

Metadata of the chapter that will be visualized online

Series Title	
Chapter Title	The Measurement of Translation Error in PISA-2006 Items: An Application of the Theory of Test Translation Error
Chapter SubTitle	
Copyright Year	2012
Copyright Holder	Springer Science + Business Media Dordrecht
Corresponding Author	Family Name Solano-Flores Particle Given Name Guillermo Suffix Division School of Education Organization University of Colorado at Boulder Address 249 UCB, 80309, Boulder, CO, USA Email Guillermo.Solano@Colorado.Edu
Author	Family Name Contreras-Niño Particle Given Name Luis Ángel Suffix Division Instituto de Investigación y Desarrollo Educativo Organization University of Baja California Address Km. 103 Carretera Tijuana-Ensenada., c.p. 22830, Ensenada, BC, Mexico Email angel@uabc.edu.mx
Author	Family Name Backhoff Particle Given Name Eduardo Suffix Division Instituto de Investigación y Desarrollo Educativo Organization University of Baja California Address Km. 103 Carretera Tijuana-Ensenada., c.p. 22830, Ensenada, BC, Mexico Email backhoff@uabc.edu.mx

Abstract We examined the translation of PISA test items based on the theory of test translation error (TTTE), which has proven to allow detection of translation errors with unprecedented levels of detail. Translation error (TE) is defined as the lack of equivalence between the original and translated versions of items on multiple translation error dimensions (TEDs) that involve design, language, and content. According to the theory, TE results not only from poor translation, but also from factors that are beyond the translators' skills (e.g., languages encode meaning in different ways). We examined the Mexican, Spanish language translation of science and mathematics PISA 2006 items. A panel comprising teachers, translators, a linguist, a test developer, and a measurement specialist examined the translation of 193 text analytical units (55 pieces of introductory text and 138 items) and identified and coded the TEs identified on ten TEDs. For each item, TE was measured as the number of different TEDs on which the review panel identified TEs. To determine which TEDs are critical to student performance, we examined the correlation between TE and item difficulty (percentage of correct answers and mean proportional score, respectively for dichotomous and non-dichotomous items) considering different sets of TEDs. The highest correlations were observed

for the sets that included the dimensions, Grammar, Semantics, Register, Information, Construct, and Culture. We also observed different magnitudes of correlations for science and mathematics items and a stronger, statistically significant correlation for translated items whose translation the review panel identified more objectionable than for the rest of the items. These results confirm that language- and content-related TEs may threaten the validity of translated items. They speak to the value using the TTTE as a formative evaluation tool that PISA countries can use to operationalize translation guidelines.

Solano-Flores, G., Contreras-Niño, L.A., & Backhoff, E. (2012). The measurement of translation error in PISA-2006 items: An application of the theory of test translation error. In Prenzel, M., Kobarg, M., Schöps, K., & Rönnebeck, S. (Eds.), *Research in the Context of the Programme for International Student Assessment*. Springer Verlag.

Chapter 8 1
The Measurement of Translation Error 2
in PISA-2006 Items: An Application 3
of the Theory of Test Translation Error 4

Guillermo Solano-Flores, Luis Ángel Contreras-Niño, 5
and Eduardo Backhoff 6

[AU1] **Abstract** We examined the translation of PISA test items based on the theory of 7
test translation error (TTTE), which has proven to allow detection of translation 8
errors with unprecedented levels of detail. Translation error (TE) is defined as the 9
lack of equivalence between the original and translated versions of items on multiple 10
translation error dimensions (TEDs) that involve design, language, and content. 11
According to the theory, TE results not only from poor translation, but also from 12
factors that are beyond the translators' skills (e.g., languages encode meaning in 13
different ways). We examined the Mexican, Spanish language translation of science 14
and mathematics PISA 2006 items. A panel comprising teachers, translators, a linguist, 15
a test developer, and a measurement specialist examined the translation of 193 text 16
analytical units (55 pieces of introductory text and 138 items) and identified and 17
coded the TEs identified on ten TEDs. For each item, TE was measured as the number 18
of different TEDs on which the review panel identified TEs. To determine which 19
TEDs are critical to student performance, we examined the correlation between TE 20
and item difficulty (percentage of correct answers and mean proportional score, 21
respectively for dichotomous and non-dichotomous items) considering different 22
sets of TEDs. The highest correlations were observed for the sets that included the 23
dimensions, Grammar, Semantics, Register, Information, Construct, and Culture. 24
We also observed different magnitudes of correlations for science and mathematics 25
items and a stronger, statistically significant correlation for translated items whose 26
translation the review panel identified more objectionable than for the rest of the items. 27

[AU2] G. Solano-Flores, Ph.D. (✉)
School of Education, University of Colorado at Boulder,
[AU3] 249 UCB, Boulder, CO 80309, USA
e-mail: Guillermo.Solano@Colorado.Edu

L.Á. Contreras-Niño, Ph.D. • E. Backhoff, Ph.D.
University of Baja California, Mexico, Km. 103 Carretera
Tijuana-Ensenada, c.p. 22830 Ensenada, BC, Mexico
e-mail: angel@uabc.edu.mx

28 These results confirm that language- and content-related TEs may threaten the validity
29 of translated items. They speak to the value using the TTTE as a formative evaluation
30 tool that PISA countries can use to operationalize translation guidelines.

31 Increased awareness of the tremendous sensitivity of tests to language (e.g., Allalouf,
32 2003; Ercikan, 1998; Ercikan, Gierl, McCreith, Puham, & Koh, 2004; Gierl, Rogers,
33 & Klingner, 1999; Grisay, 2007) in the context of international test comparisons has
34 resulted in recent years in substantial improvements of test translation and
35 adaptation procedures used by PISA (e.g., Grisay, 2003; Harkness, van de Vijver,
36 & Mohler, 2003). As part of these improvements, revised sets of test translation
37 guidelines (e.g., Halleux-Monseur, 2008; Hambleton, 1994; Hambleton, Merenda,
38 & Spielberger, 2005; van de Vijver & Poortinga, 2005) have been made available
39 for participating countries.

40 Unfortunately, whereas these revised procedures and guidelines are necessary,
41 their implementation and interpretation may not be optimal without procedures that
42 allow countries to perform detailed, systematic evaluations of their own translation
43 work. Available evidence from research on the testing of linguistically diverse pop-
44 ulations indicates that, in the absence of tools for systematically examining and
45 discussing the linguistic features of items, reviewers may not be able to detect
46 potential linguistic challenges of those items if they rely solely on their judgment
47 (Solano-Flores & Gustafson, *in press*; Solano-Flores, Trumbull, & Kwon, 2003).

48 This need for conceptual tools in test translation led us to propose a theory of test
49 translation error (TTTE; Solano-Flores, Backhoff, & Contreras-Niño, 2009) which
50 defines translation error (TE) (Note 1) as the lack of equivalence between the original
51 language version and the translated version of items. This lack of equivalence can
52 be examined along multiple dimensions that have to do with the design or visual
53 layout of the items (e.g., format, style), their linguistic features (e.g., grammar, syntax),
54 and their content (e.g., information, construct).

55 The theory postulates that error in the translation of tests is inevitable. In addition
56 to a poor translation job, TE is due to factors that are beyond the translators' skills.
57 For example, languages encode meaning differently and have different sets of gram-
58 matical rules. In addition, TE is multidimensional—an error may involve multiple
59 aspects of language (e.g., the lack of a comma is a punctuation error but it also may
60 be an error that alters the intended meaning of a sentence). Due to these reasons, and
61 given the linguistic characteristics that are typical among test items (e.g., limited
62 contextualization, high semantic load of terms, compact sentences), it is virtually
63 impossible to preserve exactly the same meaning and linguistic complexity of items
64 across languages.

65 The notion of test TE as something that cannot be eliminated but can be mini-
66 mized should be easy to understand by professionals in the educational measure-
67 ment community. As with measurement error, TE is due to multiple factors (and their
68 interaction), many of which are beyond control. According to the TTTE, effective
69 test translation can minimize, not eliminate, TE. Flawed translated items have many
70 and/or serious TEs; acceptable translated items have few and/or mild TEs.

Conventional translation review procedures focus on determining whether translated items can be accepted (e.g., Grisay, deJong, Gebhardt, Berezner, & Halleux-Monseur, 2007; Mullis, Kelly, & Haley, 1996). They reflect researchers' and evaluators' tendency to emphasize confirming evidence over disconfirming evidence in hypothesis testing (see Church, 1991; Creswell & Miller, 2000). Unlike conventional translation review procedures, a TTTE-based approach focuses on looking for evidence that disconfirms the notion that the translation of test items is correct. We contend that this approach results in more rigorous translation review procedures.

We have used the TTTE to code errors in translated items and develop measures of TE in those items. Moreover, we have been able to link TE and student performance by correlating item difficulty with measures of TE (Backhoff, Contreras-Niño, & Solano-Flores, 2011; Solano-Flores, Backhoff, & Contreras-Niño, 2006). Our findings have shown consistently that translation review based on the TTTE allows detection of TEs with a level of detail not attained with conventional test translation review procedures (Solano-Flores, Contreras-Niño, & Backhoff, 2005; Solano-Flores, Contreras-Niño, & Backhoff, 2006).

In this chapter, we show how detection and measurement of TE can contribute to improved PISA translation procedures. More specifically, we show how coding and measuring TE based on the TTTE allows identification of serious errors in PISA translated items otherwise regarded as acceptable according to conventional translation verification procedures.

Previous empirical evidence showing the sensitivity to TE of review procedures based on the theory comes from reviews of TIMSS items and relatively small samples of released PISA items (Solano-Flores et al., 2005, 2006). In this study, we reviewed a considerably larger sample of items and took into consideration the structure of many of the PISA items—assessment units consisting of one or several paragraphs with contextual information and one or more items related (see Bybee, McCrae, & Laurie, 2009).

8.1 Theoretical Framework

8.1.1 Definition of Translation Error

The theory of test translation error (TTTE) is not only about errors made in test translation, but also about errors in translated tests. According to our theory (Solano-Flores et al., 2009), *test translation* does not refer exclusively to the action of translating items but also to multiple aspects of the entire process through which translated versions of those items are created. *Translation error* does not result exclusively from poor translation job (e.g., inaccuracy of a chosen term, word-by-word translation, use of false cognates); it also results from factors that are beyond the translators' translation skills.

t1.1 **Table 8.1** Acceptability-objectionability of translated
t1.2 items according to the frequency and severity of test
t1.3 translation errors

t1.4	Mild errors	Severe errors
t1.5 Few errors	Acceptable	Questionable
t1.6 Many errors	Questionable	Objectionable

110 An example of these factors is the natural, well-known fact that no two languages
111 in the world encode meaning in the same way (see [Greenfield, 1997](#); [Nettle &](#)
112 [Romaine, 2000](#)). While the translators' job is to ensure that meaning is preserved in
113 their translations, in some cases this is accomplished at the cost of increasing the
114 amount of text. Unlike other forms of text, this is not trivial matter in tests, which
115 students usually need to respond to within certain time limits. Under these circum-
116 stances, a substantial increase in the amount of text in an item needed to express the
117 same idea as in its original version may imply more reading time and a potential
118 impact on the time students are left with to make sense of the item.

119 Another example of aspects beyond the translators' translation skills has to do with
120 the formatting of translated items. Changes in font size and style, and alterations in the
121 proportion of figures included in test items are not due to the translators' actions yet
122 affect the equivalence between the original and the translated versions of an item.

123 A third example of aspects beyond the translators' translation skills is the extent
124 to which the items reflect the culture of the target language country. While, technically,
125 the translation of an item may not be flawed, the contextual information used in it
126 may not be as familiar to the population tested in the target language as it is to the
127 population tested in the source language.

128 **8.1.2 Inevitability of Translation Error**

129 As a result of the combination of multiple factors like these, strictly speaking, a
130 translation cannot be expected to be perfect. Indeed, our findings from reviews of
131 translated items show that the majority of translated items have TEs—although they
132 are not necessarily fatally flawed ([Backhoff et al., 2011](#); [Solano-Flores et al., 2006](#);
133 [Solano-Flores et al., 2009](#)).

134 **8.1.3 Objectionability of Translated Items**

135 To what extent a translated item is objectionable or acceptable depends on the rela-
136 tion between the frequency and severity of TEs. This relationship is represented in
137 [Table 8.1](#). Acceptable translated items have few mild TEs. Questionable translated
138 items have many mild errors or few severe TEs; they may or may not affect student
139 performance depending on the nature of the TEs, the characteristics of the item, and

the characteristics of the linguistic group tested. Objectionable translated items have many and severe TEs; they are very likely to alter the intended meaning of the original item and affect student performance.

8.1.4 Translation Error Dimensions 143

Our theory postulates the existence of test translation error dimensions (TEDs), grouped in three broad categories, Design, Language, and Content. Each TED comprises several types of TE, as shown in Table 8.2. While it parallels the systems of dimensions and types of TEs used in other investigations (e.g., Backhoff et al., 2011; Solano-Flores et al., 2009), the definitions of TEDs shown in Table 8.2 and the types of TE they comprise were respectively adapted and included with the intent to meet the needs of this particular translation review project.

The TEDs, Style, Format, and Conventions, are grouped in the category, Design. These TEDs have to do with the format, editorial features, and visual layout of translated tests. Convention errors are mainly observed in multiple-choice items. TEs belonging to the category, Design tend to be mild and are unlikely to impact student performance (Note 2).

The TEDs, Grammar, Semantics, and Register, are grouped in the category, Language. These TEDs have to do with the structural and functional aspects of the language used in the translation, the preservation of meaning across languages, and the characteristics of the language usage by the target population in social and instructional contexts.

The TEDs, Information, Construct, Culture, and Origin, are grouped in the category, Content. These TEDs have to do with the ways in which information is presented and how examinees are likely to understand and make sense of items. Unlike TEs belonging to the category Design, TEs belonging to the category, Content tend to alter the structural and functional aspects of language or the ways in which examinees make sense of items. Therefore, they tend to be severe and constitute a threat to the validity of a translated item. The TED, Origin addresses the fact that examining the linguistic equivalence of items allows detection of errors not detected throughout the entire process of test development of the item (Solano-Flores, Trumbull, & Nelson-Barber, 2002). Since Origin errors are not exclusive to the translated version of test items, they are included in the list of TEDs only for conceptual purposes, to allow documentation of any anomalies identified during the process of test translation review.

8.1.5 Translation Error Multidimensionality 172

The theory postulates that test TE is multidimensional. For example, the inappropriate use of commas in *the panda eats, shoots, and leaves* (when the intended meaning is, *the panda eats shoots and leaves*) (Note 3) is both a punctuation error (Style TED) and an error that affects the meaning of the sentence (Semantics TED).

t2.1	Table 8.2 Translation error dimensions and types of translation errors (<i>italics</i>) considered in the analysis of translated PISA textual analytical units (TAUs)
t2.2	
t2.3	<i>Design dimensions</i>
t2.4	Style: The style used in the translation of the TAU is not used in printed materials in the
t2.5	country.
t2.6	• <i>Punctuation • spelling • wrong use of uppercase letter • wrong use of lowercase letter</i>
t2.7	Format: The visual layout of the translated TAU is different from the original.
t2.8	• <i>Change of size, position, or style of an illustration, table, or graph • change of justification,</i>
t2.9	<i>font, or font size of text • change of margin width • omission of graphic component</i>
t2.10	• <i>insertion of graphic component</i>
t2.11	Conventions: The translation of the TAU does not reflect item writing conventions used in the
t2.12	country.
t2.13	• <i>Inconsistent syntactical structure of stem and options • wrong use of punctuation in the item's</i>
t2.14	<i>stem • change in order of options • inconsistent syntactical structure among options • wrong</i>
t2.15	<i>use of uppercase letters in options</i>
t2.16	<i>Language dimensions</i>
t2.17	Grammar: The translation of the TAU violates grammatical rules or uses grammatical
t2.18	structures that are not common in the country.
t2.19	• <i>Literal translation • unnatural syntax of a sentence • subject-verb inconsistency • singular-</i>
t2.20	<i>plural inconsistency • wrong preposition • wrong tense • conflation of sentences</i>
t2.21	Semantics: The translation of the TAU alters its original meaning.
t2.22	• <i>Use of a false cognate • wrong translation or adaptation of an idiomatic expression</i>
t2.23	• <i>alteration of meaning • confusing translation of a sentence • multiple possible interpreta-</i>
t2.24	<i>tions of a sentence • change of gender of a character • conflation of ideas • inaccurate terms</i>
t2.25	• <i>use of terms with multiple meanings • wrong translation of a word</i>
t2.26	Register: The translation of the TAU does not reflect the terms, idiomatic expressions, and
t2.27	discursive forms used in the country.
t2.28	• <i>Use of words of low frequency in the country • wrong translation of a technical term</i>
t2.29	• <i>translation of a technical term in a way not used in the country</i>
t2.30	<i>Content dimensions</i>
t2.31	Information: The translation of the TAU alters the amount, precision, or type of information
t2.32	provided.
t2.33	• <i>Inconsistent translation of a non-technical term • change in number of times a technical term</i>
t2.34	<i>is used • insertion of technical term • insertion of a sentence or explanation • omission of a</i>
t2.35	<i>key word • omission of a technical term • omission of a sentence or explanations</i>
t2.36	Construct: The type of skill or knowledge needed to understand and respond to the TAU is
t2.37	different from the skill or knowledge needed to understand and respond to the TAU in the
t2.38	source language.
t2.39	• <i>Possible change of the item's cognitive demands • possible alteration of ways in which a task</i>
t2.40	<i>may be interpreted • wrong technical term • inconsistent translation of a technical term</i>
t2.41	• <i>undue insertion of a technical term • omission of a technical term • translation of a technical</i>
t2.42	<i>term as a non-technical term • translation of a non-technical term as a technical term</i>
t2.43	Culture: The TAU does not reflect the characteristics of the culture or the curriculum in the
t2.44	target language.
t2.45	• <i>Contextual information and situations that are uncommon in the country • measurement units</i>
t2.46	<i>not used in the country • problem posed not meaningful in the country's culture • knowledge</i>
t2.47	<i>assessed not taught in country</i>
t2.48	Origin: The TAU carries over errors from the source language version.
t2.49	• <i>Inconsistency in the content of the two source languages • conceptual errors in the design of</i>
t2.50	<i>the item • confusing directions • the answer to an item may give the clue for responding to</i>
t2.51	<i>another item within the same assessment unit</i>

8.1.6 Tension Among Translation Error Dimensions 177

Finally, the theory postulates that there is a tension between TEDs. Actions intended to avoid TE on a given TED may involve making errors on other TEDs. For example, the grammatical rules of the target language may prevent a noun from being repeated in the same sentence. In some languages, a marker needs to be used to refer to a noun in the rest of the sentence, once the noun appears in it. As a consequence, a key technical term that appears several times in the same sentence in the original version of the item appears only once in its translation. The grammatical rules of the target language need to be followed at the cost of altering the number of times that the key term appears in the sentence—which alters the amount of information provided by the item.

8.2 Methods 188

8.2.1 Sample of Assessment Units and Analytical Test Units 189

We examined 61 assessment units (one or several paragraphs with contextual information and one or more items related) from the Mexican, Spanish language version of PISA-2006. Of these 61 assessment units, 37 and 24 were respectively science and mathematics assessment units (Note 4). These 61 assessment units comprised a total of 193 text analytical units (TAUs), defined as either the introductory text or an item within an assessment unit. Of the 193 TAUs examined, 55 were introductory texts and 138 were items. Of these 138 items, 101 and 37 were respectively science and mathematics items.

8.2.2 Test Translation Review and Error Coding Procedures 198

In addition to the fact that most of the PISA 2006 items consisted of two forms of TAUs (an introductory text or an item), our coding procedure took into account that PISA items use two source languages, English and French (see [Grisay et al., 2007](#)).

We assembled a multidisciplinary translation review panel composed of three middle school teachers (Spanish, science, and mathematics); three high school teachers (Spanish, science, mathematics); one English-to-Spanish translator, and one French-to-Spanish translator (both certified by international translation professional organizations); one linguist; one test developer; and one psychometrician (measurement specialist).

The following procedure was used to review each TAU. First, the TAU in the target language (the translated item) was projected on a screen. Reviewers read the TAU and, in the case of items, responded to the item individually as if they were

211 students taking the test. This was done with the purpose of giving the reviewers the
212 opportunity to become acquainted with the content of the item and to become aware
213 of its cognitive and linguistic demands in the target language.

214 The reviewers then were asked to examine the TAU and individually record on a
215 coding form all the types of TE they thought could affect the interpretation of the
216 item. The reviewers were instructed to focus on a specific set of dimensions desig-
217 nated according to their professional background. However, they were allowed to
218 record errors on all dimensions (Table 8.3).

219 Once the reviewers finished recording their comments, the original English and
220 French versions of the TAU were projected on two additional screens. Then the
221 reviewers were asked to compare the English and French versions with the TAU in
222 the target language and to individually code any type of TE according to the list of
223 types of errors listed above for each error dimension. They also wrote their com-
224 ments on the TAU based on their experience reading and responding to it and on
225 comparing the original and translated versions.

226 For each TED, the panel discussed each reviewer's coding. Project staff facilit-
227 ated a discussion to ensure that the panel decided by consensus what errors should
228 be recorded and on which TEDs they should be coded. In the case of items, the
229 panel was asked to decide, based on the number and severity of the TEs, if the trans-
230 lated item should be classified as objectionable (i.e., an item with many and severe
231 TEs which were likely to adversely affect student performance). The review coding
232 decisions were captured on an electronic spreadsheet for further analysis.

233 **8.2.3 Data Analysis**

234 For the purpose of our analysis, we measured TE in each TAU as the number of dif-
235 ferent translation error dimensions (NDTED) on which TEs were observed in it.
236 This coarse-grain measure has proven to be sensitive to important differences in
237 translation quality among items (see Solano-Flores et al., 2005, 2006).

238 Also for the purpose of our analysis, we used the p-values of items as a measure
239 of item difficulty. Item p-value was computed as the proportion of the item's highest
240 possible score (see Adams, Bereznier, & Jakubowski, 2010), which allowed to have
241 proportional measures of difficulty for both dichotomous and partial-credit items.
242 More specifically, for dichotomous items, difficulty was computed as the proportion
243 of students who responded correctly; for partial-credit items, difficulty was com-
244 puted as the mean score of the item divided by its maximum score.

245 To examine the impact of TE on student performance, we examined the Pearson
246 correlations between NDTED and item p-value for different sets of TEDs, different
247 content areas (science and mathematics), and items that were and were not identified
248 as objectionable by the translation review panel. Impact on performance should be
249 observed as a negative correlation.

250 Given the complex interaction of the students' knowledge of the content being
251 assessed and the cognitive and linguistic demands of test items, it would be naive to

Table 8.3 Expertise provided and assigned focus on translation error dimensions by specialist: 1 = main role; 0 = adjuvant role

Expertise and Contribution	Style	Format	Conventions	Information	Grammar	Semantics	Construct	Register	Culture	Origin
t3.1										
t3.2										
t3.3	0	1	1	1	0	0	0	1	1	0
t3.4										
t3.5										
t3.6										
t3.7	1	1	0	1	1	1	0	0	0	0
t3.8										
t3.9										
t3.10	1	0	0	0	1	1	1	0	1	0
t3.11										
t3.12										
t3.13										
t3.14	0	1	1	1	0	0	1	0	0	0
t3.15										
t3.16										
t3.17	0	0	1	1	0	1	1	0	0	0
t3.18										
t3.19										

Un-corrected Proof

t4.1 **Table 8.4** Percentage of text
 t4.2 analytical units (n=193) with
 t4.3 at least one error on each of
 t4.4 the translation error
 t4.5 dimensions

Dimension	Percent	
Style	48	t4.6
Format	53	t4.7
Conventions	3	t4.8
Information	53	t4.9
Grammar	53	t4.10
Semantics	78	t4.11
Construct	35	t4.12
Register	21	t4.13
Culture	5	t4.14
Origin	41	t4.15

252 expect to observe impressively high and statistically significant correlations. Rather,
 253 we expected to observe patterns in those correlations that would indicate a systematic
 254 impact of TE on student performance, especially for language- and content-related
 255 TEDs and for items identified as objectionable by the translation review panel.

256 **8.3 Results**

257 **8.3.1 Frequency and Severity of Translation Errors**

258 We observed TEs on at least one dimension for almost all (96%) of the TAUs. Of the 138
 259 TAUs which consisted of items, 26 were identified by the committee as objectionable.

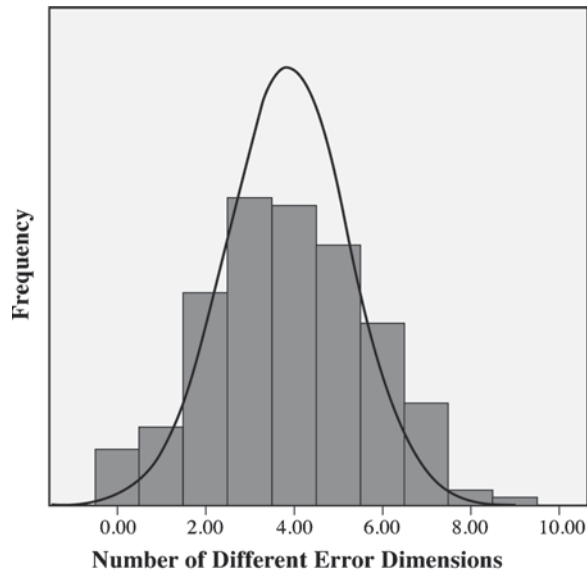
260 Table 8.4 shows the percentage of TAUs identified as having at least one error on
 261 each of the TEDs. As indicated above, many of these errors are not likely to bias test
 262 results and many are even difficult to be noticed by individuals who have no experience
 263 reviewing test translations. On the other hand, there are TEs that may potentially
 264 threaten the validity of test items. Such is the case for errors on the Semantics,
 265 Grammar, and Information dimensions, which were observed respectively in 78%,
 266 53% and 53% of the TAUs.

267 On average, a TAU had errors on 3.9 different dimensions (s.d.=1.834). As
 268 Fig. 8.1 shows, the number of different dimensions in which error was observed had
 269 a normal frequency distribution.

270 **8.3.2 Translation Error and Item Difficulty**

271 As Table 8.5 shows, Pearson correlation coefficients of $-.059$ and $-.117$ between
 272 NDTED and p-value were observed respectively when all dimensions were consid-
 273 ered and when the three language dimensions (Grammar, Semantics, and Register)
 274 and three of the four content dimensions (Information, Construct, and Culture)

Fig. 8.1 Frequency distribution of text analytical units by number of different error dimensions



Comparison	Correlation	
By set of dimensions (138)		t5.1
All dimensions	-.059	t5.2
Language and content dimensions ^a	-.117	t5.3
By content area ^a (138)		t5.4
Science (n=101)	-.115	t5.5
Mathematics (n=37)	-.213	t5.6
By objectionability (language and content dimensions) ^a (138)		t5.7
Non-objectionable items (n=112)	-.084	t5.8
Objectionable items (n=26)	-.404**	t5.9
Sample and subsample sizes in parentheses		t5.10
^a Includes the three language dimensions (Grammar, Semantics, and Register) and three of the four content dimensions (Information, Construct, and Culture)		t5.11
**Significant at p=.01 (2-tailed)		t5.12

were considered. (As mentioned above, since Origin errors are common to both the source and language versions of items, they were not included in the analyses). This difference supports findings from previous test translation reviews that design dimensions (Style, Format, and Conventions) are unlikely to affect student performance whereas language and content dimensions tend to have a greater impact on student performance and may potentially threaten the validity of translated items.

275
276
277
278
279
280

281 Correlation coefficients of $-.115$ and $-.213$ between NDTED and p-value were
282 observed respectively for the science and mathematics items. These results are consistent
283 with findings from other translation reviews, in which we (e.g., [Solano-Flores,
284 Backhoff, & Contreras-Niño, 2005](#)) have observed higher correlations between
285 NDTED and item difficulty for mathematics than science items.

286 Correlation coefficients of $-.084$ and $-.404$ (significant) were observed respec-
287 tively for acceptable and objectionable items. This considerable difference indicates
288 that the review procedure allows identification of items which have sets of errors
289 that are likely to seriously impact student performance. This finding is important,
290 considering that the number of items identified as objectionable (26) constitute
291 about 19% of the 138 items examined.

292 **8.4 Summary and Conclusions**

293 The theory of test translation error (TTTE; [Solano-Flores et al., 2009](#)) postulates the
294 existence of translation error dimensions (TEDs; e.g., Semantics, Construct, Grammar)
295 and views translation error (TE) as multidimensional (a translation error can
296 belong to several TEDs). It also postulates that a tension exists between TEDs
297 (i.e., in translating a test item, avoiding error on one dimension may produce error
298 on other dimensions). Accordingly, error-free test translation is impossible; effective
299 test translation minimizes but does not eliminate error. The theory also postulates
300 that while items usually have multiple TEs, most of them are mild and even unnoticeable.
301 Objectionable translated items have many and severe TEs and are likely to pose
302 serious linguistic challenges to examinees who are given the translated version of
303 a test.

304 In this chapter, we report the results of our review of the Spanish language
305 Mexican version of PISA-2006 science and mathematics text analytical units
306 (TAUs). Consistent with results from our review of the Spanish Mexican translation
307 of TIMSS-1995 ([Solano-Flores et al., 2005](#)) and the Spanish Mexican translation of
308 PISA-2003 ([Backhoff et al., 2011](#)), our results show that translation reviews based
309 on the TTTE are highly sensitive to TE.

310 The results also confirm previous findings that student performance tends to
311 be resilient to TE on design-related TEDs and sensitive to TE on language- and
312 content-related TEDs. Also, items whose translation was identified as objectionable
313 by the review panel correlated higher with item difficulty than items whose transla-
314 tion was not identified as objectionable—a finding that speaks to the sensitivity of
315 TTTE-based judgmental review procedures.

316 A limitation of our analyses of correlations of measures of TE and item difficulty
317 stems from the fact that we did not account for the effect of TE observed in the
318 introductory text of assessment units. Future research should explore models for
319 examining this relationship.

320 Unlike other approaches created to examine translation quality, the TTE focuses
321 on disconfirming (rather than confirming) evidence that the translation of test items

is correct. In addition, because they use multidisciplinary review panels which discuss the linguistic features of the items at length, TTTE-based coding procedures are sensitive to TE with a level of precision and detail not attained with conventional approaches.

Experienced test translators who have attended our workshops on the use of the TTTE and the methods described in this chapter (e.g., Backhoff, Solano-Flores, & Contreras-Niño, 2010; Solano-Flores et al., 2010) react initially with skepticism when we report our findings. They find it difficult to believe that items translated according to available translation guidelines have multiple TEs. It is not until they observe the discussions of the review panels examining specific translated items that they appreciate the level of sensitivity of the theory and our coding procedures to the nuances of language in translated items.

As with measurement error, TE cannot be entirely eliminated, but it can be minimized. As our results show, a theoretical perspective that assumes error inevitability in test translation is more sensitive to the complexities of language in translated PISA items and can contribute to the improvement of future PISA translation procedures. We hope that, in the future, PISA participating countries use our approach as a tool for operationalizing PISA translation procedures and formatively evaluating their own translation work.

[AU4] **Author's Note** 341

Portions of this paper are originally from a paper presented at the PISA Research Conference, 14–16 September 2009, Kiel, Germany. The investigation reported in this paper was commissioned and funded by the National Institute for Educational Evaluation (INEE), Mexico, and conducted through a contract with the Autonomous University of Baja California (UABC), Mexico. The opinions expressed are not necessarily those of the funding agency. Contact author: Guillermo Solano-Flores, guillermo.solano@colorado.edu.

Notes 349

- Note 1. While *translation error* (in singular) is used here to refer to lack of equivalence between the original language version and the translated version of an item, *translation errors* (in plural) or *a translation error* are used to refer to specific instances or types of translation error (e.g., the inaccurate translation of a term or an inappropriate use of punctuation). 350-353
- Note 2. Of course, there are exceptions. For example, an alteration in the proportion of the length of the axes in a graph showing a functional relationship may make the line of the function look steeper in the translated item than in the original—which may affect how the examinee interprets the function. 354-357
- Note 3. The example is based the story told by Lynne Truss (2004) at the beginning of her well-known book on punctuation, *Eats, shoots, and leaves*. 358-359
- Note 4. One of the science assessment units and 17 of the mathematics assessment units consisted of a stand-alone item with no introductory text. 360-361

362 **References**

- 363 Adams, R., Bereznar, A., & Jakubowski, M. (2010). *Analysis of PISA 2006 preferred items ranking*
 364 *using the percent-correct method* (OECD Education Working Papers, No. 46). OECD
 365 Publishing. Retrieved June 7, 2011, <http://dx.doi.org/10.1787/5km4psmntkq5-en>
- 366 Allalouf, A. (2003). Revising translated differential item functioning items as a tool for improving
 367 cross-lingual assessment. *Applied Measurement in Education*, 16(1), 55–73.
- 368 Backhoff, E., Contreras-Niño, L. A., & Solano-Flores, G. (2011). *The theory of test translation*
 369 *error and the TIMSS and PISA international test comparisons*. Mexico: National Institute for
 370 Educational Evaluation [Sp.].
- 371 Backhoff, E., Solano-Flores, G., & Contreras-Niño, L. A. (2010, February 18–19). *Analysis of the*
 372 *Mexican Spanish language translation of PISA-2006*. Presentation at the Ibero-American
 373 Seminar on the theory of test translation error in international comparisons. National Ministry
 374 of Education and National Institute for Educational Evaluation, Mexico City, Mexico.
- 375 Bybee, R., McCrae, B., & Laurie, R. (2009). PISA 2006: An assessment of scientific literacy.
 376 *Journal of Research in Science Teaching*, 46(8), 865–883.
- 377 Church, B. (1991). An examination of the effect that commitment to a hypothesis has on auditors’
 378 evaluations of confirming and disconfirming evidence. *Contemporary Accounting Research*,
 379 7(2), 513–534.
- 380 Creswell, J. W., & Miller, D. L. (2000). Determining validity in qualitative inquiry. *Theory into*
 381 *Practice*, 39(3), 124–130.
- 382 Ercikan, K. (1998). Translation effects in international assessment. *International Journal of*
 383 *Educational Research*, 29, 543–553.
- 384 Ercikan, K., Gierl, M. J., McCreith, T., Puham, G., & Koh, K. (2004). Comparability of bilingual
 385 versions of assessments: Sources of incomparability of English and French versions of Canada’s
 386 national achievement tests. *Applied Measurement in Education*, 17(3), 301–321.
- 387 Gierl, M. J., Rogers, W. T., & Klinger, D. (1999, April). *Using statistical and judgmental reviews*
 388 *to identify and interpret translation DIF*. Paper presented at the annual meeting of the National
 389 Council on Measurement in Education, Montreal, QC.
- 390 Greenfield, P. M. (1997). You can’t take it with you: Why ability assessments don’t cross cultures.
 391 *American Psychologist*, 52(10), 1115–1124.
- 392 Grisay, A. (2003). Translation procedures in OECD/PISA 2000 international assessment. *Language*
 393 *Testing*, 20(2), 225–240.
- 394 Grisay, A. (2007). *The challenge of adapting PISA materials into non Indo-European languages:*
 395 *Some evidence from a brief exploration of language issues in Chinese and Arabic*. OECD Core
 396 A Consortium.
- 397 Grisay, A., de Jong, J. H., Gebhardt, E., Bereznar, A., & Halleux-Monseur, B. (2007). Translation
 398 equivalence across PISA countries. *Journal of Applied Measurement*, 8(3), 249–266.
- 399 Halleux-Monseur, B. (2008). *Translation, adaptation and verification of test material in OECD*
 400 *international surveys*. Paris: Directorate for Education, Institutional Management in Higher
 401 Education Governing Board, OECD.
- 402 Hambleton, R. K. (1994). Guidelines for adapting educational and psychological tests: A progress
 403 report. *European Journal of Psychological Assessment*, 10(3), 229–244.
- 404 Hambleton, R. K., Merenda, P. F., & Spielberger, C. D. (Eds.). (2005). *Adapting educational and*
 405 *psychological tests for cross-cultural assessment*. Mahwah, NJ: Lawrence Erlbaum Associates,
 406 Publishers.
- 407 Harkness, J., van de Vijver, F. J. R., & Mohler, P. (Eds.). (2003). *Cross-cultural survey methods*.
 408 Hoboken, NJ: Wiley.
- 409 Mullis, I. V. S., Kelly, D. L., & Haley, K. (1996). Translation Verification Procedures. In M. O.
 410 Martin & I. V. S. Mullis (Eds.), *Third international mathematics and science study: Quality*
 411 *assurance in data collection*. Chestnut Hill, MA: Boston College.
- 412 Nettle, D., & Romaine, S. (2000). *Vanishing voice: The extinction of the world’s languages*. New
 413 York: Oxford University Press.

	Solano-Flores, G., Backhoff, E., & Contreras-Niño, L. A. (2006). <i>Methodology for evaluating the quality of test translations in international test comparisons: The case of Mexico, TIMSS-1995</i> . [Sp.]. Mexico: National Institute for Educational Evaluation (INEE).	414 415 416
	<u>Solano-Flores, G., Backhoff, E., & Contreras-Niño, L. A. (2009). Theory of test translation error. <i>International Journal of Testing</i>, 9, 78–91.</u>	417 418
	Solano-Flores, G., Backhoff, E., & Contreras-Niño, L. A. (2010, February 18–19). <i>Test translation review sessions: A demonstration</i> . Presentation at the Ibero-American Seminar on the theory of test translation error in international comparisons. National Ministry of Education and National Institute for Educational Evaluation, Mexico City, Mexico.	419 420 421 422
	<u>Solano-Flores, G., Contreras-Niño, L. A., & Backhoff, E. (2005, April 12–14). <i>The Mexican translation of TIMSS-95: Test translation lessons from a post-mortem study</i>. Paper presented at the 2005 annual meeting of the National Council on Measurement in Education. Montreal, Quebec, Canada.</u>	423 424 425 426
	<u>Solano-Flores, G., Contreras-Niño, L. A., & Backhoff, E. (2006). Test translation and adaptation: Lessons learned and recommendations for countries participating in TIMSS, PISA, and other international comparisons. <i>REDIE: Electronic Journal of Educational Research</i>, 8(2). [Sp.] http://redie.uabc.mx/vol8no2/contents-solano2.html</u>	427 428 429 430
[AU5]	Solano-Flores, G., & Gustafson, M. (In Press). Assessment of English language learners: A critical, probabilistic, systemic view. In M. Simon, K. Ercikan, & M. Rousseau (Eds.), <i>Improving large scale assessment in education: Theory, issues, and practice</i> . Taylor & Francis, Routledge.	431 432 433
	Solano-Flores, G., Trumbull, E., & Kwon, M. (2003, April 21–25). <i>The metrics of linguistic complexity and the metrics of student performance in the testing of English language learners</i> . Symposium paper presented at the 2003 Annual Meeting of the American Evaluation Research Association. Chicago.	434 435 436 437
	<u>Solano-Flores, G., Trumbull, E., & Nelson-Barber, S. (2002). Concurrent development of dual language assessments: An alternative to translating tests for linguistic minorities. <i>International Journal of Testing</i>, 2(2), 107–129.</u>	438 439 440
	Truss, L. (2004). <i>Eats, shoots & leaves</i> . New York: Gotham Books.	441
	van de Vijver, F. J. R., & Poortinga, Y. H. (2005). Conceptual and methodological issues in adapting tests. In R. K. Hambleton, P. F. Merenda, & C. D. Spielberger (Eds.), <i>Adapting educational and psychological tests for cross-cultural assessment</i> . Mahwah, NJ: Lawrence Erlbaum Associates, Publishers.	442 443 444 445

Author's Proof

Please see the Author's responses in the Word document sent along with the proof

Author Queries

Chapter No.: 8 0001528456

Queries	Details Required	Author's Response
AU1	Please provide keywords	
AU2	Please confirm the corresponding author.	
AU3	Please provide the name of the department in the affiliations.	
AU4	Please confirm the placement of Author's Note.	
AU5	Please update Reference Solano–Flores and Gustafson (in press).	

Uncorrected Proof