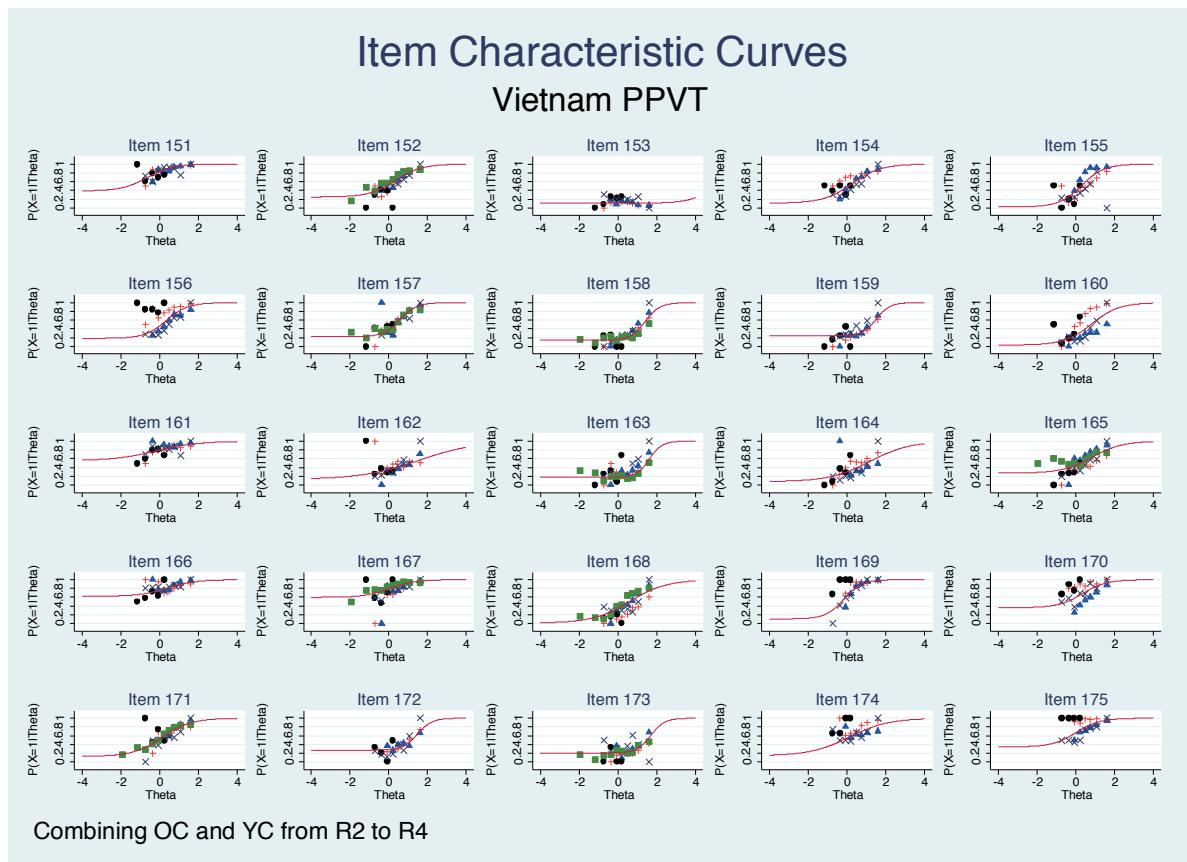
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM





EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM





## Appendix D. Item analysis performed by country and main languages

**Table 1.** Item fit and DIF analysis for Amharic

Item	Anchor For R4	Poor Fit	N					DIF					Deleted
			OC R2	OC R3	YC R2	YC R3	YC R4	OC R2	OC R3	YC R2	YC R3	YC R4	
1			26	40	677	287	0						
2			26	40	677	287	0						
3			26	40	677	287	0						
4	X		26	40	677	287	919						
5			26	40	677	287	0	x		o			
6			26	40	677	287	0						
7	X		26	40	677	287	919			o	o		
8	X		26	40	676	287	919						
9	X		26	40	677	287	919						
10		x	26	40	676	284	0	x	x	x	x		D
11	X		26	40	677	287	919						
12	X		26	40	675	287	919						
13			45	43	812	356	0						
14			45	43	811	356	0						
15			45	43	811	356	0						
16	X		45	43	811	356	919		o				
17			45	43	812	356	0						
18			45	43	812	356	0	x	x				
19			45	43	812	356	0		x				
20			45	43	811	356	0		o	x			
21			45	43	812	356	0		x				
22		x	45	43	808	356	0	x	x	x	x		D
23	X		45	43	812	356	919						
24	X		45	43	810	356	919						
25	X		93	91	854	513	919						
26	X		93	91	852	513	919						
27			93	91	854	512	0	x	o	o	o		
28	X		93	91	855	513	919						
29	X		93	91	850	513	919						
30	X		93	91	852	513	919		o				
31	X		93	91	855	513	919						
32			93	91	854	513	0						
33	X		93	91	853	513	919						
34			91	91	851	513	0						
35	X		93	91	854	513	919						
36			93	91	855	513	0	x					
37			110	101	630	573	0						
38			112	101	631	573	0						
39	X		111	101	630	573	919						
40			112	101	627	573	0						
41			112	101	631	573	0			o			
42		x	111	101	628	569	0	x	x	x	x		D
43			109	100	625	567	0	x					
44			111	101	630	573	0		x				
45	X		112	101	631	573	919						
46	X		112	101	631	572	919						
47	X		112	101	631	573	919						
48			111	101	631	573	0						
49			128	117	385	649	0			x			
50			130	117	388	653	0		x				





EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

Item	Anchor For R4	Poor Fit	N					DIF					Deleted
			OC R2	OC R3	YC R2	YC R3	YC R4	OC R2	OC R3	YC R2	YC R3	YC R4	
159			268	369	21	284	0				o		
160		x	268	369	21	284	0	x	x	x	x		D
161		x	268	369	21	284	0	x	x	x	x		D
162			268	369	21	284	0						
163	X		268	369	21	284	919			x			
164			268	369	21	284	0			x			
165			268	369	21	284	0			x			
166		x	268	369	21	284	0	x	x	x	x		D
167		x	267	369	21	284	0	x	x	x	x		D
168			265	369	21	282	0			x			
169		x	226	357	14	270	0	x	x	x	x		D
170			226	357	14	270	0			x			
171			226	357	14	270	0			x			
172			226	357	14	270	0		x				
173			225	357	14	268	0			x			
174			226	357	14	270	0						
175	x		226	357	14	270	0	x	x	x	x		D
176			226	357	14	270	0			x			
177			226	357	14	270	0			x			
178			225	357	14	270	0			x			
179	X		224	357	14	270	919						
180			226	357	14	270	0			x			
181			205	345	14	259	0			x			
182			205	345	14	259	0	x		x			
183			205	345	14	259	0			x			
184			205	345	14	259	0	x		x			
185	x		190	345	13	259	0	x	x	x	x		D
186			205	345	14	259	0			x	x		
187	x		205	345	14	259	0	x	x	x	x		D
188			202	345	13	258	0		x				
189			205	345	14	259	0			x			
190			205	345	14	259	0			x			
191			205	345	14	259	0	x		x			
192			200	344	14	259	0			x			
193			193	342	14	249	0			x			
194	x		192	342	14	249	0	x	x	x	x		D
195			193	342	14	249	0			x			
196			193	342	14	249	0			x			
197	X		193	342	14	249	919				o		
198			193	342	14	249	0			x			
199			193	342	14	249	0			x			
200			192	342	14	249	0			x			
201			193	342	14	248	0			x			
202	x		193	342	14	249	0	x	x	x	x		D
203	x		192	342	14	249	0	x	x	x	x		D
204			193	342	14	248	0			x			















EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

Item	Anchor For R4	Poor Fit	N					DIF					Item
			OC R2	OC R3	3	4	5	OC R2	OC R3	YC R2	YC R3	YC R4	
156			93	155	5	56	0	x	x				
157		x	78	144	4	45	0	x	x	x	x		D
158		x	79	144	4	45	0	x	x	x	x		D
159		x	78	144	4	45	0	x	x	x	x		D
160			78	144	4	45	0		x	x			
161				79	144	4	45	0	x	x	x	x	D
162					78	144	4	45	0			x	
163	X			76	144	4	45	338		x			
164					78	144	4	45	0				
165						78	144	4	45	0		x	
166						78	144	4	45	0	x	x	x
167		x				78	144	4	45	0	x	x	x
168						78	144	4	45	0			
169		x				64	130	3	34	0	x		x
170			x			64	130	3	34	0	x	x	x
171			x			64	130	3	34	0	x	x	x
172		x				64	130	3	34	0	x	x	x
173			x			62	130	3	34	0	x	x	x
174			x			64	130	3	34	0	x	x	x
175		x				64	130	3	34	0	x	x	x
176			x			64	130	3	34	0	x	x	x
177			x			64	130	3	34	0	x	x	x
178						64	130	3	34	0	x		
179	X	x				63	130	3	34	338	x		x
180		x				64	130	3	34	0	x	x	x
181		x				53	127	2	33	0	x	x	x
182			x			53	127	2	33	0	x	x	x
183		x				52	127	2	33	0	x	x	x
184		x				53	127	2	33	0	x	x	x
185		x				52	127	2	33	0	x	x	x
186		x				53	127	2	33	0	x	x	x
187		x				53	127	2	33	0	x	x	x
188			x			53	127	2	33	0	x	x	x
189		x				53	127	2	33	0	o		o
190		x				53	127	2	33	0	x	x	x
191		x				53	127	2	33	0	x	x	x
192			x			52	127	2	33	0	x	x	x
193		x				50	125	2	30	0		x	
194		x				50	125	2	30	0	x	x	x
195		x				50	125	2	30	0	x	x	x
196		x				50	125	2	30	0	x	x	x
197	X					50	125	2	30	337		x	
198		x				50	125	2	30	0	x	x	x
199		x				50	125	2	30	0	x	x	x
200		x				50	125	2	30	0	x	x	x
201		x				50	125	2	30	0	x	x	x
202		x				50	125	2	30	0	x	x	x
203		x				50	125	2	30	0	x	x	x
204		x				50	125	2	30	0	x		

















## Appendix E. Item parameter for all the equated scales estimated

**Table 1.** Item parameters for Amharic

Item	Item difficulty	Item discrimination	Item guessing
1	-2.09	0.79	0.66
2	-2.52	1.28	0.34
3	-2.69	1.18	0.46
4	-2.15	1.36	0.21
5	-1.23	0.75	0.18
6	-2.38	1.13	0.26
7	-1.11	0.85	0.14
8	-0.97	2.15	0.25
9	-1.05	1.36	0.33
11	-1.81	1.04	0.12
12	-0.96	0.93	0.41
13	-1.78	1.42	0.41
14	-1.55	0.60	0.38
15	-1.44	1.28	0.26
16	-1.29	0.80	0.52
17	-1.03	1.63	0.23
18	0.28	1.15	0.30
19	-0.59	1.22	0.61
20	-0.12	0.77	0.09
21	-1.68	0.58	0.19
23	-0.41	0.84	0.12
24	-0.80	0.77	0.27
25	-1.17	1.75	0.29
26	-0.87	1.74	0.20
28	-0.90	1.43	0.35
29	-1.31	0.63	0.09
30	-0.73	0.72	0.17
31	-0.61	1.42	0.21
32	-0.92	1.83	0.28
33	-0.67	1.29	0.32
34	-0.11	0.61	0.21
35	-0.43	0.83	0.22
36	-0.42	0.60	0.29
37	0.22	0.88	0.17
38	-1.00	1.83	0.53
39	-0.87	0.62	0.05
40	1.06	1.06	0.17
41	-0.22	0.63	0.18
43	0.46	0.81	0.15
44	0.92	0.52	0.08
45	-0.35	0.79	0.25
46	-0.38	1.74	0.10
47	-1.16	1.27	0.10
48	-1.16	0.69	0.45
49	0.20	1.14	0.36
50	-0.67	1.03	0.31
52	0.26	0.99	0.14
53	1.07	0.89	0.21
54	0.67	0.75	0.33
55	0.33	1.23	0.23
56	0.16	1.26	0.09
57	-0.10	1.13	0.25
58	0.82	0.78	0.11
59	-1.44	0.89	0.24

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

Item	Item difficulty	Item discrimination	Item guessing
60	-0.27	1.66	0.13
62	-1.35	0.90	0.41
63	-1.03	1.14	0.35
64	-0.71	1.49	0.10
66	1.49	0.80	0.28
67	0.33	0.67	0.29
68	-0.56	0.93	0.20
70	0.06	0.93	0.30
71	0.05	0.52	0.10
72	-0.19	0.90	0.34
74	-0.41	1.51	0.12
75	-0.17	1.24	0.25
76	0.23	1.05	0.11
77	-0.42	1.27	0.33
78	-0.12	1.02	0.23
79	-1.19	0.96	0.26
80	0.31	1.36	0.18
81	1.38	0.55	0.14
82	-0.12	1.25	0.15
83	0.55	0.66	0.15
84	0.09	0.56	0.09
85	-0.88	0.76	0.22
86	2.07	0.70	0.06
87	0.27	0.94	0.20
88	1.10	0.76	0.04
89	1.59	0.64	0.18
90	0.08	1.22	0.17
91	-0.14	0.75	0.12
92	1.05	0.94	0.13
93	-0.52	1.11	0.28
94	0.55	1.09	0.38
95	0.38	1.23	0.32
96	0.87	0.90	0.37
97	1.20	0.68	0.10
98	0.96	0.96	0.15
99	0.39	0.78	0.25
100	-0.78	1.03	0.50
101	0.10	0.99	0.05
103	-0.24	1.02	0.22
104	-0.21	1.50	0.45
105	-0.56	0.65	0.36
106	1.18	0.61	0.11
107	0.78	0.50	0.39
108	0.43	1.22	0.43
109	-0.96	0.80	0.48
110	0.74	1.05	0.05
111	0.62	0.99	0.20
113	2.71	0.36	0.18
114	-0.34	0.63	0.31
115	-0.82	0.80	0.20
116	-0.28	1.09	0.31
117	0.90	0.64	0.05
118	-0.57	1.26	0.23
120	0.73	0.88	0.14
121	-0.83	0.63	0.37
122	0.42	1.23	0.10
123	0.73	1.05	0.08
124	0.89	1.22	0.06
125	1.16	0.72	0.10

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

Item	Item difficulty	Item discrimination	Item guessing
126	1.42	0.71	0.22
128	-0.07	0.96	0.21
129	0.55	0.92	0.07
130	0.06	1.67	0.12
131	1.17	0.82	0.04
132	0.33	1.10	0.74
133	1.09	0.79	0.21
134	0.02	1.19	0.38
135	0.69	0.98	0.09
136	-0.23	1.44	0.22
137	0.53	0.98	0.18
138	1.10	1.34	0.22
140	1.19	0.84	0.06
141	0.18	1.42	0.19
142	0.92	1.17	0.10
146	1.20	0.90	0.14
147	-1.00	0.84	0.40
150	1.95	0.78	0.10
151	0.10	0.85	0.19
152	2.03	0.97	0.27
154	1.33	0.87	0.11
157	-0.44	0.99	0.42
158	2.62	0.83	0.09
159	2.08	0.86	0.07
162	0.30	0.75	0.15
163	0.90	1.24	0.15
164	0.89	0.99	0.43
165	1.19	1.03	0.54
168	0.71	0.76	0.38
170	-0.54	0.75	0.21
171	1.22	0.82	0.44
172	2.70	1.03	0.25
173	2.56	1.08	0.16
174	0.10	0.94	0.49
176	0.90	0.58	0.12
177	1.59	0.80	0.19
178	0.98	1.10	0.73
179	0.14	1.39	0.13
180	1.89	0.89	0.08
181	0.33	0.79	0.17
182	0.96	0.79	0.19
183	0.52	0.68	0.70
184	-0.11	0.78	0.47
186	-0.16	0.74	0.43
188	1.73	0.66	0.15
189	1.31	0.70	0.28
190	-0.67	0.56	0.49
191	0.20	1.15	0.10
192	1.03	1.20	0.41
193	2.05	0.79	0.18
195	0.12	1.28	0.17
196	0.10	0.40	0.24
197	0.72	1.37	0.11
198	-0.59	0.77	0.41
199	0.79	0.64	0.53
200	0.94	1.26	0.51
201	2.34	1.76	0.37
204	0.76	1.12	0.04
300	0.89	0.89	0.20

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

<b>Item</b>	<b>Item difficulty</b>	<b>Item discrimination</b>	<b>Item guessing</b>
301	0.66	0.74	0.25
303	-0.36	0.76	0.33
304	-0.17	0.71	0.34
305	0.51	0.90	0.14
306	-1.21	1.05	0.24
307	-0.45	1.45	0.17
308	0.09	0.67	0.09
309	-1.01	0.74	0.90
310	0.03	1.01	0.21
311	0.82	0.76	0.14
312	0.22	1.35	0.33
313	0.04	0.72	0.22
314	1.27	1.03	0.24
315	-0.48	1.76	0.22
316	-0.58	0.84	0.04
317	0.25	2.63	0.34
318	0.38	2.18	0.23
319	-0.28	0.79	0.23

**Table 2.** Item parameters for Oromifa

Item	Item difficulty	Item discrimination	Item guessing
1	-2.27	1.02	0.46
2	-2.57	1.21	0.38
3	-2.31	1.28	0.53
4	-2.18	1.04	0.27
5	0.46	0.57	0.18
6	-2.02	1.37	0.27
7	-0.82	1.02	0.20
8	-0.70	1.51	0.20
9	-0.94	1.08	0.23
11	-1.78	1.03	0.18
12	-0.92	0.74	0.31
13	-1.43	1.14	0.34
14	-1.66	0.87	0.37
15	-1.14	0.97	0.30
16	-1.49	0.73	0.31
17	-0.72	1.23	0.19
19	-0.72	1.14	0.49
20	-0.75	1.19	0.14
21	-1.44	0.73	0.22
22	-0.91	1.15	0.28
24	-0.25	0.60	0.21
25	-0.88	1.62	0.33
26	-0.60	1.22	0.23
28	-0.29	0.99	0.34
29	-0.80	1.02	0.23
30	-0.57	0.79	0.32
31	-0.31	0.89	0.18
32	-0.77	1.34	0.30
33	-0.39	1.18	0.35
35	0.51	0.91	0.26
36	-0.09	0.76	0.32
37	1.73	0.58	0.17
38	-1.52	1.11	0.44
39	0.06	0.74	0.17
43	0.26	0.97	0.17
46	0.25	1.37	0.16
47	-1.04	1.13	0.22
48	-1.60	0.81	0.26
49	1.12	0.79	0.21
50	-0.56	1.04	0.47
51	0.83	0.74	0.25
52	0.39	1.06	0.12
53	-0.23	1.03	0.14
54	1.09	0.71	0.33
55	0.39	0.94	0.23
56	0.96	0.99	0.19
57	-0.06	1.20	0.32
58	1.58	0.65	0.25
59	-0.65	1.14	0.32
60	0.52	1.81	0.21
62	-1.06	1.41	0.24
63	-1.36	0.89	0.22
64	-0.52	1.20	0.20
67	1.13	0.66	0.34
68	-0.19	0.81	0.26
69	0.61	0.70	0.55
70	1.13	0.70	0.24

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

Item	Item difficulty	Item discrimination	Item guessing
71	1.07	0.85	0.24
72	0.38	0.62	0.42
73	-1.55	0.78	0.52
74	-0.09	1.49	0.14
75	0.03	1.06	0.18
76	1.52	1.25	0.18
77	-0.47	0.91	0.39
78	0.61	0.58	0.19
79	-0.90	0.58	0.18
80	0.85	0.79	0.09
81	1.25	1.00	0.13
82	0.92	1.31	0.21
83	1.03	0.93	0.30
84	0.63	0.80	0.32
86	0.72	0.96	0.15
87	1.11	1.15	0.16
88	2.50	0.74	0.19
89	1.25	1.34	0.21
90	-0.07	1.19	0.25
91	0.78	0.93	0.08
92	1.71	0.75	0.16
93	-0.15	0.72	0.41
94	0.34	0.55	0.30
95	0.20	0.97	0.29
97	1.42	0.97	0.06
98	0.82	0.76	0.15
99	2.04	0.70	0.24
101	-0.48	1.04	0.34
103	-0.18	0.99	0.24
106	1.13	1.25	0.21
107	0.08	0.64	0.30
108	-0.20	0.67	0.61
109	-0.53	1.11	0.14
110	1.31	1.03	0.03
111	0.46	1.16	0.17
112	-0.07	0.65	0.20
114	0.44	0.72	0.51
115	0.28	0.73	0.28
117	1.73	0.93	0.12
120	1.08	0.93	0.11
122	0.19	0.65	0.25
123	0.88	0.66	0.14
125	1.37	0.91	0.16
126	0.17	0.97	0.41
129	1.62	0.72	0.32
130	0.12	1.14	0.22
131	1.82	0.94	0.21
134	-0.07	0.93	0.35
135	0.58	1.01	0.26
136	-0.23	1.35	0.30
139	0.05	0.84	0.22
142	1.06	1.03	0.16
145	0.24	1.45	0.18
148	0.01	0.69	0.40
149	0.55	0.78	0.28
151	0.32	0.81	0.16
152	0.28	1.02	0.40
154	1.75	0.62	0.10
156	0.76	0.82	0.12

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

<b>Item</b>	<b>Item difficulty</b>	<b>Item discrimination</b>	<b>Item guessing</b>
160	1.73	0.74	0.08
162	1.35	0.99	0.13
163	0.82	0.84	0.22
164	1.18	0.84	0.13
165	0.26	0.74	0.18
168	0.97	0.93	0.26
169	-0.18	1.58	0.27
178	-0.28	1.22	0.35
179	0.77	0.75	0.16
189	1.24	0.69	0.33
193	1.21	0.60	0.35
197	0.96	0.82	0.12
204	1.31	0.89	0.24
300	1.00	0.69	0.20
301	0.31	0.65	0.19
302	0.94	0.73	0.27
303	0.87	0.71	0.21
304	1.29	1.00	0.08
305	0.49	0.74	0.34
306	0.76	0.70	0.20
307	1.22	0.62	0.14
308	0.63	1.07	0.13
309	-0.97	1.17	0.34
310	0.46	0.73	0.26
311	0.48	0.85	0.12
312	0.17	0.77	0.30
313	1.04	0.76	0.33
314	0.32	0.85	0.24
315	0.78	0.81	0.31
316	1.48	0.71	0.24

**Table 3.** Item parameters for Tigrinya

Item	Item difficulty	Item discrimination	Item guessing
1	-2.15	1.39	0.71
2	-2.76	0.98	0.38
3	-2.43	1.17	0.42
4	-1.93	1.17	0.23
5	-0.94	0.96	0.43
6	-1.95	1.21	0.55
8	-0.96	1.11	0.15
9	-0.75	1.87	0.34
11	-1.26	0.89	0.51
12	-1.47	0.70	0.44
13	-1.94	1.01	0.22
14	-1.85	0.78	0.51
15	-1.61	1.22	0.20
16	-1.51	0.88	0.27
17	-0.92	1.29	0.10
18	0.49	1.00	0.22
19	-1.07	0.68	0.22
20	-0.71	1.17	0.16
21	-1.25	0.70	0.13
22	-1.04	1.07	0.17
24	-0.79	0.69	0.17
25	-0.88	1.78	0.33
26	-0.67	1.56	0.15
27	-0.30	1.06	0.21
28	-0.87	1.31	0.29
29	-1.18	0.78	0.13
30	-0.86	1.16	0.26
31	-0.49	0.98	0.12
32	-0.89	1.43	0.13
33	-0.79	1.26	0.28
35	-0.68	0.73	0.07
36	-0.98	0.99	0.18
37	0.87	0.90	0.08
38	-1.21	1.13	0.30
39	-0.81	1.97	0.53
40	0.59	0.55	0.20
41	-0.02	1.03	0.15
42	0.46	0.69	0.06
43	0.36	0.96	0.17
44	0.15	1.01	0.20
45	-0.49	0.74	0.07
46	0.02	0.95	0.06
47	-0.91	1.76	0.08
48	-1.01	1.18	0.21
49	0.25	1.04	0.11
50	-0.35	0.98	0.14
51	-0.03	1.50	0.27
52	0.11	0.99	0.12
53	0.64	0.79	0.06
55	0.65	0.86	0.22
56	-0.21	0.80	0.18
57	0.13	0.92	0.14
58	0.34	0.63	0.10
59	-0.41	1.51	0.69
60	-0.11	1.30	0.06
63	-1.12	0.94	0.20
64	-0.75	1.18	0.11

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

Item	Item difficulty	Item discrimination	Item guessing
66	0.45	0.90	0.15
67	0.32	1.13	0.29
68	0.59	0.89	0.49
70	0.63	0.91	0.34
71	0.11	0.52	0.50
72	0.92	0.65	0.40
74	0.19	1.02	0.14
75	-0.55	0.46	0.32
76	0.94	1.05	0.16
78	-0.35	0.69	0.24
79	-0.94	1.05	0.21
80	0.38	1.16	0.11
81	0.92	0.85	0.08
82	0.24	0.67	0.11
83	0.20	1.24	0.19
84	1.22	0.54	0.24
85	-0.14	0.69	0.22
86	1.47	0.90	0.12
87	0.59	1.00	0.29
88	2.01	0.63	0.18
89	1.65	0.53	0.10
90	0.57	1.10	0.24
91	0.11	0.80	0.32
93	0.59	0.94	0.58
94	-0.16	0.74	0.33
95	0.21	1.28	0.10
96	0.75	1.04	0.14
97	1.69	1.16	0.13
98	2.03	0.97	0.25
99	1.11	0.89	0.47
100	-1.38	0.43	0.47
101	1.34	0.62	0.32
102	1.64	0.93	0.15
103	-0.37	0.82	0.17
106	1.15	0.97	0.22
107	0.44	0.70	0.37
108	0.08	0.93	0.45
109	-1.07	0.54	0.51
110	0.76	1.04	0.02
111	0.77	0.78	0.11
112	0.22	0.76	0.22
113	1.13	0.97	0.16
114	-0.74	0.66	0.26
115	0.05	0.68	0.48
116	-0.60	1.08	0.25
117	1.28	1.51	0.22
120	0.54	1.01	0.22
121	0.78	0.70	0.40
122	1.58	1.05	0.16
123	1.11	0.71	0.10
124	1.62	0.91	0.07
125	0.74	0.96	0.09
128	0.66	0.94	0.11
129	1.21	0.47	0.18
130	0.73	1.07	0.25
131	0.65	0.57	0.07
133	0.69	0.70	0.07
134	0.04	0.71	0.64
135	0.78	0.98	0.20

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

<b>Item</b>	<b>Item difficulty</b>	<b>Item discrimination</b>	<b>Item guessing</b>
136	0.67	0.94	0.21
137	0.05	0.67	0.15
138	1.25	0.70	0.08
139	-0.42	0.64	0.79
142	1.54	1.11	0.19
146	2.77	0.61	0.14
147	-0.33	0.63	0.52
148	0.26	0.60	0.57
149	1.05	1.17	0.25
150	1.63	0.68	0.30
151	0.58	0.76	0.35
152	1.10	0.92	0.21
154	1.06	0.54	0.12
155	0.86	1.59	0.54
156	0.64	1.04	0.12
157	0.35	1.14	0.34
158	3.53	0.82	0.36
159	2.37	0.72	0.07
160	0.55	0.43	0.04
161	1.06	0.71	0.51
162	1.17	0.93	0.25
163	1.49	1.15	0.07
164	1.52	0.90	0.08
165	0.81	0.73	0.19
166	0.30	0.89	0.25
167	0.58	0.77	0.46
173	1.85	0.75	0.06
179	1.28	0.55	0.23
180	2.51	0.58	0.11
181	1.90	1.19	0.39
183	1.12	0.83	0.45
185	2.19	1.23	0.07
188	1.85	0.55	0.10
197	1.84	0.74	0.12
204	1.44	0.67	0.14
300	-1.54	0.57	0.21
301	1.12	0.88	0.07
302	1.15	1.32	0.31
303	2.14	0.86	0.44
304	0.39	0.90	0.15
305	-0.34	0.95	0.08
306	0.63	0.69	0.09
307	0.47	0.91	0.09
308	-0.76	1.06	0.14
309	-0.40	0.68	0.95
310	0.87	1.15	0.07
311	0.83	1.34	0.23
312	-0.32	1.08	0.11

**Table 4.** Item parameters for Telugu (India)

Item	Item difficulty	Item discrimination	Item guessing
1	-3.43	0.54	0.59
2	-3.53	0.71	0.23
3	-3.30	1.32	0.24
4	-2.67	1.12	0.37
5	-1.17	0.62	0.13
6	-3.59	0.70	0.23
7	-2.60	0.78	0.17
8	-0.81	0.72	0.13
9	-2.13	0.81	0.15
10	-0.92	0.60	0.29
11	-1.33	1.67	0.51
12	-1.65	1.22	0.25
13	-1.19	0.96	0.27
14	-1.41	1.13	0.46
15	-1.81	0.99	0.20
16	-0.90	1.48	0.45
17	-1.44	1.14	0.20
18	-0.29	1.01	0.23
19	-0.26	1.12	0.46
20	-0.25	1.35	0.26
21	0.16	0.60	0.15
22	-1.03	0.62	0.05
23	-0.41	0.76	0.19
24	-0.15	0.48	0.16
25	-1.39	1.50	0.35
26	-1.42	1.25	0.10
27	0.47	1.29	0.13
28	-1.12	1.25	0.23
29	-1.14	0.45	0.11
30	0.42	0.48	0.45
31	-0.59	1.25	0.18
32	-0.69	0.65	0.29
33	-0.88	1.25	0.31
34	0.00	1.22	0.29
35	-0.05	1.60	0.33
36	-1.36	0.71	0.03
37	0.49	0.68	0.25
38	-2.00	0.89	0.16
39	-0.01	0.88	0.22
40	0.46	1.66	0.18
41	-0.39	1.07	0.25
43	1.04	0.60	0.15
44	1.43	0.54	0.17
45	2.27	0.98	0.25
46	0.42	1.16	0.17
47	-1.51	1.23	0.07
48	-1.82	0.70	0.32
49	1.82	1.01	0.11
50	0.97	0.70	0.19
51	0.51	1.44	0.20
52	-0.38	1.13	0.19
55	-0.05	0.94	0.30
56	0.39	1.17	0.21
57	1.53	0.70	0.32
58	0.53	0.66	0.05
59	-1.09	0.93	0.13
60	-0.36	0.84	0.15

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

<b>Item</b>	<b>Item difficulty</b>	<b>Item discrimination</b>	<b>Item guessing</b>
61	0.04	2.50	0.23
62	0.67	1.19	0.17
63	-0.43	0.54	0.28
64	-0.77	1.71	0.25
65	1.63	1.00	0.28
67	2.01	0.72	0.32
68	-0.43	1.59	0.31
69	-1.06	0.68	0.09
70	-0.35	1.44	0.35
71	-0.86	0.84	0.15
72	0.64	1.55	0.25
73	1.63	0.34	0.18
74	-1.13	1.38	0.16
75	-0.18	0.88	0.16
76	0.89	1.10	0.18
77	-0.25	1.16	0.30
78	-0.78	1.32	0.21
79	-0.04	0.61	0.23
80	0.29	1.04	0.16
81	0.81	0.80	0.21
82	0.09	1.64	0.21
83	0.08	1.15	0.22
84	0.95	0.77	0.23
86	1.81	1.01	0.16
87	0.40	0.36	0.33
88	-0.04	1.13	0.10
90	0.55	1.47	0.28
91	0.27	0.72	0.14
92	1.25	1.04	0.30
93	-1.16	0.99	0.19
96	-0.10	1.01	0.10
97	1.18	0.97	0.19
98	1.42	0.86	0.32
99	-0.11	0.71	0.12
100	0.87	1.77	0.07
101	-0.31	1.08	0.10
103	1.76	0.93	0.11
105	0.17	1.33	0.29
106	0.51	1.03	0.20
108	0.57	0.49	0.38
110	-0.84	0.58	0.23
111	0.68	1.40	0.21
112	0.84	1.51	0.34
113	1.98	1.03	0.26
114	0.87	0.95	0.13
115	0.08	0.87	0.41
116	0.38	1.20	0.26
117	1.22	1.10	0.05
118	0.52	1.20	0.20
119	0.49	1.83	0.29
120	1.12	0.93	0.08
121	0.63	1.02	0.25
122	0.92	0.86	0.24
123	0.48	1.28	0.31
125	0.77	0.73	0.20
126	-0.91	1.37	0.09
129	0.60	0.94	0.24
130	0.69	1.07	0.10
131	-0.56	0.46	0.13

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

Item	Item difficulty	Item discrimination	Item guessing
133	-0.65	1.16	0.18
134	1.76	1.29	0.07
138	1.61	1.04	0.11
140	0.92	0.62	0.18
141	1.28	0.95	0.23
142	0.78	0.80	0.19
143	1.00	1.85	0.31
144	1.03	1.31	0.07
149	0.92	0.69	0.20
150	0.59	0.86	0.22
151	1.55	0.60	0.32
154	0.72	0.74	0.22
155	0.68	1.66	0.24
156	0.64	1.08	0.25
157	0.05	1.54	0.38
160	1.60	0.97	0.18
161	0.28	0.47	0.16
162	2.20	0.38	0.26
163	1.67	1.09	0.17
165	1.19	1.98	0.36
168	1.16	0.99	0.19
169	1.51	1.66	0.18
170	1.81	0.87	0.03
174	1.88	1.70	0.17
175	1.00	1.08	0.10
176	1.36	0.54	0.06
178	0.89	2.03	0.26
179	0.02	1.20	0.16
180	0.88	0.76	0.05
183	0.13	0.55	0.16
185	1.23	0.97	0.24
189	1.54	0.93	0.31
190	0.59	0.86	0.06
191	0.46	1.06	0.05
192	0.20	1.28	0.10
194	1.52	1.04	0.14
196	1.37	0.52	0.21
197	1.21	0.73	0.14
199	1.02	0.76	0.14
200	1.82	1.20	0.08
205	-3.02	0.80	0.53
211	-1.87	1.03	0.28
212	-1.00	0.97	0.10
225	0.25	0.92	0.26
228	-0.76	0.95	0.12
234	-0.08	0.50	0.29
235	-0.41	1.39	0.23
244	1.24	0.92	0.20
246	2.17	0.94	0.21
247	-0.29	1.11	0.12
249	0.44	0.74	0.09
251	-1.64	1.03	0.27
253	-0.66	0.83	0.27
261	0.10	1.03	0.41
266	-1.84	0.95	0.27
267	-0.34	1.77	0.55
270	1.95	0.87	0.20
271	1.40	0.79	0.50
277	-1.94	0.83	0.39

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

<b>Item</b>	<b>Item difficulty</b>	<b>Item discrimination</b>	<b>Item guessing</b>
286	-0.30	0.93	0.17
288	-0.37	0.75	0.22
289	1.72	1.08	0.39
299	-0.06	1.04	0.19
302	1.21	0.71	0.37
304	-0.06	0.89	0.81
305	0.02	0.84	0.20
307	0.92	1.22	0.22
312	1.03	0.55	0.39
315	0.14	1.33	0.21
318	-0.93	1.02	0.19
321	1.91	0.54	0.08
322	-1.22	0.95	0.48
324	1.37	1.08	0.17
325	0.55	0.87	0.22
326	-0.56	0.68	0.11
328	0.72	0.89	0.25
331	0.39	0.58	0.95
334	0.29	1.02	0.20
338	0.48	1.26	0.80
339	0.97	0.98	0.21
340	-1.19	0.46	0.23
341	2.20	0.57	0.72
342	0.94	0.63	0.10
343	1.91	1.51	0.42
345	0.36	0.65	0.87
346	0.51	0.95	0.18
347	0.39	0.98	0.80
349	0.08	1.09	0.40
352	1.20	1.31	0.51
355	1.13	0.49	0.31
356	2.17	0.61	0.28
359	-0.43	1.10	0.13
362	3.14	0.97	0.21
365	0.57	0.63	0.18
370	-0.90	0.92	0.57
372	0.28	0.92	0.25
374	-0.76	0.97	0.22
375	1.51	0.82	0.22
379	1.25	2.01	0.44
383	-0.53	0.91	0.12
384	1.32	0.54	0.08
387	0.33	0.65	0.66
388	-0.41	1.19	0.28
391	0.78	0.87	0.59
395	0.11	0.77	0.21
396	0.37	0.82	0.42
397	3.22	0.79	0.07
401	1.95	0.76	0.34
403	0.22	0.43	0.56
404	1.60	1.03	0.20
405	1.74	2.57	0.37
406	2.16	1.15	0.47
407	1.18	1.24	0.37
408	1.36	1.81	0.24
500	1.13	1.10	0.07
501	0.37	1.92	0.21
502	2.27	0.91	0.03
503	0.74	0.62	0.19

EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

Item	Item difficulty	Item discrimination	Item guessing
504	2.63	0.96	0.38
505	0.29	1.30	0.42
506	-0.03	1.20	0.36
507	1.28	0.91	0.54
508	0.03	0.85	0.19
509	0.98	0.93	0.38
510	1.05	1.66	0.11
511	0.69	1.54	0.24
512	1.57	0.70	0.15
513	0.94	0.57	0.59
514	1.68	0.79	0.08
515	1.22	1.39	0.23
516	1.42	1.00	0.15
517	-0.06	0.46	0.74

**Table 5.** Item parameters for Vietnamese (Vietnam)

Item	Item difficulty	Item discrimination	Item guessing
4	-2.54	1.74	0.36
5	-1.94	0.85	0.04
7	-2.64	1.38	0.13
8	-1.74	1.44	0.33
9	-2.28	1.40	0.14
10	-2.04	0.63	0.12
11	-1.67	1.24	0.52
12	-1.49	1.18	0.46
13	-1.99	1.60	0.63
14	-0.76	0.92	0.65
15	-1.75	2.76	0.66
16	-1.47	1.04	0.25
17	-1.88	1.18	0.05
18	0.00	0.88	0.35
19	-0.97	1.26	0.73
20	-1.51	0.64	0.03
21	-0.87	0.90	0.72
22	-1.12	0.74	0.33
23	-1.81	1.19	0.07
24	-1.89	1.07	0.05
25	-1.37	0.67	0.05
26	-1.63	2.05	0.60
27	-1.42	2.32	0.56
28	-2.37	1.31	0.04
29	-1.38	1.46	0.61
30	-1.75	1.28	0.52
31	-1.84	1.32	0.04
32	-2.42	1.41	0.17
33	-2.14	1.03	0.04
34	-1.41	1.70	0.37
35	-1.59	1.23	0.02
36	-0.75	1.39	0.44
37	-0.46	0.65	0.13
38	-2.29	1.54	0.33
39	-0.67	1.83	0.31
40	-0.30	1.55	0.21
42	0.56	0.79	0.04
43	-0.99	1.64	0.29
44	-1.04	0.86	0.02
45	-0.96	1.14	0.10
46	-0.93	2.25	0.44
47	-1.52	2.03	0.47
48	-0.74	1.25	0.29
49	-1.99	1.11	0.20
50	-0.64	1.24	0.23
51	-0.87	1.69	0.55
52	-0.67	1.22	0.08
55	-0.67	1.37	0.33
56	-0.23	0.61	0.05
57	-1.52	0.80	0.09
58	-0.52	1.65	0.37
59	-0.98	0.77	0.03
60	-0.71	1.93	0.11
61	-1.40	0.78	0.29
63	-0.86	2.86	0.46
64	-1.37	1.23	0.09
65	-1.37	0.43	0.09

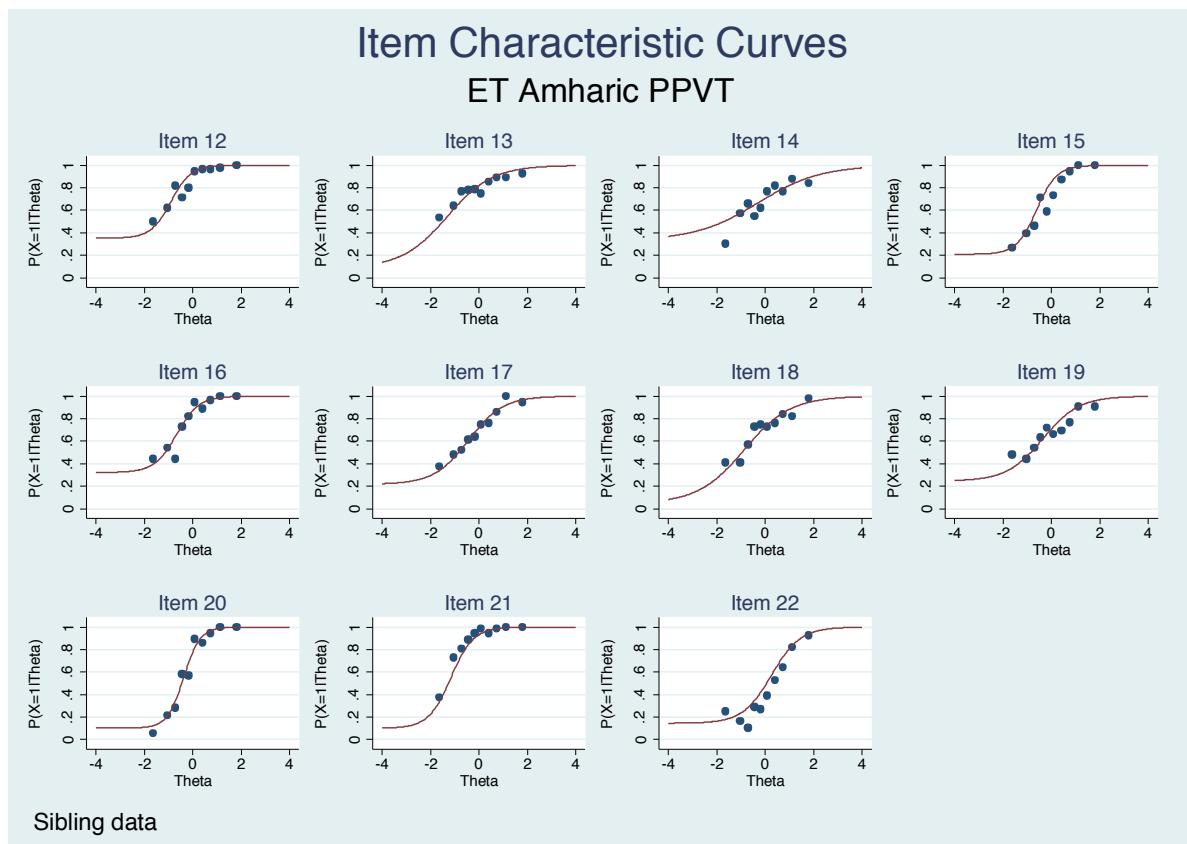
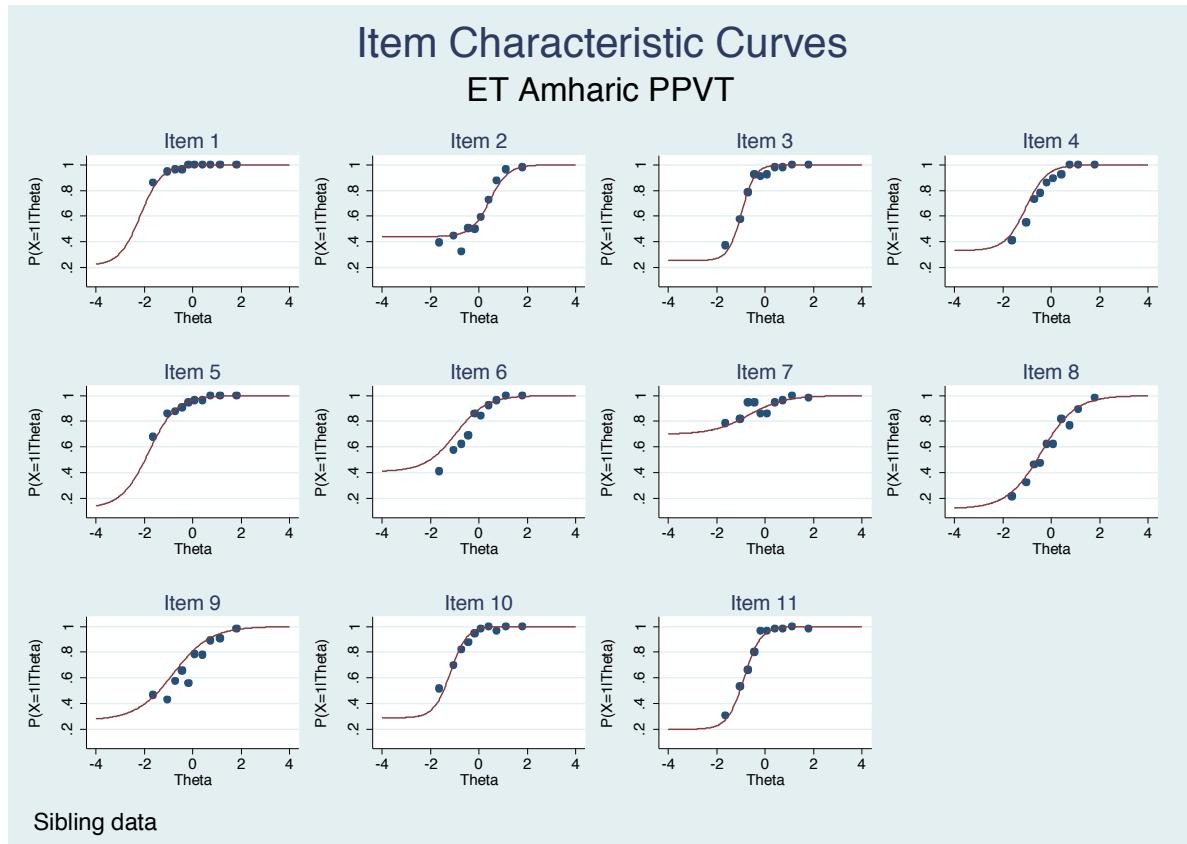
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

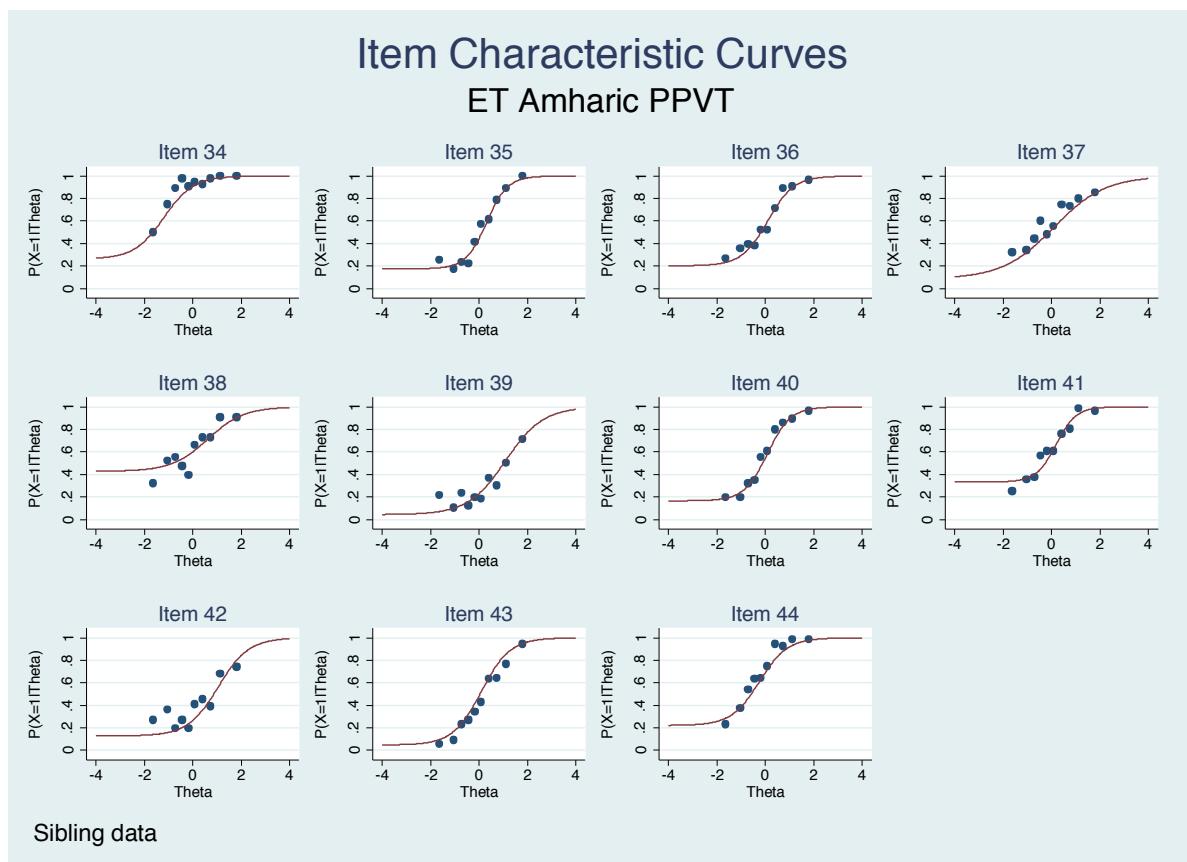
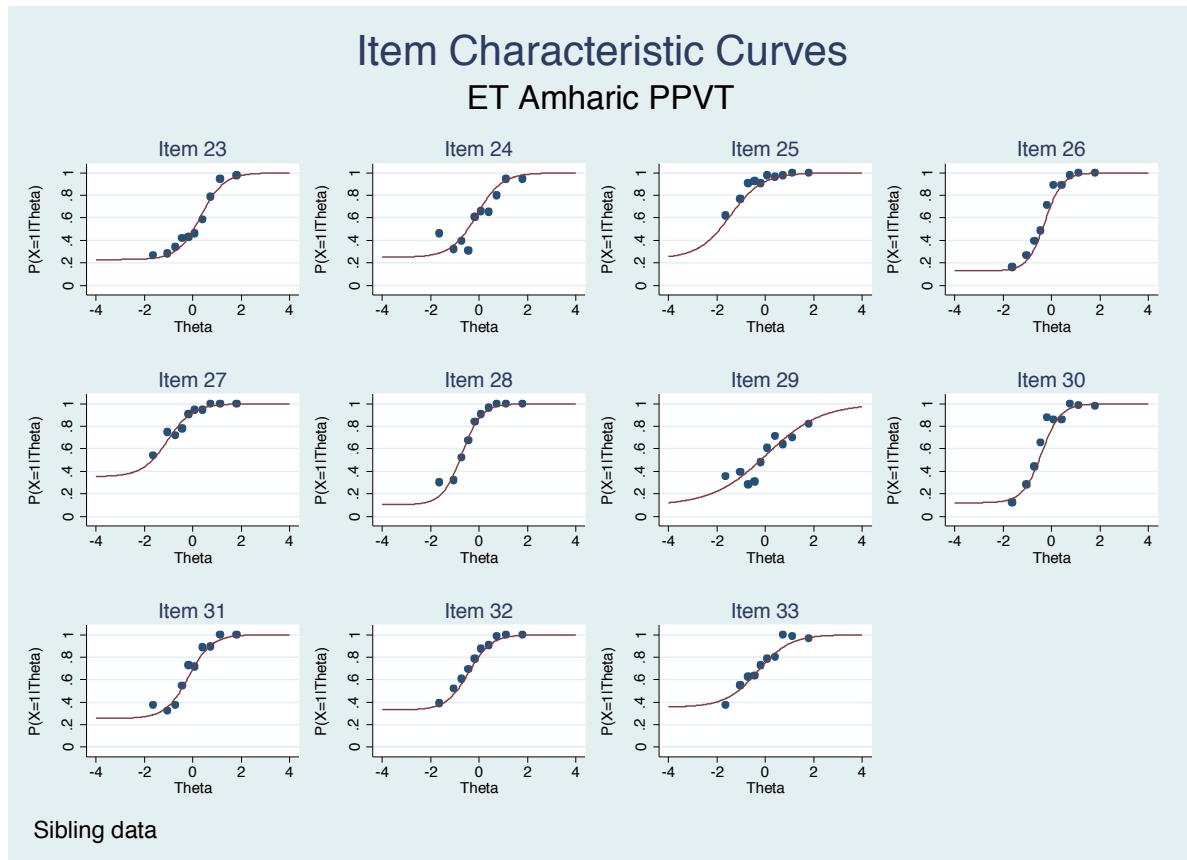
Item	Item difficulty	Item discrimination	Item guessing
66	-0.27	0.83	0.09
67	0.05	0.55	0.31
68	-0.92	1.61	0.48
69	-1.13	1.17	0.25
70	-1.14	1.07	0.11
72	-0.78	1.79	0.59
73	1.74	0.67	0.09
74	-1.19	1.66	0.27
75	-1.33	1.41	0.21
76	-0.40	1.12	0.09
77	-0.88	1.47	0.41
78	-1.22	0.95	0.03
79	-1.00	1.52	0.40
80	-0.50	1.41	0.19
81	-0.58	1.14	0.11
82	-0.78	1.20	0.09
83	-0.43	1.66	0.20
84	-0.53	0.78	0.04
86	2.29	0.70	0.14
87	0.79	1.16	0.16
88	-0.63	0.62	0.05
89	1.07	0.62	0.25
90	-0.93	0.75	0.03
92	0.73	1.19	0.24
93	1.55	0.50	0.49
94	0.92	1.32	0.19
95	0.25	1.79	0.25
96	-0.05	1.73	0.24
97	-0.31	1.11	0.09
98	0.15	1.00	0.13
99	-0.60	0.92	0.30
100	-0.01	1.39	0.13
101	-1.18	0.87	0.14
102	1.62	0.62	0.10
103	-0.95	1.20	0.24
104	-1.27	0.99	0.14
105	-0.49	1.12	0.12
106	-0.18	0.95	0.08
107	-0.34	1.14	0.36
108	-0.67	0.83	0.07
109	0.54	1.25	0.20
110	0.46	1.06	0.60
111	0.13	1.15	0.13
112	-0.36	0.72	0.20
114	-0.52	1.10	0.31
115	0.22	0.80	0.02
116	-0.27	1.24	0.13
117	0.34	1.71	0.17
120	0.12	1.03	0.05
121	-0.34	1.11	0.06
122	0.24	0.80	0.25
123	-0.21	1.16	0.13
125	0.34	0.90	0.07
126	0.63	1.02	0.18
128	0.02	0.97	0.04
130	-0.16	1.18	0.07
131	0.14	0.79	0.05
133	-0.96	0.81	0.28
134	-0.78	1.21	0.15

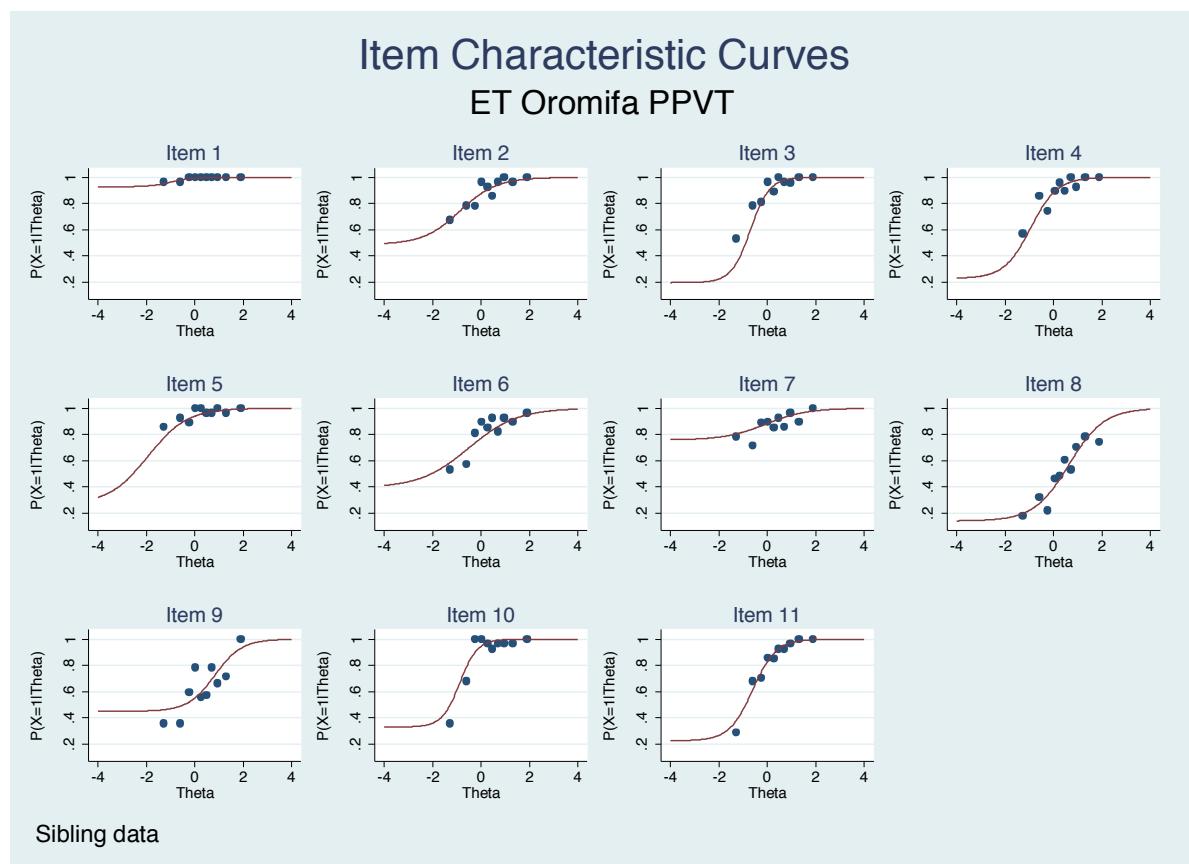
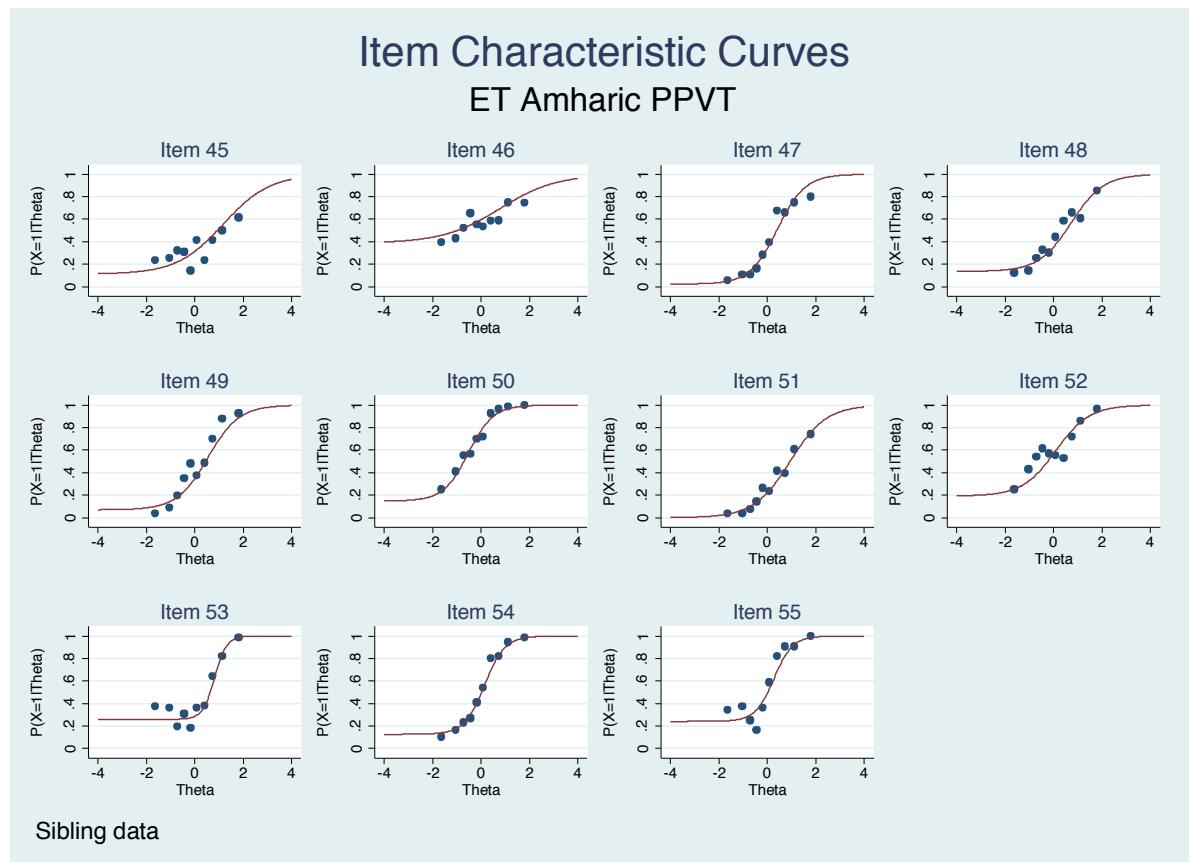
EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM

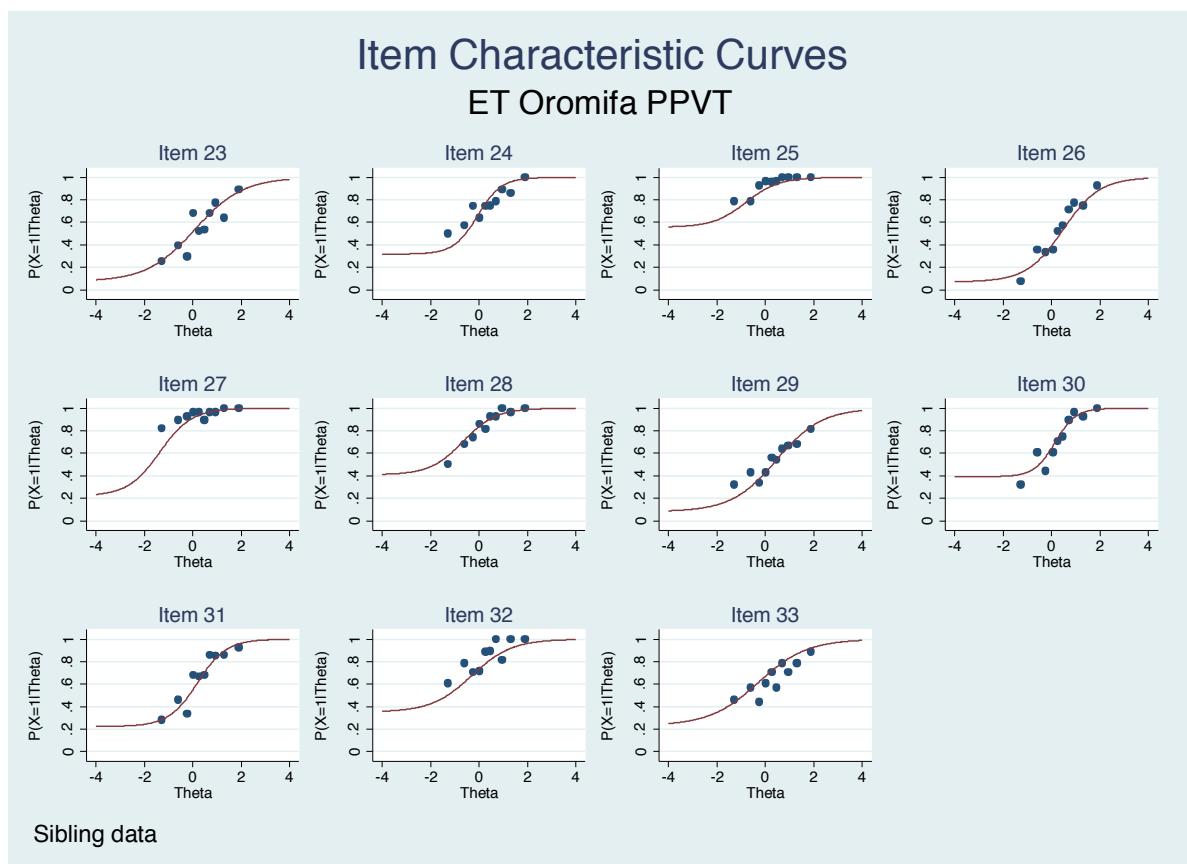
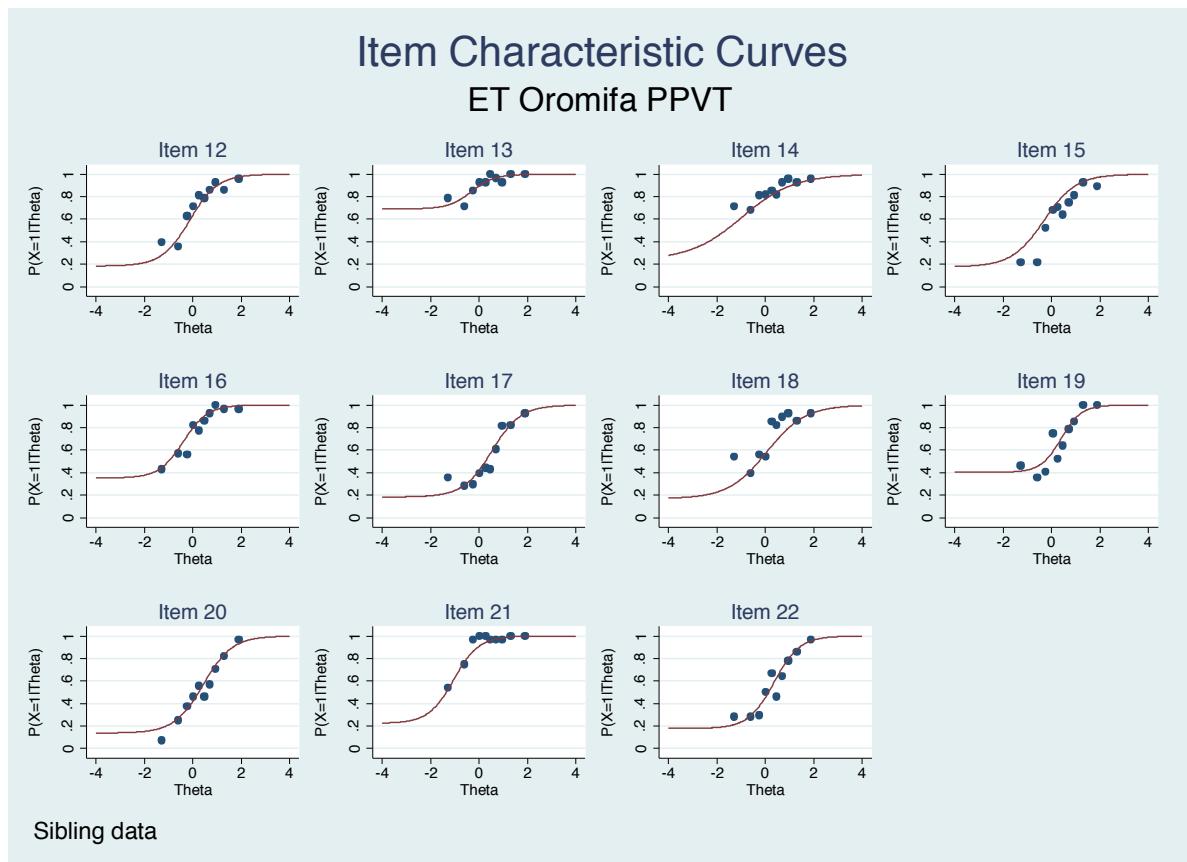
<b>Item</b>	<b>Item difficulty</b>	<b>Item discrimination</b>	<b>Item guessing</b>
135	0.45	1.31	0.31
137	-0.75	0.93	0.08
138	0.95	1.22	0.13
139	0.30	0.72	0.11
140	1.07	0.69	0.20
141	0.17	0.58	0.13
142	0.76	0.68	0.05
144	-0.15	1.13	0.26
146	0.99	1.37	0.26
149	0.38	0.72	0.11
150	1.92	0.50	0.07
152	0.26	0.83	0.24
154	0.30	0.88	0.10
155	0.69	1.14	0.10
156	0.50	2.27	0.38
157	0.57	1.17	0.17
158	1.44	1.50	0.15
159	1.49	1.21	0.20
160	1.92	0.60	0.06
162	1.73	0.59	0.23
163	1.60	1.69	0.18
164	1.20	0.54	0.05
165	1.04	0.86	0.37
167	-0.29	0.73	0.47
168	0.69	0.75	0.06
170	0.07	0.74	0.26
171	-0.02	0.72	0.10
172	1.16	0.83	0.08
173	1.75	1.39	0.19
174	0.57	1.65	0.50
175	-0.01	0.49	0.18
176	0.35	1.11	0.37
178	0.52	1.03	0.15
179	0.88	1.84	0.69
180	1.19	0.99	0.19
185	0.84	0.88	0.13
188	1.15	0.66	0.07
189	0.19	1.27	0.46
190	0.11	1.00	0.62
194	0.25	1.08	0.39
196	0.38	0.93	0.69
197	0.32	0.72	0.14
199	0.44	0.84	0.67
200	-0.10	0.63	0.30
203	-0.63	1.03	0.09
204	0.24	1.03	0.31
300	-0.23	1.12	0.30
301	0.21	1.24	0.07
302	0.09	1.51	0.14
303	0.70	0.60	0.14
304	0.18	0.88	0.23
305	-0.37	0.80	0.13
306	0.27	0.74	0.21
307	2.81	0.94	0.18
308	0.78	0.95	0.27
309	0.20	0.87	0.02
310	1.56	0.59	0.11

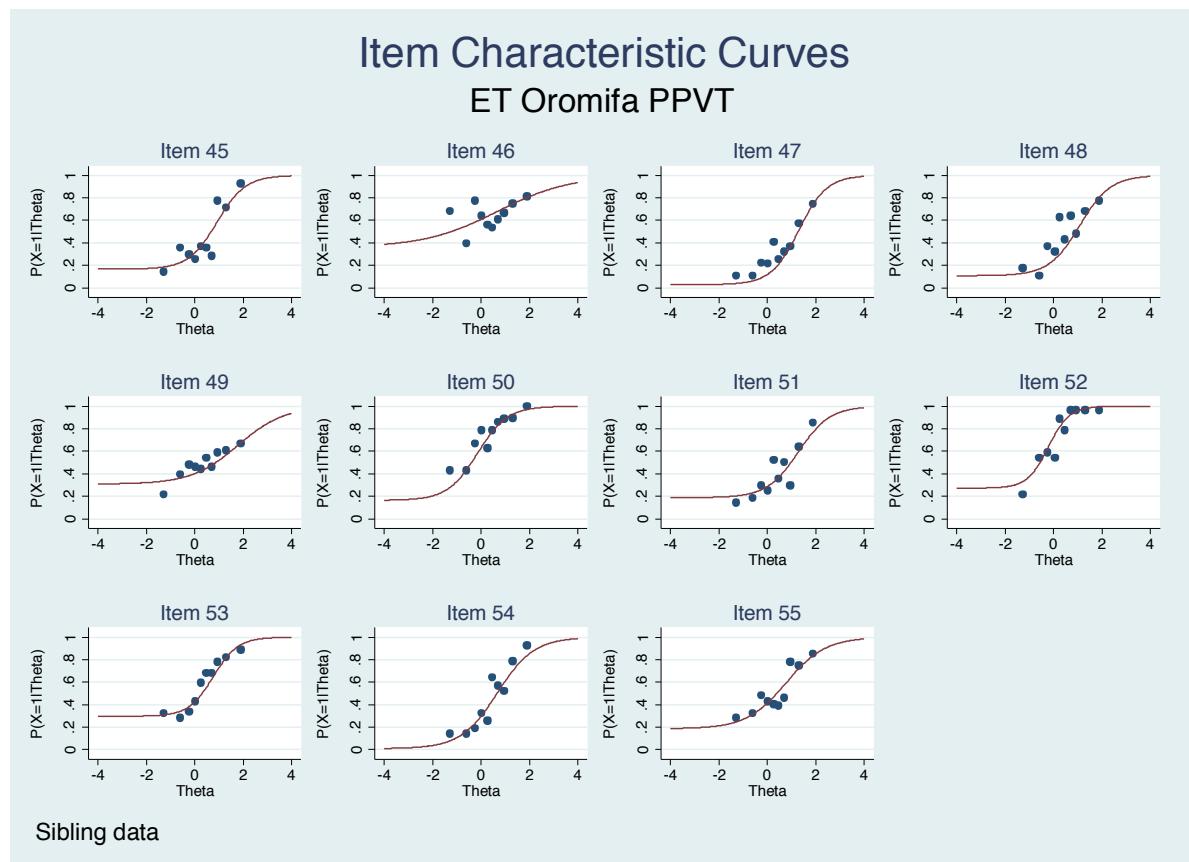
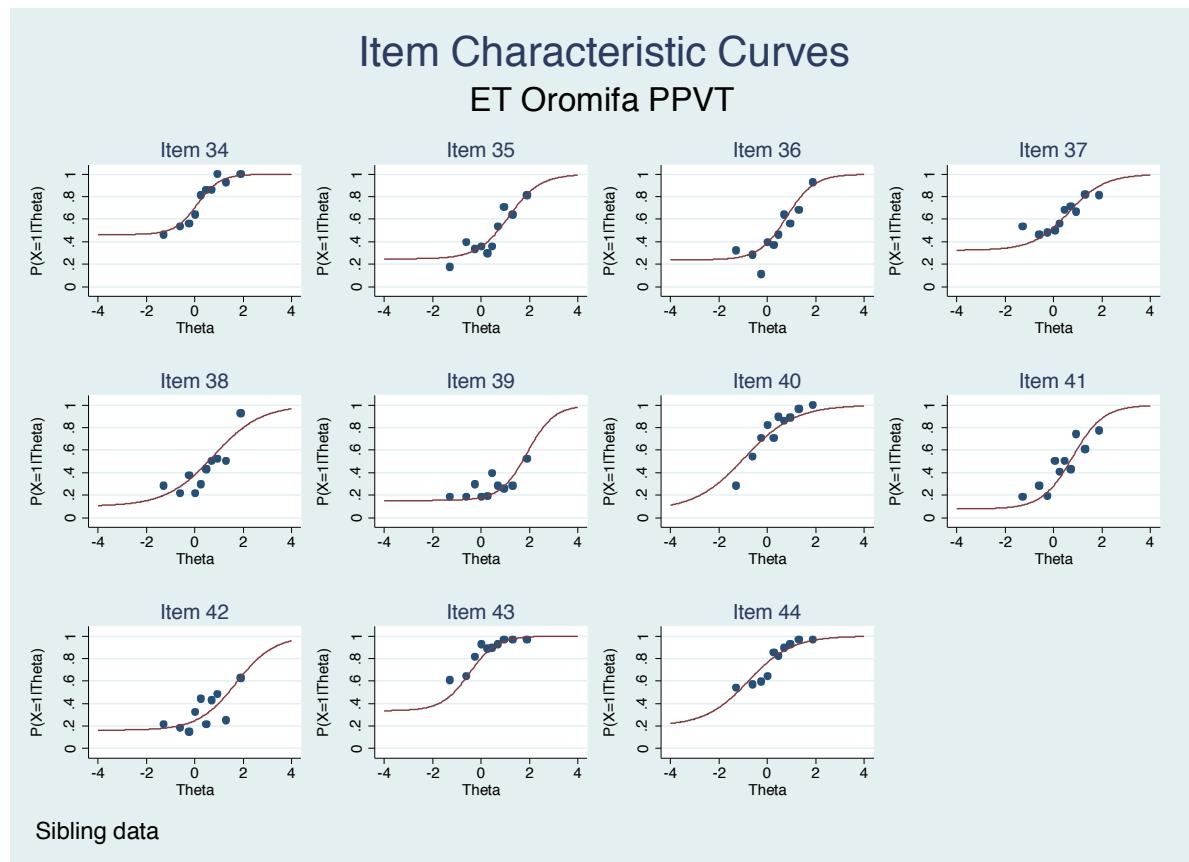
## Appendix F. ICC curves for IRT equating analysis performed with siblings in Round 4

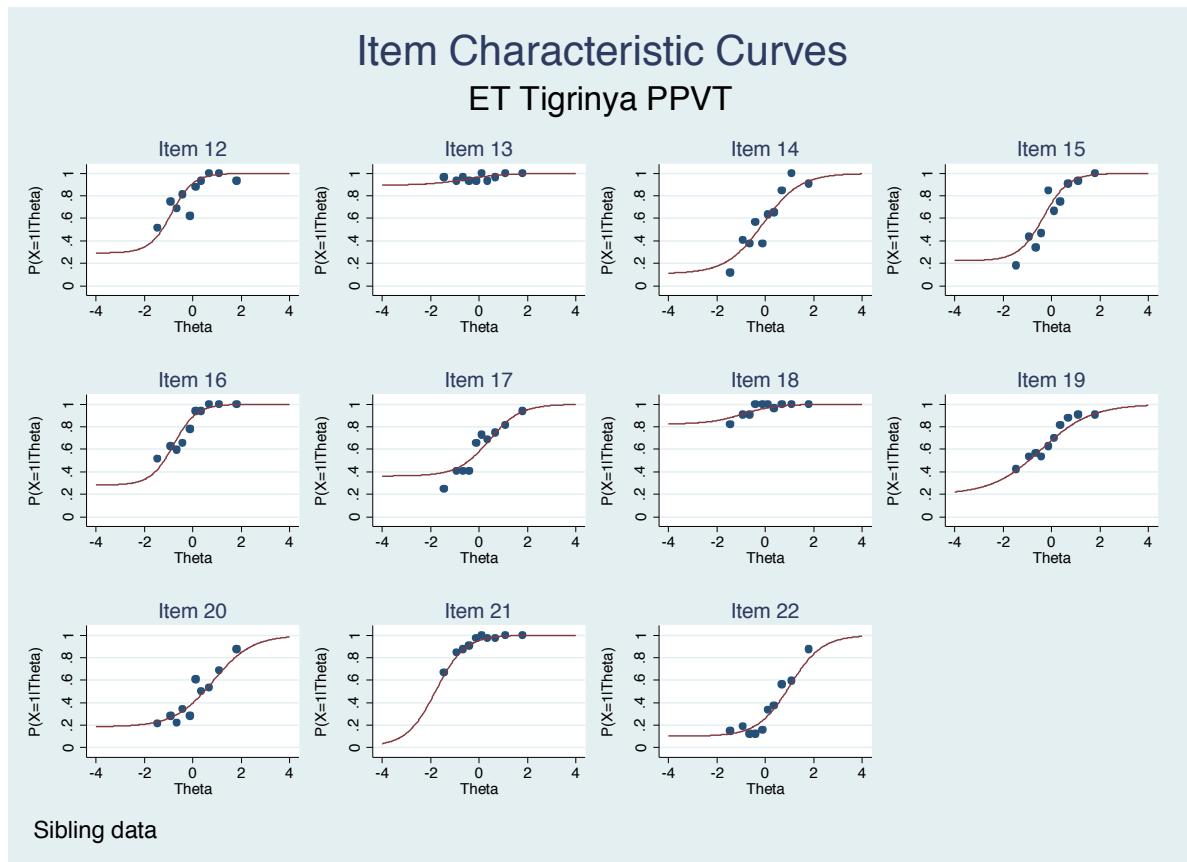
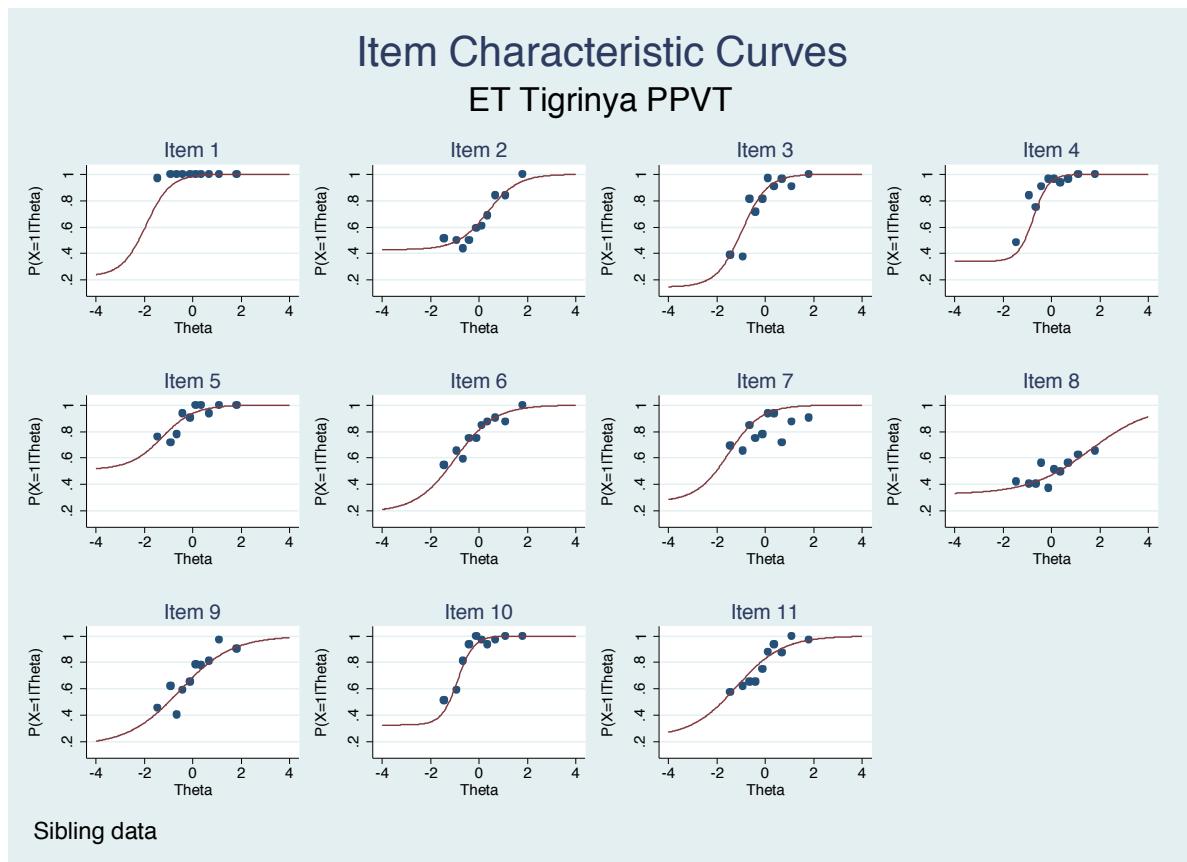


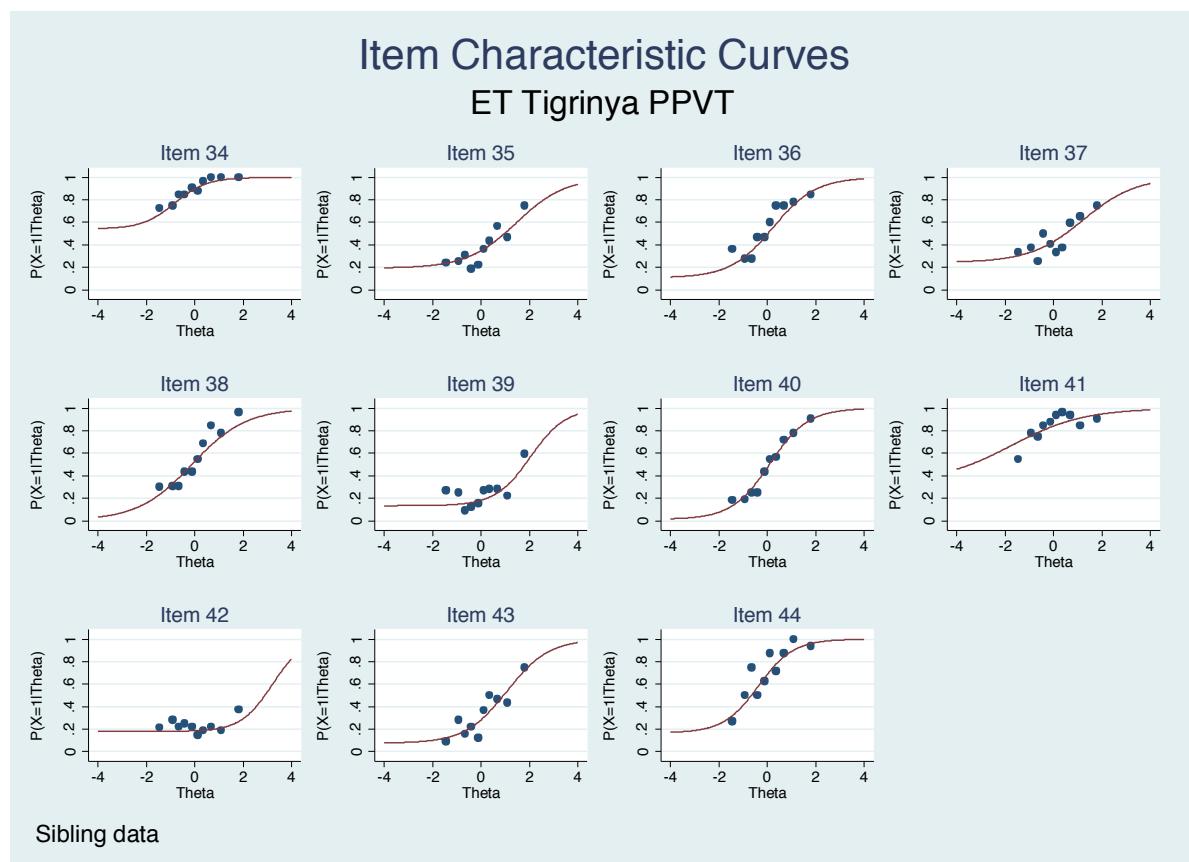
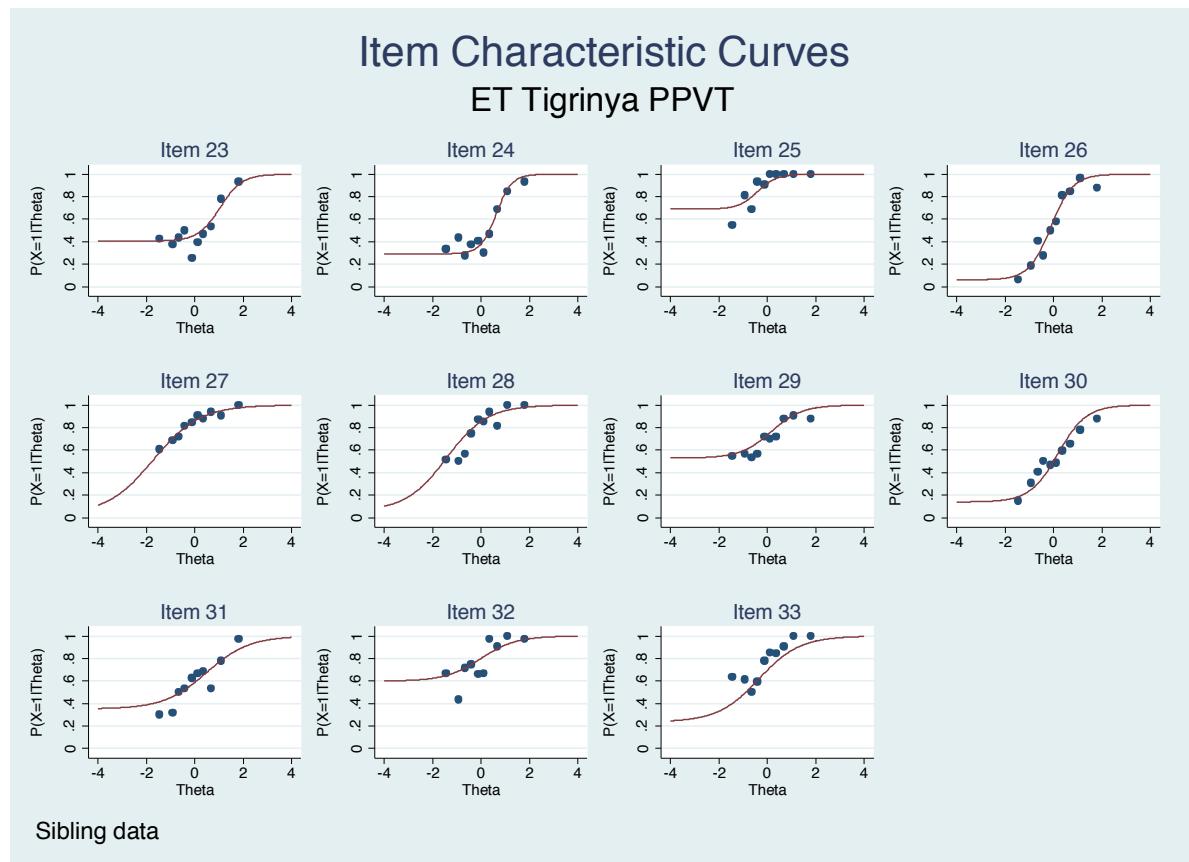


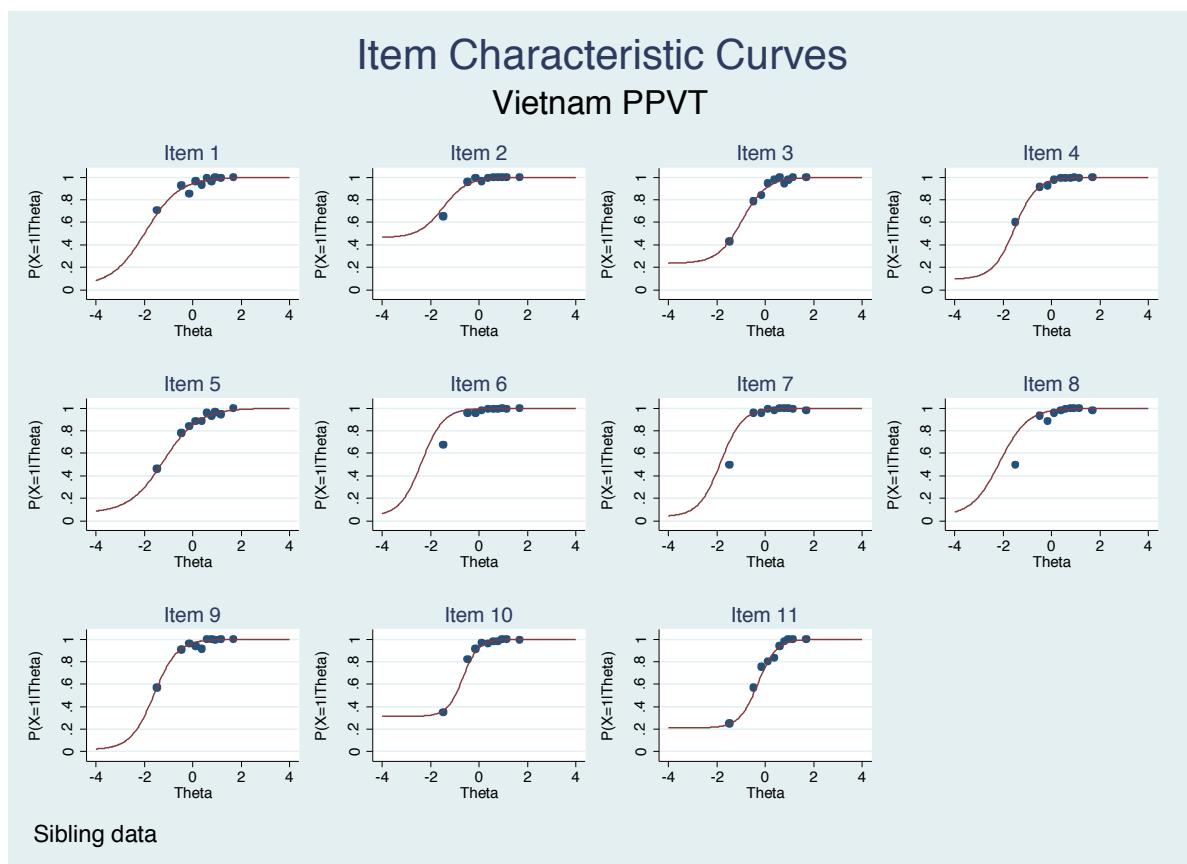
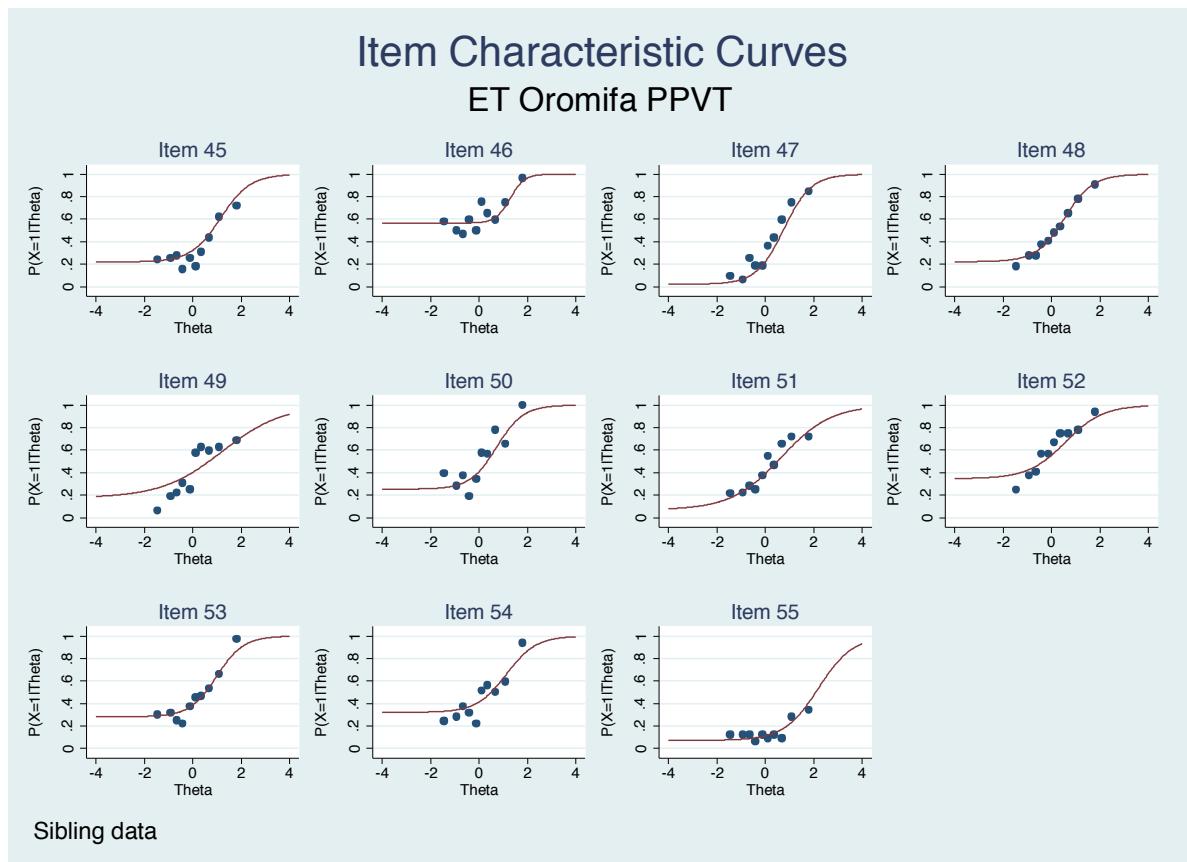


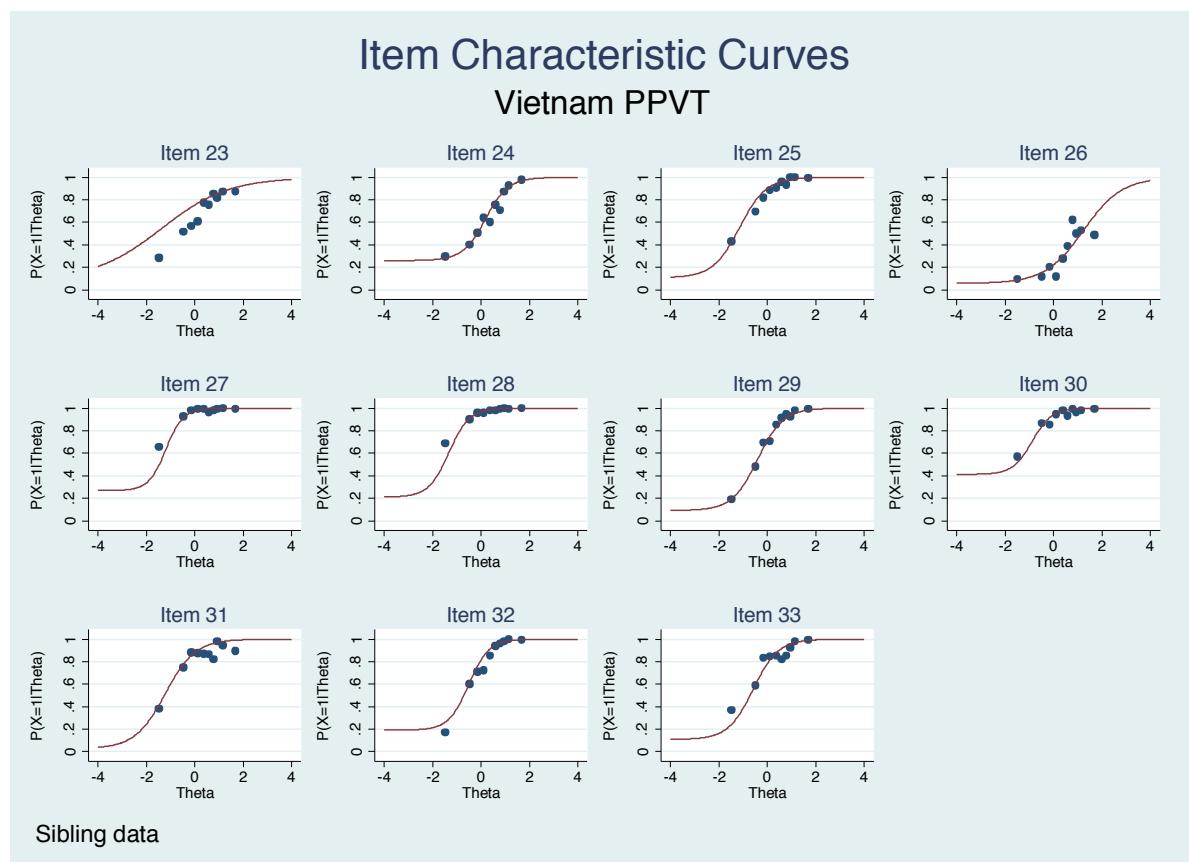
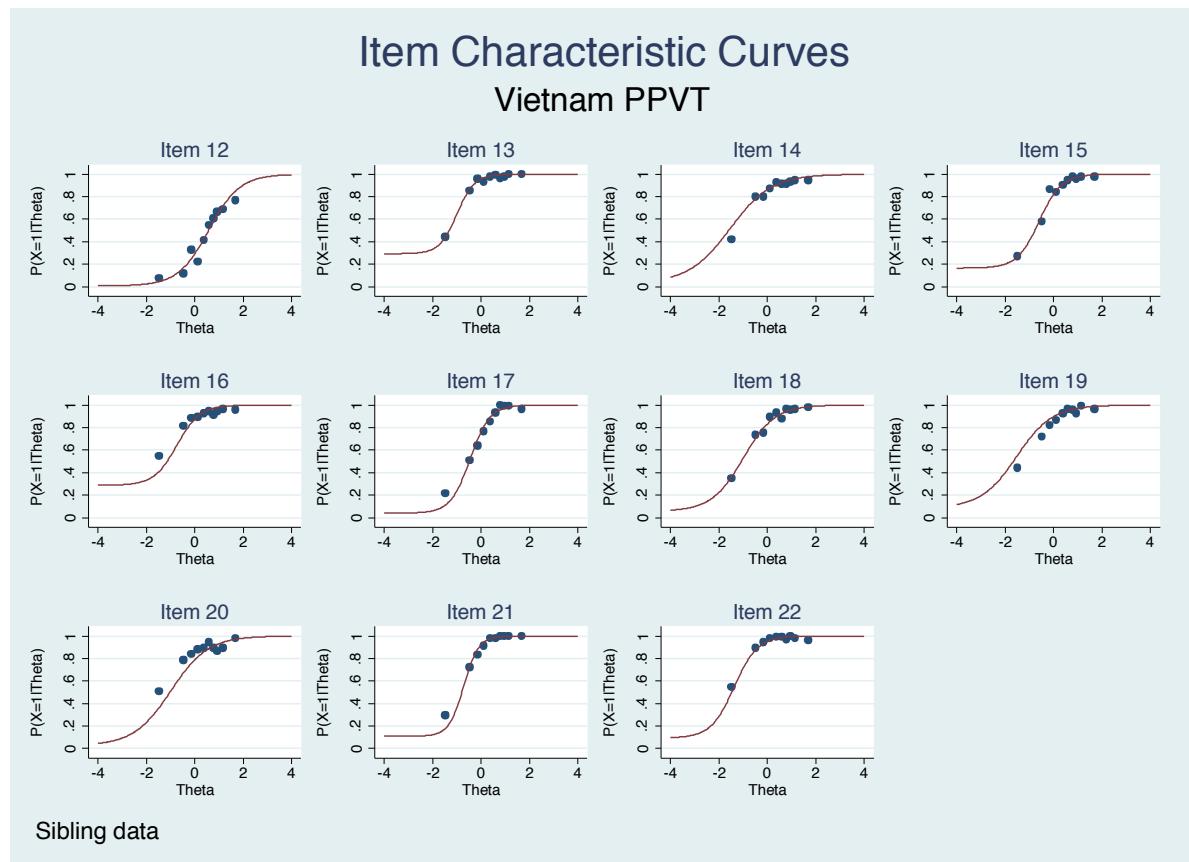


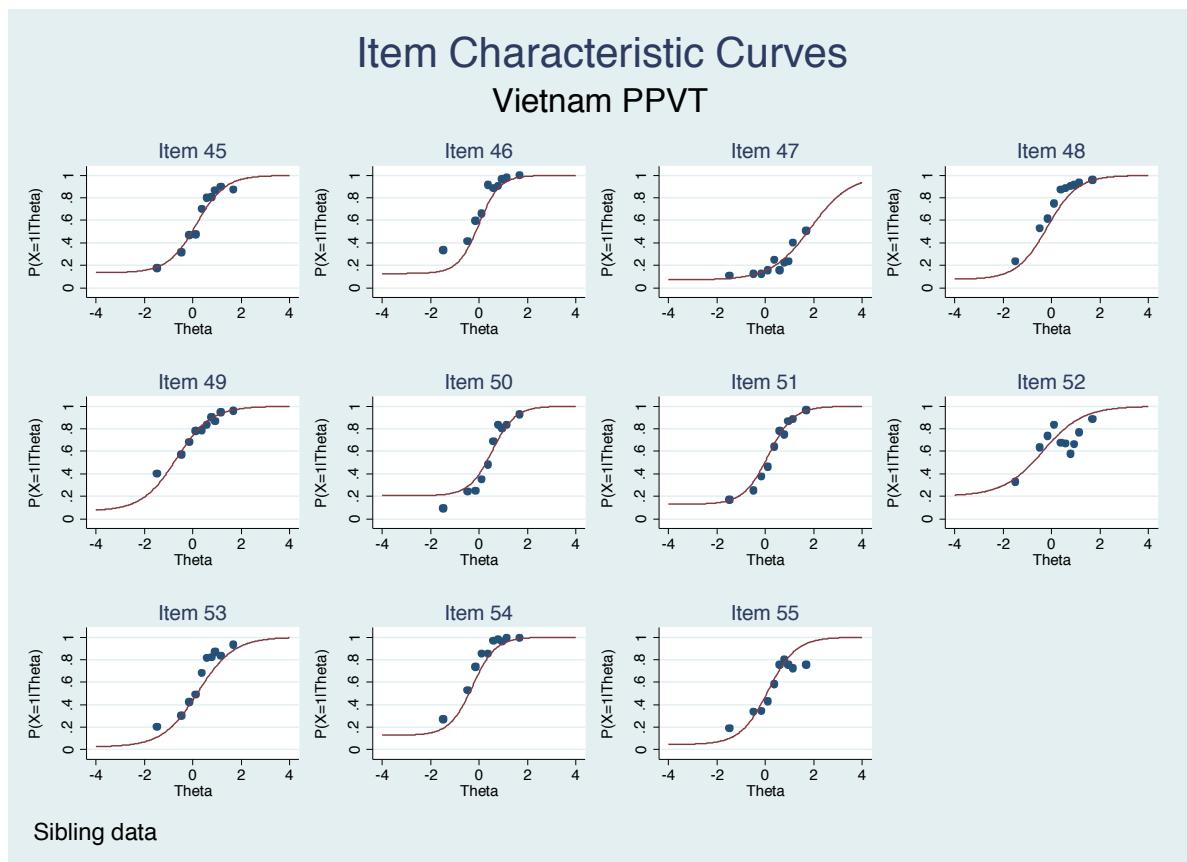
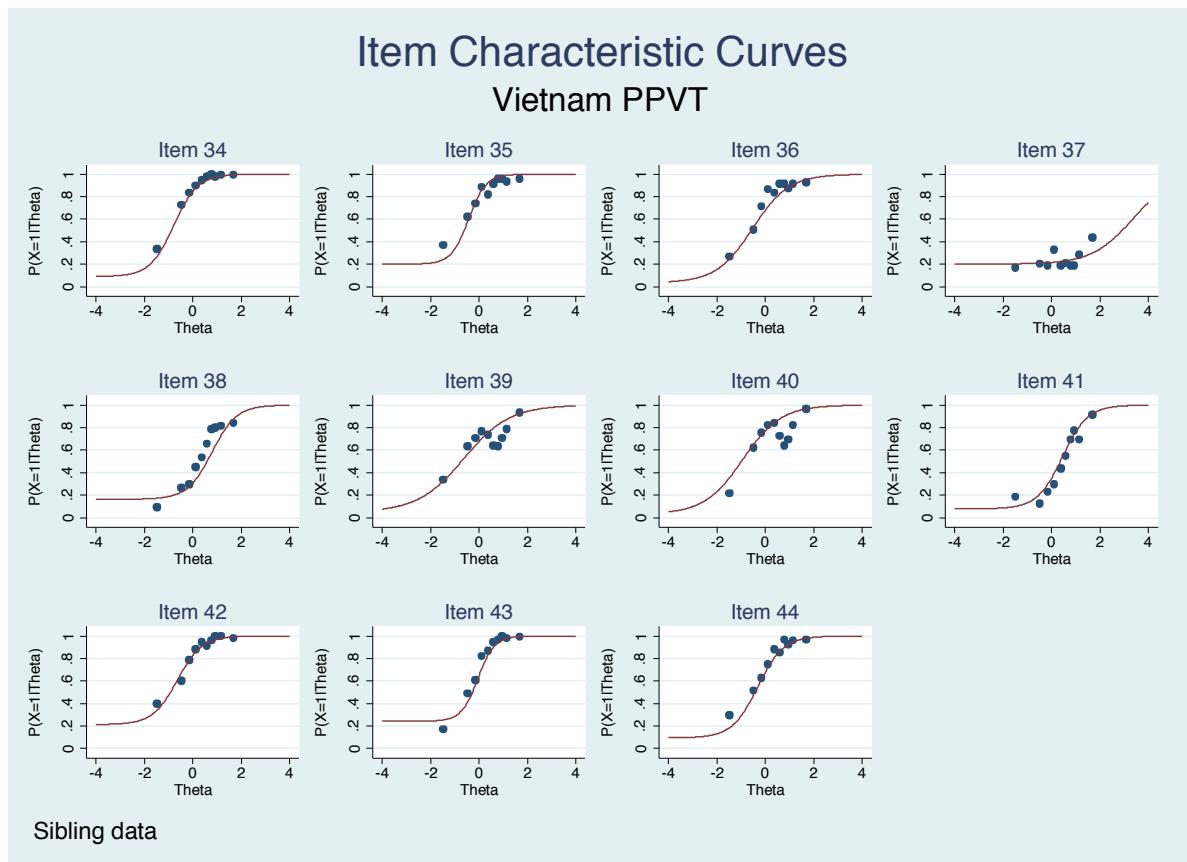




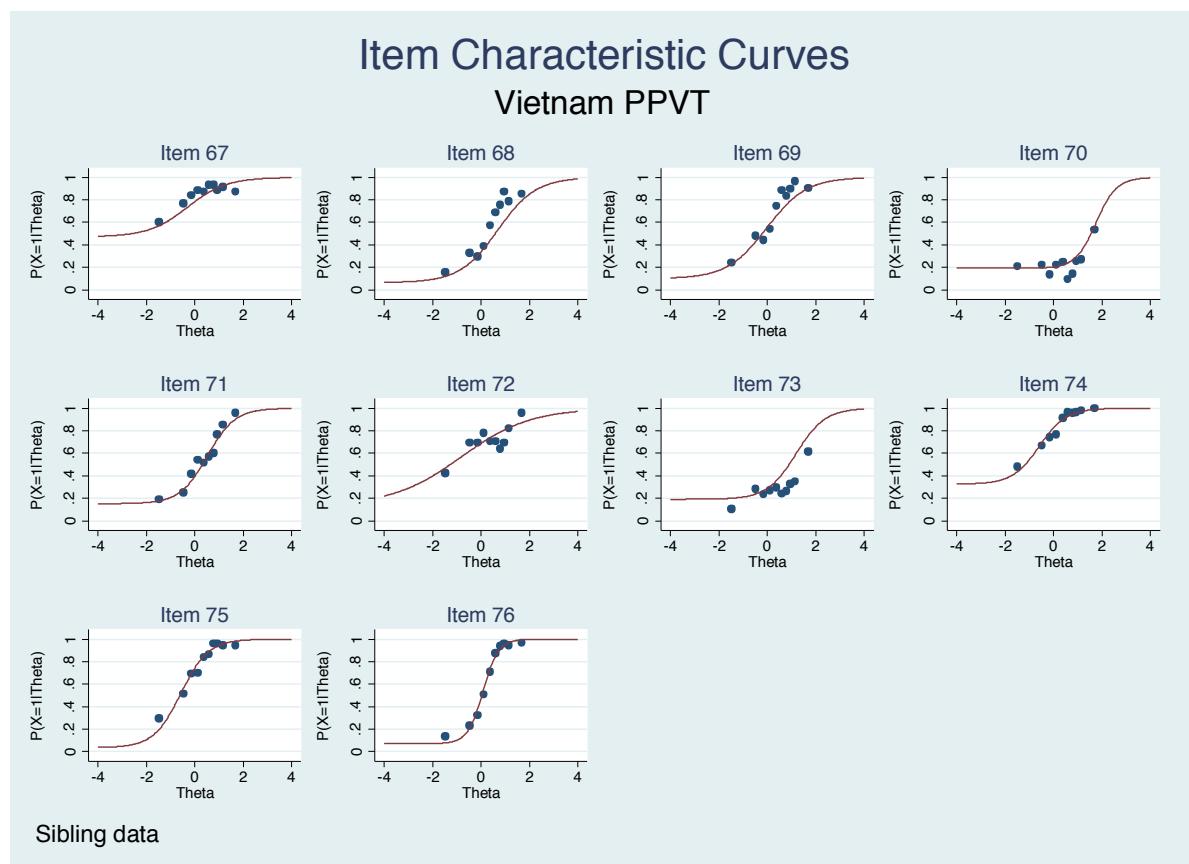
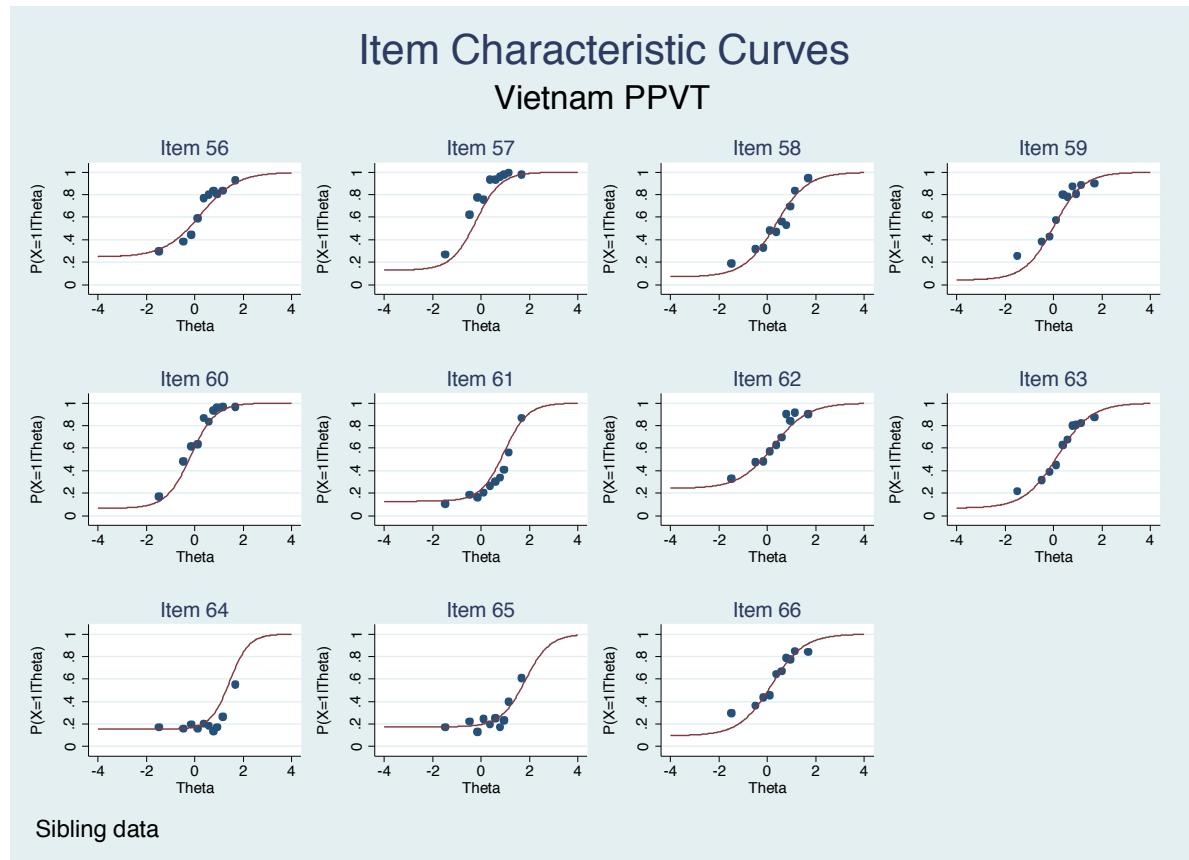








EQUATING TEST SCORES FOR RECEPTIVE VOCABULARY ACROSS  
ROUNDS AND COHORTS IN ETHIOPIA, INDIA AND VIETNAM



## Appendix G. Item parameters used to equate sibling scores with main survey sample

**Table 1.** Item parameters for Amharic

Anchor Item	Item discrimination	Item Difficulty	Item Guessing
4	1.36	-2.15	0.21
8	2.15	-0.97	0.25
9	1.36	-1.05	0.33
11	1.04	-1.81	0.12
12	0.93	-0.96	0.41
23	0.84	-0.41	0.12
24	0.77	-0.80	0.27
25	1.75	-1.17	0.29
26	1.74	-0.87	0.20
28	1.43	-0.90	0.35
29	0.63	-1.31	0.09
31	1.42	-0.61	0.21
33	1.29	-0.67	0.32
35	0.83	-0.43	0.22
39	0.62	-0.87	0.05
45	0.79	-0.35	0.25
46	1.74	-0.38	0.10
47	1.27	-1.16	0.10
52	0.99	0.26	0.14
55	1.23	0.33	0.23
57	1.13	-0.10	0.25
59	0.89	-1.44	0.24
60	1.66	-0.27	0.13
63	1.14	-1.03	0.35
64	1.49	-0.71	0.10
71	0.52	0.05	0.10
74	1.51	-0.41	0.12
75	1.24	-0.17	0.25
77	1.27	-0.42	0.33
79	0.96	-1.19	0.26
80	1.36	0.31	0.18
84	0.56	0.09	0.09
88	0.76	1.10	0.04
90	1.22	0.08	0.17
92	0.94	1.05	0.13
101	0.99	0.10	0.05
103	1.02	-0.24	0.22
106	0.61	1.18	0.11
107	0.50	0.78	0.39
120	0.88	0.73	0.14
129	0.92	0.55	0.07
151	0.85	0.10	0.19
179	1.39	0.14	0.13

**Table 2.** Item parameters for Oromifa

Anchor Item	Item discrimination	Item Difficulty	Item Guessing
8	1.51	-0.70	0.20
9	1.08	-0.94	0.23
25	1.62	-0.88	0.33
26	1.22	-0.60	0.23
31	0.89	-0.31	0.18
33	1.18	-0.39	0.35
39	0.74	0.06	0.17
47	1.13	-1.04	0.22
57	1.20	-0.06	0.32
63	0.89	-1.36	0.22
84	0.80	0.63	0.32
91	0.93	0.78	0.08
92	0.75	1.71	0.16
101	1.04	-0.48	0.34
110	1.03	1.31	0.03
120	0.93	1.08	0.11

**Table 3.** Item parameters for Tigrinya

Anchor Item	Item discrimination	Item Difficulty	Item Guessing
4	1.17	-1.93	0.23
8	1.11	-0.96	0.15
9	1.87	-0.75	0.34
11	0.89	-1.26	0.51
16	0.88	-1.51	0.27
25	1.78	-0.88	0.33
28	1.31	-0.87	0.29
33	1.26	-0.79	0.28
59	1.51	-0.41	0.69
60	1.30	-0.11	0.06
74	1.02	0.19	0.14
78	0.69	-0.35	0.24
82	0.67	0.24	0.11
84	0.54	1.22	0.24
103	0.82	-0.37	0.17
106	0.97	1.15	0.22
110	1.04	0.76	0.02
120	1.01	0.54	0.22
129	0.47	1.21	0.18
130	1.07	0.73	0.25
131	0.57	0.65	0.07
151	0.76	0.58	0.35

**Table 4.** Item parameters for Vietnam

Anchor Item	Item discrimination	Item Difficulty	Item Guessing
5	0.85	-1.94	0.04
12	1.18	-1.49	0.46
28	1.31	-2.37	0.04
31	1.32	-1.84	0.04
33	1.03	-2.14	0.04
35	1.23	-1.59	0.02
39	1.83	-0.67	0.31
40	1.55	-0.30	0.21
43	1.64	-0.99	0.29
48	1.25	-0.74	0.29
57	0.80	-1.52	0.09
59	0.77	-0.98	0.03
60	1.93	-0.71	0.11
64	1.23	-1.37	0.09
65	0.43	-1.37	0.09
70	1.07	-1.14	0.11
74	1.66	-1.19	0.27
75	1.41	-1.33	0.21
76	1.12	-0.40	0.09
77	1.47	-0.88	0.41
78	0.95	-1.22	0.03
80	1.41	-0.50	0.19
81	1.14	-0.58	0.11
82	1.20	-0.78	0.09
83	1.66	-0.43	0.20
84	0.78	-0.53	0.04
87	1.16	0.79	0.16
88	0.62	-0.63	0.05
90	0.75	-0.93	0.03
96	1.73	-0.05	0.24
97	1.11	-0.31	0.09
98	1.00	0.15	0.13
100	1.39	-0.01	0.13
106	0.95	-0.18	0.08
108	0.83	-0.67	0.07
109	1.25	0.54	0.20
111	1.15	0.13	0.13
112	0.72	-0.36	0.20
115	0.80	0.22	0.02
116	1.24	-0.27	0.13
120	1.03	0.12	0.05
122	0.80	0.24	0.25
123	1.16	-0.21	0.13
125	0.90	0.34	0.07
128	0.97	0.02	0.04
130	1.18	-0.16	0.07
138	1.22	0.95	0.13
152	0.83	0.26	0.24
158	1.50	1.44	0.15
167	0.73	-0.29	0.47
168	0.75	0.69	0.06
171	0.72	-0.02	0.10
173	1.39	1.75	0.19
178	1.03	0.52	0.15
180	0.99	1.19	0.19

# Equating Test Scores for Receptive Vocabulary Across Rounds and Cohorts in Ethiopia, India and Vietnam



An International Study of Childhood Poverty

## About Young Lives

Young Lives is an international study of childhood poverty, involving 12,000 children in 4 countries over 15 years. It is led by a team in the Department of International Development at the University of Oxford in association with research and policy partners in the 4 study countries: Ethiopia, India, Peru and Vietnam.

Through researching different aspects of children's lives, we seek to improve policies and programmes for children.

---

## Young Lives Partners

Young Lives is coordinated by a small team based at the University of Oxford, led by Professor Jo Boyden.

- *Ethiopian Development Research Institute, Ethiopia*
  - *Pankhurst Development Research and Consulting plc, Ethiopia*
  - *Centre for Economic and Social Studies, Hyderabad, India*
  - *Sri Padmavathi Mahila Visvavidyalayam (Women's University), Andhra Pradesh, India*
  - *Grupo de Análisis para el Desarrollo (GRADE), Peru*
  - *Instituto de Investigación Nutricional (IIN), Peru*
  - *Centre for Analysis and Forecasting, Vietnamese Academy of Social Sciences, Vietnam*
  - *General Statistics Office, Vietnam*
  - *Oxford Department of International Development, University of Oxford, UK*
- 

## Contact:

**Young Lives**

Oxford Department of  
International Development,  
University of Oxford,  
Mansfield Road,  
Oxford OX1 3TB, UK

Tel: +44 (0)1865 281751

Email: [younglives@younglives.org.uk](mailto:younglives@younglives.org.uk)

Website: [www.younglives.org.uk](http://www.younglives.org.uk)

