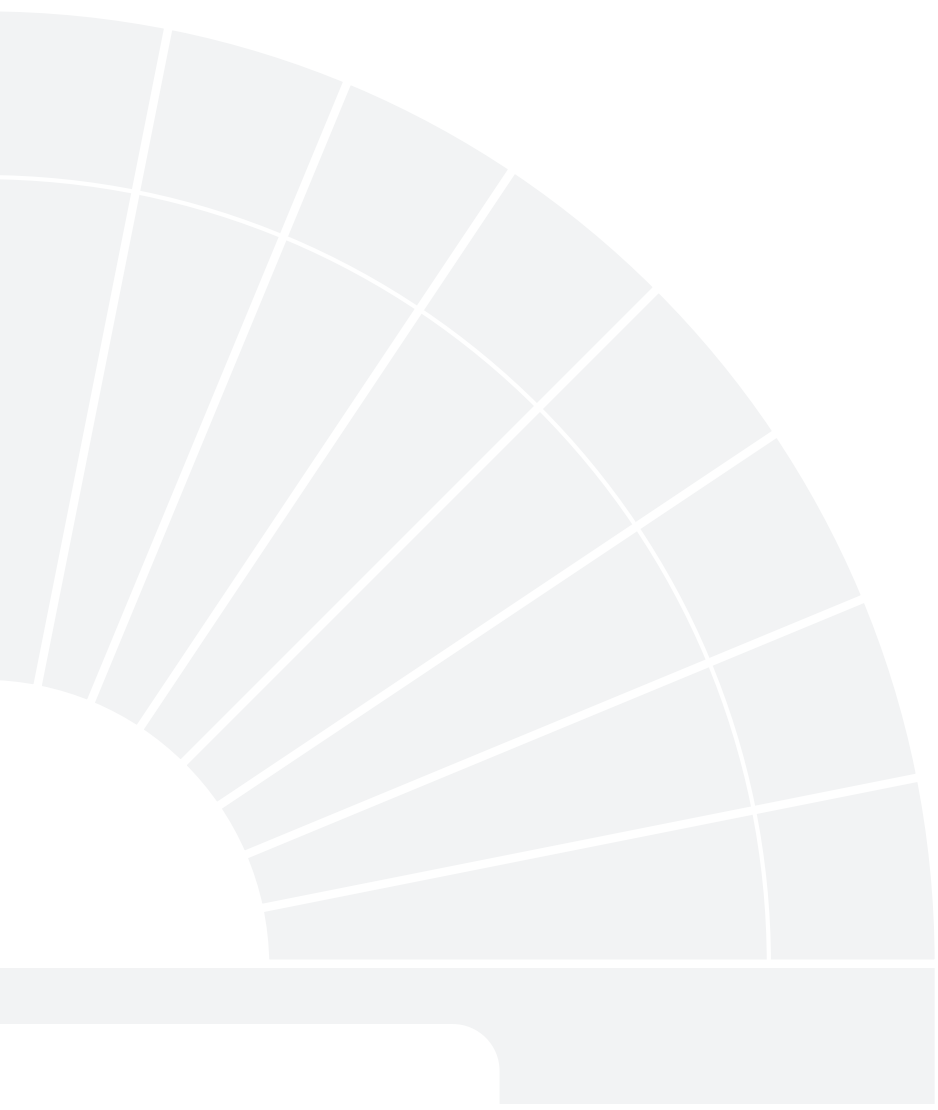# The Test-Curriculum Matching Analysis: Mathematics

C

When comparing student achievement across countries, it is important that the comparisons be as fair as possible. LSS has worked toward this goal in a number of ways, including providing detailed procedures for standardizing the population definitions, sampling, test translations, test administration, scoring, and database formation. Similar to the procedures used for developing the original TIMSS instruments, developing the TIMSS 1999 tests involved a series of reviews by representatives of the participating countries, experts in mathematics, and testing specialists.[1] The National Research Coordinators (NRCs) from each country formally approved the TIMSS 1999 tests, thus accepting them as being sufficiently fair to compare their students' mathematics achievement with that of students from other countries.

Although the tests were developed to represent a set of agreed-upon mathematics content, differences among the curricula of participating countries result in various topics being taught at different grades. To restrict test items to topics included in the curricula of all participating countries and covered in the same sequence would severely limit test coverage and restrict the research questions that the study is designed to address. The tests, therefore, inevitably have some items measuring topics unfamiliar to some students in some countries.

The Test-Curriculum Matching Analysis (TCMA) was conducted to investigate the appropriateness of the TIMSS 1999 mathematics test for the eighth-grade students in the participating countries. TCMA also shows how student performance for individual countries varies when based only on the test questions that are judged to be relevant to their own curricula.[2]

To gather data about the extent to which the TIMSS 1999 tests were relevant to the curricula of the participating countries, each NRC reported whether each item was in that country's intended curriculum at the grade tested. The NRC was asked to choose a person or persons who were very familiar with the curriculum at the grade tested to make this determination. Since an item might be in the curriculum for some but not all students in a country, an item was determined appropriate if it was in the intended curriculum for more than 50 percent of the students. The NRCs had considerable flexibility in selecting items and may have considered items inappropriate for other reasons. All participating countries returned the information for analysis.

Exhibits C.1 and C.2 present the TCMA results for the TIMSS 1999 tests. Exhibit C.1 shows the average percent correct for each country on items selected as appropriate and on the test as a whole. Exhibit C.2 shows the standard errors corresponding to the percentages presented in Exhibit C.1.

C1–C2

---

[1]  See Appendix A for more information on test development.

[2]  Because there may also be curriculum areas covered in some countries that are not covered by the TIMSS 1999 tests, the TCMA does not provide complete information about how well the tests cover the curricula of the countries.

In Exhibit C.1, the last row of the exhibit indicates that the countries varied substantially in the number of items (score points) identified as appropriate.[3] The percentages ranged from 100 percent (169 score points) in Chinese Taipei, the Slovak Republic, Latvia (LSS), the United States, Lithuania, Moldova, and Indonesia to 58 percent (98 score points) in Chile. Thirty-four of the 38 countries indicated that the items representing three-quarters or more of the score points (127 out of a possible 169) were appropriate.

Since most countries indicated that some items were not included in their intended curriculum at the grade tested, the data were analyzed to determine whether the inclusion of these items had any effect on the international performance comparisons.[4]

The first column in Exhibit C.1 shows the average percent correct on all test items for each country. The countries are presented in order of their overall performance based on overall percent correct, from highest to lowest. To interpret this exhibit, reading across a row provides the average percent correct for the students in that country on the items selected by each of the countries listed across the top of the exhibit. For example, Singapore, where the average percent correct was 77 percent on its own set of items, also had 78 percent correct for the items selected by Korea, 77 percent for the items selected by Hong Kong, and so forth. The column for a country listed across the top shows how each of the other countries performed on the subset of items selected as appropriate for its own students. Using the set of items selected by Finland as an example, on average 77 percent of these items were answered correctly by students in Singapore, 73 percent by students in Korea, 72 percent by those in Hong Kong, and so forth. The shaded diagonal element in the exhibit shows how each country performed on the subset of items that it selected based on its own curriculum. Thus, Finnish students averaged 56 percent correct on the set of items identified by Finland for the analysis.

The international averages of each country's selected items are presented across the second to the last row of the exhibit. They show that the selection of items for the participating countries varied somewhat in average difficulty, ranging from 48 to 54 percent. Despite these differences, the overall picture presented by Exhibit C.1 reveals that different item selections do not make a major difference in how well countries perform relative to one another. The items selected by some countries were more difficult than those selected by others. The relative performance of countries on various item selections did vary somewhat, but generally not in a statistically significant manner.[5]

---

[3] Of the 162 items in the test, some items were assigned more score points than others. In particular, some items had two parts, and some extended-response items were scored on a two-point scale. The total number of score points available for analysis was 169. The TCMA uses score points in order to give the same weight to items given them in test scoring.

[4] It should be noted that the performance levels presented in Exhibit C.1 are based on average percents, which are different from the average scale scores that are presented in Chapter 1.

[5] Small differences in performance shown in this exhibit are not statistically significant. The standard errors for the estimated average percent correct statistics are in Exhibit C.2. It can be said with 95 percent confidence that the value for the entire population falls between the sample estimate plus or minus two standard errors.

Comparing the diagonal element for a country with the overall average percent correct shows the difference between performance on the subset of items chosen as appropriate and performance on the test as a whole. In general, there were only small increases in each country's performance on its own subset of items. To illustrate, the average percent correct for Singapore was 77 percent. The diagonal element shows that Singaporean students had the same percent correct (77 percent) based on the smaller set of items selected as they did overall. All countries had a difference of less than five percentage points between the two performance measures, with the largest difference four percent for the Netherlands (65 percent compared with 61 percent).

It is clear that the selection of items does not have a major effect on the general relationship among countries. Countries that had substantially higher or lower relative performance on all items also had higher or lower relative performance on the different sets of items selected for the TCMA. For example, Singapore had the highest average percent correct on the test as a whole and on most of the different item selections, with Korea, Hong Kong, and Chinese Taipei among the four highest-performing countries in all cases. Although there are some changes in the ordering of countries based on the items selected for the TCMA, most of these differences are within the boundaries of sampling error. As an example, consider the 149 score points selected by Jordan. The Jordanian students did better on these items than on the test as a whole, with 41 percent correct on these items, on average, compared with 38 percent correct on all items. However, most other countries also did better on these particular items, with an international average of 54 percent correct on the items selected by Jordan. All 30 countries that performed better than Jordan on the overall test also performed better on the items selected by Jordan.

The TCMA results provide evidence that the TIMSS 1999 mathematics test provides a reasonable basis for comparing achievement of the participating countries. This result is not unexpected, since making the test as fair as possible was a major consideration in test development. The fact that the majority of countries indicated that most items were appropriate for their students means that the different average percent correct estimates were based on essentially the same items. Insofar as countries rejected items that would be difficult for their students, these items tended to be difficult for students in other countries as well. The analysis shows that omitting such items tends to improve the results for that country, but also tends to improve the results for all other countries, so that the overall pattern of results is largely unaffected.

**Exhibit C.1   Average Percent Correct for Test-Curriculum Matching Analysis – Mathematics**

Based on Subsets of Items Specially Identified by Each Country as Addressing its Curriculum
(See Exhibit C.2 for corresponding standard errors)

TIMSS 1999 · 8th grade Mathematics

**Instructions:** Read **across** the row to compare that country's performance based on the test items included by each of the countries across the top. Read **down** the column under a country name to compare the performance of the country down the left on the items included by the country listed on the top. Read along the **diagonal** to compare performance for each different country based on its own decisions about the test items to include.

| Countries | Average Percent Correct on All Items | Singapore | Korea, Rep. of | Hong Kong, SAR | Chinese Taipei | Japan | Belgium (Flemish) | Hungary | Slovak Republic | Netherlands | Slovenia | Canada | Russian Federation | Australia | Malaysia | Czech Republic | Finland | Bulgaria | Latvia (LSS) | United States | England | New Zealand | Lithuania | Italy | Cyprus | Romania | Moldova | Israel | Thailand | Macedonia, Rep. of | Tunisia | Jordan | Turkey | Iran, Islamic Rep. | Indonesia | Chile | Philippines | Morocco | South Africa |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Singapore | 77 (1.3) | 77 | 78 | 77 | 77 | 77 | 77 | 78 | 77 | 77 | 77 | 77 | 78 | 77 | 78 | 77 | 77 | 78 | 77 | 77 | 78 | 77 | 77 | 77 | 77 | 77 | 77 | 78 | 78 | 77 | 77 | 78 | 77 | 77 | 77 | 80 | 79 | 77 | 77 |
| Korea, Rep. of | 73 (0.4) | 72 | 73 | 72 | 72 | 72 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 73 | 72 | 71 | 73 | 73 | 73 | 73 | 74 | 72 | 76 | 73 | 73 | 72 | 73 | 73 | 74 | 72 |
| Hong Kong, SAR | 72 (1.1) | 72 | 73 | 72 | 72 | 72 | 72 | 73 | 72 | 75 | 72 | 72 | 73 | 72 | 73 | 72 | 72 | 73 | 72 | 72 | 73 | 73 | 72 | 72 | 71 | 72 | 72 | 73 | 73 | 74 | 72 | 74 | 73 | 72 | 72 | 73 | 73 | 73 | 72 |
| Chinese Taipei | 72 (0.9) | 72 | 73 | 72 | 72 | 72 | 72 | 72 | 72 | 74 | 72 | 72 | 72 | 72 | 72 | 72 | 72 | 73 | 72 | 72 | 73 | 73 | 72 | 72 | 70 | 72 | 72 | 73 | 73 | 73 | 72 | 74 | 73 | 72 | 72 | 73 | 72 | 72 | 72 |
| Japan | 72 (0.4) | 71 | 71 | 71 | 72 | 72 | 72 | 72 | 72 | 74 | 72 | 72 | 72 | 72 | 72 | 72 | 72 | 71 | 72 | 72 | 73 | 72 | 72 | 72 | 72 | 72 | 72 | 71 | 74 | 72 | 72 | 74 | 72 | 72 | 72 | 73 | 72 | 71 | 71 |
| Belgium (Flemish) | 66 (0.8) | 65 | 66 | 65 | 66 | 66 | 66 | 67 | 66 | 68 | 66 | 66 | 66 | 67 | 66 | 66 | 66 | 67 | 66 | 66 | 68 | 67 | 66 | 66 | 68 | 66 | 66 | 66 | 67 | 66 | 66 | 69 | 67 | 66 | 66 | 66 | 67 | 66 | 65 |
| Hungary | 61 (0.9) | 61 | 62 | 61 | 61 | 61 | 61 | 62 | 61 | 63 | 61 | 61 | 63 | 62 | 62 | 61 | 62 | 62 | 61 | 61 | 62 | 62 | 61 | 61 | 59 | 61 | 61 | 62 | 62 | 62 | 61 | 65 | 62 | 62 | 61 | 61 | 63 | 61 | 61 |
| Slovak Republic | 61 (1.1) | 61 | 61 | 61 | 61 | 61 | 61 | 62 | 61 | 64 | 61 | 61 | 62 | 62 | 62 | 61 | 60 | 62 | 61 | 61 | 62 | 61 | 61 | 61 | 59 | 61 | 60 | 61 | 62 | 62 | 61 | 64 | 61 | 62 | 61 | 61 | 63 | 60 | 61 |
| Netherlands | 61 (1.8) | 60 | 60 | 60 | 61 | 60 | 61 | 60 | 61 | 65 | 60 | 61 | 60 | 62 | 60 | 61 | 61 | 60 | 61 | 61 | 63 | 62 | 61 | 61 | 55 | 60 | 61 | 60 | 60 | 61 | 60 | 64 | 61 | 60 | 61 | 61 | 61 | 60 | 59 |
| Slovenia | 60 (0.7) | 59 | 60 | 59 | 60 | 59 | 60 | 59 | 60 | 62 | 59 | 60 | 58 | 60 | 59 | 60 | 60 | 59 | 60 | 60 | 59 | 60 | 60 | 59 | 57 | 60 | 60 | 59 | 59 | 59 | 58 | 63 | 59 | 60 | 59 | 59 | 62 | 60 | 59 |
| Canada | 59 (0.5) | 58 | 58 | 58 | 58 | 59 | 59 | 59 | 59 | 62 | 58 | 59 | 58 | 58 | 58 | 58 | 56 | 60 | 59 | 58 | 61 | 60 | 60 | 59 | 60 | 58 | 59 | 59 | 59 | 59 | 58 | 62 | 59 | 59 | 58 | 59 | 60 | 59 | 58 |
| Russian Federation | 58 (1.3) | 58 | 59 | 58 | 58 | 58 | 58 | 59 | 58 | 59 | 58 | 58 | 60 | 58 | 59 | 58 | 56 | 60 | 58 | 58 | 58 | 59 | 58 | 58 | 57 | 58 | 58 | 59 | 60 | 59 | 58 | 62 | 59 | 59 | 58 | 57 | 60 | 58 | 57 |
| Australia | 58 (1.2) | 57 | 58 | 57 | 58 | 57 | 58 | 57 | 57 | 61 | 57 | 58 | 57 | 59 | 58 | 57 | 58 | 57 | 58 | 58 | 60 | 59 | 58 | 58 | 52 | 57 | 58 | 58 | 57 | 57 | 57 | 61 | 58 | 57 | 58 | 57 | 59 | 57 | 56 |
| Malaysia | 57 (1.1) | 57 | 58 | 57 | 57 | 57 | 57 | 58 | 58 | 59 | 57 | 57 | 58 | 58 | 58 | 57 | 57 | 58 | 57 | 57 | 59 | 58 | 57 | 57 | 55 | 57 | 57 | 58 | 58 | 58 | 57 | 60 | 58 | 58 | 57 | 60 | 59 | 57 | 57 |
| Czech Republic | 57 (1.1) | 57 | 57 | 57 | 57 | 57 | 57 | 57 | 57 | 60 | 56 | 57 | 58 | 58 | 58 | 57 | 57 | 58 | 57 | 57 | 57 | 58 | 57 | 58 | 54 | 57 | 57 | 58 | 58 | 57 | 57 | 60 | 58 | 58 | 57 | 56 | 57 | 57 | 56 |
| Finland | 56 (0.7) | 55 | 56 | 55 | 56 | 55 | 56 | 56 | 56 | 59 | 55 | 56 | 55 | 57 | 56 | 55 | 56 | 55 | 56 | 56 | 59 | 58 | 56 | 56 | 50 | 56 | 56 | 56 | 56 | 55 | 55 | 60 | 56 | 56 | 56 | 56 | 57 | 56 | 54 |
| Bulgaria | 54 (1.4) | 54 | 55 | 54 | 54 | 54 | 54 | 55 | 54 | 56 | 54 | 54 | 55 | 55 | 55 | 54 | 54 | 55 | 54 | 54 | 54 | 54 | 54 | 55 | 53 | 54 | 54 | 55 | 55 | 55 | 54 | 57 | 55 | 55 | 54 | 54 | 55 | 54 | 54 |
| Latvia (LSS) | 53 (0.8) | 52 | 53 | 52 | 53 | 53 | 53 | 54 | 53 | 55 | 53 | 53 | 54 | 53 | 54 | 53 | 53 | 54 | 53 | 53 | 53 | 53 | 53 | 53 | 52 | 53 | 53 | 53 | 54 | 53 | 53 | 56 | 53 | 53 | 53 | 52 | 54 | 53 | 52 |
| United States | 52 (0.9) | 52 | 53 | 52 | 52 | 53 | 52 | 53 | 53 | 55 | 52 | 53 | 54 | 53 | 54 | 53 | 52 | 52 | 53 | 52 | 54 | 54 | 53 | 53 | 48 | 52 | 52 | 52 | 54 | 53 | 51 | 56 | 52 | 52 | 53 | 52 | 54 | 53 | 51 |
| England | 50 (1.1) | 49 | 50 | 49 | 49 | 49 | 50 | 49 | 50 | 54 | 50 | 51 | 49 | 51 | 49 | 50 | 50 | 49 | 50 | 50 | 53 | 52 | 50 | 50 | 43 | 49 | 50 | 50 | 50 | 50 | 49 | 53 | 51 | 49 | 50 | 50 | 50 | 50 | 48 |
| New Zealand | 50 (1.2) | 49 | 50 | 49 | 49 | 50 | 50 | 50 | 49 | 53 | 49 | 50 | 49 | 51 | 50 | 49 | 50 | 50 | 49 | 49 | 52 | 51 | 49 | 50 | 44 | 49 | 49 | 50 | 49 | 49 | 49 | 53 | 50 | 49 | 49 | 49 | 51 | 49 | 48 |
| Lithuania | 49 (1.0) | 49 | 50 | 49 | 49 | 49 | 49 | 50 | 49 | 51 | 49 | 49 | 51 | 50 | 50 | 49 | 50 | 49 | 49 | 49 | 52 | 50 | 49 | 49 | 47 | 49 | 49 | 50 | 51 | 50 | 49 | 53 | 50 | 50 | 49 | 50 | 50 | 49 | 48 |
| Italy | 49 (0.8) | 48 | 49 | 48 | 49 | 49 | 49 | 48 | 48 | 48 | 49 | 49 | 48 | 50 | 48 | 49 | 49 | 49 | 49 | 49 | 50 | 49 | 49 | 49 | 45 | 49 | 49 | 49 | 48 | 47 | 48 | 51 | 48 | 47 | 48 | 47 | 49 | 47 | 48 |
| Cyprus | 47 (0.4) | 46 | 47 | 46 | 46 | 47 | 47 | 48 | 47 | 49 | 47 | 47 | 47 | 48 | 48 | 47 | 46 | 47 | 47 | 47 | 49 | 48 | 47 | 47 | 45 | 47 | 47 | 49 | 48 | 47 | 47 | 50 | 48 | 48 | 47 | 48 | 49 | 47 | 46 |
| Romania | 47 (1.2) | 46 | 47 | 46 | 47 | 47 | 47 | 48 | 47 | 48 | 47 | 47 | 48 | 47 | 48 | 47 | 47 | 49 | 47 | 47 | 49 | 47 | 47 | 47 | 45 | 47 | 47 | 48 | 48 | 47 | 47 | 49 | 48 | 48 | 47 | 47 | 48 | 47 | 47 |
| Moldova | 46 (0.9) | 46 | 46 | 46 | 46 | 46 | 46 | 45 | 46 | 47 | 46 | 45 | 47 | 46 | 46 | 46 | 45 | 46 | 46 | 46 | 47 | 46 | 46 | 46 | 42 | 46 | 45 | 45 | 46 | 46 | 45 | 48 | 46 | 46 | 46 | 46 | 46 | 45 | 46 |
| Israel | 45 (1.1) | 44 | 46 | 45 | 45 | 45 | 45 | 45 | 45 | 47 | 45 | 45 | 46 | 46 | 46 | 45 | 45 | 45 | 45 | 45 | 46 | 46 | 45 | 45 | 38 | 45 | 46 | 46 | 46 | 44 | 40 | 48 | 45 | 46 | 45 | 45 | 46 | 45 | 44 |
| Thailand | 45 (1.1) | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 45 | 47 | 45 | 45 | 46 | 46 | 46 | 45 | 45 | 45 | 45 | 45 | 46 | 46 | 45 | 45 | 41 | 45 | 45 | 46 | 46 | 41 | 40 | 48 | 45 | 46 | 45 | 45 | 46 | 45 | 44 |
| Macedonia, Rep. of | 40 (0.8) | 40 | 40 | 39 | 40 | 40 | 40 | 41 | 40 | 42 | 40 | 40 | 40 | 40 | 40 | 40 | 39 | 41 | 40 | 40 | 40 | 41 | 40 | 40 | 38 | 40 | 40 | 40 | 41 | 41 | 40 | 43 | 41 | 40 | 40 | 39 | 40 | 40 | 39 |
| Tunisia | 39 (0.5) | 40 | 40 | 39 | 39 | 39 | 38 | 41 | 39 | 39 | 39 | 40 | 41 | 40 | 40 | 39 | 39 | 41 | 39 | 39 | 40 | 40 | 39 | 40 | 39 | 40 | 39 | 40 | 41 | 41 | 40 | 42 | 40 | 39 | 39 | 41 | 41 | 40 | 40 |
| Jordan | 38 (0.6) | 38 | 38 | 38 | 39 | 39 | 38 | 38 | 38 | 38 | 38 | 38 | 38 | 38 | 37 | 38 | 38 | 38 | 38 | 38 | 39 | 38 | 37 | 38 | 35 | 38 | 38 | 38 | 39 | 39 | 38 | 41 | 39 | 39 | 38 | 37 | 39 | 38 | 37 |
| Turkey | 37 (0.9) | 36 | 38 | 36 | 37 | 37 | 37 | 37 | 37 | 38 | 37 | 37 | 38 | 38 | 37 | 37 | 36 | 38 | 37 | 37 | 38 | 37 | 37 | 38 | 37 | 38 | 37 | 37 | 39 | 37 | 38 | 40 | 38 | 36 | 37 | 37 | 38 | 38 | 36 |
| Iran, Islamic Rep. | 36 (0.6) | 36 | 37 | 35 | 36 | 36 | 36 | 37 | 36 | 37 | 36 | 36 | 36 | 36 | 37 | 36 | 35 | 37 | 36 | 36 | 37 | 36 | 36 | 36 | 34 | 36 | 36 | 37 | 37 | 36 | 36 | 39 | 36 | 36 | 36 | 37 | 37 | 35 | 35 |
| Indonesia | 35 (0.8) | 34 | 35 | 34 | 35 | 35 | 35 | 35 | 35 | 36 | 35 | 35 | 35 | 35 | 35 | 35 | 34 | 35 | 35 | 35 | 35 | 35 | 35 | 35 | 32 | 35 | 35 | 35 | 35 | 35 | 34 | 37 | 35 | 35 | 35 | 34 | 35 | 35 | 34 |
| Chile | 31 (0.8) | 31 | 31 | 30 | 31 | 31 | 31 | 31 | 31 | 33 | 31 | 31 | 31 | 32 | 31 | 31 | 31 | 31 | 31 | 31 | 32 | 32 | 31 | 31 | 28 | 31 | 31 | 31 | 31 | 31 | 31 | 34 | 32 | 31 | 31 | 32 | 32 | 31 | 30 |
| Philippines | 26 (0.7) | 26 | 27 | 26 | 26 | 26 | 26 | 27 | 27 | 28 | 26 | 27 | 26 | 26 | 26 | 26 | 26 | 27 | 26 | 26 | 26 | 27 | 26 | 27 | 27 | 26 | 26 | 27 | 27 | 27 | 26 | 29 | 27 | 27 | 26 | 27 | 27 | 26 | 26 |
| Morocco | 24 (0.2) | 23 | 23 | 23 | 24 | 23 | 24 | 23 | 24 | 25 | 23 | 24 | 24 | 24 | 24 | 23 | 23 | 24 | 24 | 24 | 24 | 24 | 24 | 24 | 22 | 24 | 24 | 23 | 24 | 24 | 23 | 26 | 24 | 24 | 24 | 24 | 24 | 23 | 23 |
| South Africa | 21 (0.6) | 20 | 21 | 20 | 21 | 20 | 21 | 21 | 21 | 21 | 20 | 21 | 21 | 21 | 21 | 20 | 20 | 21 | 21 | 21 | 21 | 21 | 21 | 21 | 20 | 20 | 21 | 21 | 21 | 20 | 20 | 23 | 21 | 21 | 21 | 21 | 21 | 20 | 20 |
| **International Avg.** | **51 (0.2)** | **50** | **51** | **50** | **51** | **51** | **51** | **51** | **51** | **53** | **51** | **51** | **51** | **51** | **51** | **51** | **51** | **51** | **51** | **52** | **52** | **52** | **51** | **51** | **48** | **51** | **51** | **51** | **51** | **51** | **50** | **54** | **51** | **51** | **51** | **51** | **52** | **51** | **50** |
| Number of Items (Score Points) Identified* | 169 | 159 | 127 | 152 | 169 | 160 | 166 | 142 | 169 | 120 | 162 | 165 | 126 | 159 | 141 | 162 | 140 | 135 | 169 | 169 | 136 | 156 | 169 | 161 | 99 | 153 | 169 | 128 | 147 | 134 | 145 | 149 | 165 | 150 | 169 | 98 | 126 | 150 | 141 |

*Of the 162 items in the mathematics test, some items had two parts and some extended–response items were scored on a two-point scale, resulting in 169 total score points.

( ) Standard errors for the average percent of correct responses on all items appear in parentheses. Because results are rounded to the nearest whole number, some totals may appear inconsistent.

Because population coverage falls below 65% Latvia is annotated LSS for Latvian Speaking Schools only.

SOURCE: IEA Third International Mathematics and Science Study (TIMSS), 1998-1999.

# Exhibit C.2  Standard Errors for the Test-Curriculum Matching Analysis – Mathematics

TIMSS 1999
8th grade
Mathematics

| Country | Average Percent Correct on All Items* |
|---|---|
| Singapore | 77 (1.3) |
| Korea, Rep. of | 73 (0.4) |
| Hong Kong, SAR | 72 (1.1) |
| Chinese Taipei | 72 (0.9) |
| Japan | 72 (0.4) |
| Belgium (Flemish) | 66 (0.8) |
| Hungary | 61 (0.9) |
| Slovak Republic | 61 (1.1) |
| Netherlands | 61 (1.8) |
| Slovenia | 60 (0.7) |
| Canada | 59 (0.5) |
| Russian Federation | 58 (1.3) |
| Australia | 58 (1.2) |
| Malaysia | 57 (1.1) |
| Czech Republic | 57 (1.1) |
| Finland | 56 (0.7) |
| Bulgaria | 54 (1.4) |
| Latvia (LSS) | 53 (0.8) |
| United States | 52 (0.9) |
| England | 50 (1.1) |
| New Zealand | 50 (1.2) |
| Lithuania | 49 (1.0) |
| Italy | 49 (0.8) |
| Cyprus | 47 (0.4) |
| Romania | 47 (1.2) |
| Moldova | 46 (0.9) |
| Israel | 45 (0.8) |
| Thailand | 45 (1.1) |
| Macedonia, Rep. of | 40 (0.8) |
| Tunisia | 39 (0.5) |
| Jordan | 38 (0.6) |
| Turkey | 37 (0.9) |
| Iran, Islamic Rep. | 36 (0.6) |
| Indonesia | 35 (0.8) |
| Chile | 31 (0.8) |
| Philippines | 26 (0.7) |
| Morocco | 24 (0.2) |
| South Africa | 21 (0.6) |
| **International Avg.** | **51 (0.2)** |

Number of Items (Score Points) Identified*: 169

Column country headers (rotated, left to right): Singapore, Korea Rep. of, Hong Kong SAR, Chinese Taipei, Japan, Belgium (Flemish), Hungary, Slovak Republic, Netherlands, Slovenia, Canada, Russian Federation, Australia, Malaysia, Czech Republic, Finland, Bulgaria, Latvia (LSS), United States, England, New Zealand, Lithuania, Italy, Cyprus, Romania, Moldova, Israel, Thailand, Macedonia Rep. of, Tunisia, Jordan, Turkey, Iran Islamic Rep., Indonesia, Chile, Philippines, Morocco, South Africa.

Number of Items per column: 159, 127, 152, 169, 160, 166, 142, 169, 120, 162, 165, 126, 159, 141, 162, 140, 135, 169, 169, 136, 156, 169, 161, 99, 153, 169, 128, 147, 134, 145, 149, 165, 150, 169, 98, 126, 150, 141

*Of the 162 items in the mathematics test, some items had two parts and some extended-response items were scored on a two-point scale, resulting in 169 total score points.

( ) Standard errors for the average percent of correct responses appear in parentheses. The matrix contains standard errors corresponding to the average percent of correct responses based on TCMA subsets of items, as displayed in Exhibit C.1. Because results are rounded to the nearest whole number, some totals may appear inconsistent.

Because population coverage falls below 65% Latvia is annotated LSS for Latvian Speaking Schools only.

SOURCE: IEA Third International Mathematics and Science Study (TIMSS), 1998-1999.