# Appendix B

## THE TEST–CURRICULUM MATCHING ANALYSIS

When comparing student achievement across countries, it is important that the comparisons be as "fair" as possible. TIMSS has worked towards 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. Developing the TIMSS tests involved the interaction of experts in the sciences with representatives of the participating countries and testing specialists.[1] The National Research Coordinators (NRCs) from each country formally approved the TIMSS test, thus accepting it as being sufficiently fair to compare their students' science achievement with that of students from other countries.

Although the TIMSS test was developed to represent a set of agreed-upon science content areas, there are differences among the curricula of participating countries that result in various science topics being taught at different grades. To restrict test items not only to those topics in the curricula of all countries but also to those covered in the same sequence in all participating countries would severely limit test coverage and restrict the research questions about international differences that TIMSS is designed to address. The TIMSS tests, therefore, inevitably contain some items measuring topics unfamiliar to some students in some countries.

The Test-Curriculum Matching Analysis (TCMA) was developed and conducted to investigate the appropriateness of the TIMSS science test for seventh- and eighth-grade students in the participating countries, and to show how student performance for individual countries varied when based only on the test questions that were judged to be relevant to their own curriculum.[2]

To gather data about the extent to which the TIMSS tests were relevant to the curriculum of the participating countries, TIMSS asked the NRC of each country to report whether or not each item was in their country's intended curriculum at each of the two grades being tested. The NRC was asked to choose a person or persons who were very familiar with the curricula at the grades being tested to make the 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% of the students. The NRCs had considerable flexibility in selecting items and may have considered items inappropriate for other reasons. All participating countries except Thailand returned the information for analysis.

Tables B.1 and B.2 present the TCMA results for the eighth and seventh grades, respectively. The first row of each table indicates that at both grades the countries varied substantially in the number of items considered appropriate. At the eighth

---

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

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

grade, more than half of the countries indicated that items representing three-quarters or more of the score points (110 out of a possible 146) were appropriate,[3] with the percent ranging from 100% in Spain, Iceland, and the United States to approximately 40% in Korea (59 score points) and French-speaking Belgium (58 score points). Fewer items were selected at the seventh grade, but nearly half of the countries selected at least 60%, with several selecting at least three-quarters of the score points. All items were selected at the seventh grade as well as the eighth grade in both the United States and Iceland. At the seventh grade there were also several countries, including Korea and Japan, which retained about 30% or less. That lower percentages of items were selected for the TCMA at the seventh grade is consistent with the instrument-development process, which put more emphasis on the upper-grade curriculum.

Since most countries indicated that some items were not included in their intended curricula at the two grades tested, the question becomes whether the inclusion of these items had any effect on the international performance comparisons.[4] The TCMA results provide a method for answering this question, providing evidence that it is reasonable to make cross-national comparisons on the basis of the TIMSS science test.

Each of the first columns in Tables B.1 and B.2 shows the overall average percent correct for each country (as discussed in Chapter 2 and reproduced here for convenience in making comparisons). The countries are presented in the order of their overall performance, from highest to lowest. To interpret these tables, reading across a row provides the average percent correct for the students in the country identified by that row on the items selected by each of the countries named across the top of the table. For example, at the eighth grade, Singapore, where the average percent correct was 72% on its own set of items, also had 72% for the items selected by Korea, 73% for those selected by Japan, 69% for those selected by the Czech Republic, and so forth. The column for a country shows how each of the other countries performed on the subset of items selected for its own students. Using the set of items selected by Hong Kong as an example, on average, 71% of these items were answered correctly by the Singaporean students, 65% by the Korean students, 66% by the Japanese, and so forth. The shaded diagonal elements in each table show how each country performed on the subset of items that it selected based on its own curriculum. Thus, the Hong Kong students themselves averaged 59% correct responses on the items identified by Hong Kong for the analysis.

---

[3] Of the 135 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 and others on a three-point scale. The total number of score points available for analysis was 146. The TCMA uses the score points in order to give the same weight to items that they received in the test scoring.

[4] It should be noted that the performance levels presented in Tables B.1 and B.2 are based on average percents correct as was done in Chapter 2, which is different from the average scale scores that were presented in Chapter 1. The cost and delay of scaling would have been prohibitive for the TCMA analyses.

**Table B.1** Test-Curriculum Matching Analysis Results - Science - Upper Grade (Eighth Grade)*

Average Percent Correct Based on Subsets of Items Specially Identified by Each Country as Addressing Its Curriculum (See Table B.3 for corresponding standard errors)

**Instructions:** Read *across* the row to compare how a 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.

| Country | Average Percent Correct on All Items 146** | Singapore | Korea | Japan | Czech Republic | Netherlands | Bulgaria | Slovenia | Austria | England | Hungary | Belgium (Fl) | Australia | Slovak Republic | Sweden | Canada | Ireland | United States | Russian Federation | Germany | New Zealand | Norway | Hong Kong | Israel | Switzerland | Spain | Scotland | France | Iceland | Greece | Denmark | Belgium (Fr) | Latvia (LSS) | Portugal | Romania | Lithuania | Iran, Islamic Rep. | Cyprus | Kuwait | Colombia | South Africa |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *(Number of Score Points Included)* | | 109 | 59 | 86 | 136 | 102 | 112 | 140 | 131 | 124 | 129 | 98 | 133 | 129 | 125 | 121 | 90 | 146 | 96 | 129 | 126 | 111 | 68 | 102 | 105 | 146 | 97 | 73 | 146 | 111 | 70 | 58 | 113 | 133 | 99 | 120 | 87 | 78 | 131 | 112 | 74 |
| **Singapore** | 70 (1.0) | 72 | 72 | 73 | 69 | 72 | 70 | 69 | 70 | 71 | 69 | 69 | 70 | 70 | 71 | 71 | 72 | 70 | 70 | 70 | 72 | 70 | 71 | 71 | 72 | 70 | 71 | 71 | 70 | 71 | 72 | 70 | 68 | 70 | 71 | 71 | 69 | 72 | 70 | 72 | 73 |
| **Korea** | 66 (0.3) | 66 | 67 | 68 | 66 | 67 | 66 | 66 | 66 | 66 | 65 | 67 | 65 | 66 | 66 | 66 | 67 | 65 | 68 | 66 | 66 | 67 | 65 | 66 | 68 | 66 | 67 | 67 | 66 | 65 | 67 | 67 | 65 | 65 | 65 | 68 | 65 | 66 | 66 | 67 | 65 |
| **Japan** | 65 (0.3) | 66 | 68 | 67 | 65 | 67 | 66 | 66 | 66 | 66 | 65 | 64 | 65 | 66 | 66 | 66 | 67 | 65 | 68 | 64 | 67 | 66 | 65 | 66 | 67 | 64 | 67 | 67 | 65 | 65 | 67 | 68 | 65 | 65 | 64 | 68 | 64 | 66 | 66 | 65 | 65 |
| **Czech Republic** | 64 (0.8) | 65 | 70 | 67 | 64 | 66 | 65 | 64 | 65 | 63 | 62 | 64 | 64 | 66 | 66 | 63 | 67 | 64 | 67 | 63 | 64 | 65 | 64 | 62 | 67 | 64 | 63 | 65 | 64 | 65 | 67 | 68 | 62 | 64 | 63 | 66 | 64 | 66 | 64 | 65 | 65 |
| *Netherlands* | 62 (1.0) | 63 | 66 | 66 | 63 | 66 | 63 | 62 | 63 | 61 | 62 | 64 | 62 | 63 | 64 | 62 | 64 | 62 | 64 | 63 | 62 | 64 | 60 | 60 | 66 | 62 | 63 | 64 | 62 | 63 | 66 | 67 | 61 | 62 | 62 | 64 | 61 | 65 | 62 | 63 | 61 |
| *Bulgaria* | 62 (1.0) | 64 | 65 | 65 | 63 | 64 | 64 | 62 | 63 | 61 | 62 | 61 | 61 | 63 | 63 | 61 | 63 | 62 | 66 | 62 | 63 | 63 | 62 | 62 | 65 | 62 | 62 | 64 | 62 | 62 | 64 | 66 | 60 | 62 | 60 | 63 | 61 | 64 | 61 | 63 | 64 |
| *Slovenia* | 62 (0.5) | 62 | 65 | 64 | 62 | 64 | 62 | 62 | 63 | 61 | 61 | 63 | 61 | 63 | 63 | 61 | 63 | 62 | 64 | 61 | 63 | 63 | 62 | 60 | 64 | 61 | 62 | 60 | 60 | 60 | 63 | 64 | 59 | 60 | 61 | 62 | 61 | 63 | 60 | 61 | 63 |
| *Austria* | 61 (0.7) | 63 | 65 | 65 | 61 | 63 | 63 | 62 | 63 | 61 | 62 | 63 | 61 | 63 | 63 | 61 | 64 | 61 | 64 | 62 | 63 | 63 | 63 | 60 | 64 | 61 | 62 | 60 | 60 | 61 | 63 | 64 | 59 | 61 | 60 | 62 | 62 | 63 | 61 | 63 | 63 |
| **England** | 61 (0.6) | 62 | 62 | 63 | 62 | 64 | 62 | 62 | 62 | 62 | 62 | 63 | 61 | 62 | 63 | 61 | 63 | 61 | 62 | 61 | 63 | 62 | 63 | 61 | 64 | 61 | 62 | 61 | 59 | 60 | 64 | 63 | 59 | 61 | 61 | 63 | 61 | 64 | 60 | 63 | 61 |
| *Hungary* | 61 (0.6) | 61 | 63 | 64 | 60 | 62 | 61 | 61 | 61 | 60 | 62 | 61 | 61 | 62 | 62 | 61 | 61 | 61 | 64 | 61 | 61 | 62 | 60 | 61 | 63 | 61 | 59 | 62 | 60 | 60 | 61 | 63 | 59 | 61 | 60 | 63 | 61 | 64 | 60 | 62 | 61 |
| **Belgium (Fl)** | 60 (1.1) | 61 | 62 | 64 | 61 | 64 | 60 | 60 | 61 | 60 | 61 | 62 | 60 | 60 | 62 | 60 | 63 | 60 | 62 | 61 | 61 | 63 | 63 | 60 | 64 | 60 | 61 | 62 | 60 | 61 | 65 | 65 | 59 | 60 | 60 | 62 | 59 | 63 | 60 | 62 | 59 |
| *Australia* | 60 (0.7) | 61 | 61 | 62 | 60 | 62 | 60 | 60 | 61 | 60 | 61 | 60 | 61 | 61 | 61 | 60 | 61 | 60 | 62 | 60 | 61 | 61 | 60 | 60 | 62 | 60 | 60 | 60 | 59 | 59 | 61 | 64 | 59 | 59 | 59 | 62 | 59 | 63 | 59 | 61 | 59 |
| **Slovak Republic** | 59 (0.6) | 60 | 63 | 63 | 60 | 61 | 60 | 60 | 61 | 58 | 60 | 59 | 59 | 61 | 60 | 59 | 62 | 59 | 63 | 60 | 60 | 60 | 61 | 59 | 62 | 59 | 60 | 60 | 59 | 60 | 61 | 62 | 58 | 59 | 59 | 62 | 60 | 61 | 59 | 60 | 61 |
| *Sweden* | 59 (0.6) | 58 | 60 | 63 | 58 | 61 | 59 | 59 | 59 | 58 | 59 | 58 | 59 | 59 | 61 | 60 | 60 | 59 | 61 | 59 | 60 | 59 | 60 | 60 | 63 | 59 | 60 | 58 | 59 | 59 | 62 | 63 | 57 | 58 | 58 | 61 | 58 | 61 | 59 | 60 | 59 |
| *Canada* | 59 (0.5) | 59 | 58 | 61 | 60 | 61 | 58 | 57 | 58 | 58 | 58 | 59 | 59 | 59 | 60 | 61 | 60 | 59 | 60 | 59 | 61 | 60 | 60 | 59 | 62 | 58 | 58 | 59 | 58 | 60 | 61 | 62 | 56 | 58 | 58 | 61 | 57 | 61 | 58 | 59 | 59 |
| **Ireland** | 58 (0.9) | 59 | 60 | 61 | 59 | 60 | 57 | 58 | 59 | 59 | 58 | 58 | 58 | 58 | 59 | 59 | 62 | 58 | 59 | 58 | 60 | 59 | 60 | 59 | 61 | 58 | 59 | 57 | 58 | 58 | 60 | 60 | 57 | 58 | 58 | 60 | 57 | 61 | 58 | 60 | 57 |
| **United States** | 58 (1.0) | 59 | 58 | 60 | 58 | 60 | 58 | 58 | 59 | 59 | 59 | 58 | 57 | 58 | 60 | 60 | 61 | 58 | 61 | 58 | 60 | 58 | 59 | 58 | 61 | 58 | 58 | 57 | 58 | 58 | 62 | 63 | 57 | 58 | 57 | 60 | 56 | 61 | 58 | 60 | 58 |
| **Russian Federation** | 58 (0.8) | 59 | 61 | 62 | 59 | 61 | 59 | 58 | 59 | 58 | 59 | 58 | 58 | 60 | 60 | 60 | 61 | 59 | 65 | 59 | 60 | 60 | 60 | 58 | 61 | 58 | 57 | 60 | 58 | 59 | 62 | 63 | 57 | 58 | 57 | 59 | 56 | 59 | 58 | 60 | 61 |
| *Germany* | 58 (0.8) | 58 | 60 | 62 | 58 | 62 | 58 | 58 | 59 | 58 | 58 | 58 | 58 | 59 | 59 | 59 | 61 | 58 | 61 | 59 | 58 | 60 | 59 | 58 | 61 | 58 | 57 | 60 | 57 | 57 | 61 | 63 | 57 | 57 | 57 | 60 | 57 | 60 | 57 | 59 | 58 |
| **New Zealand** | 58 (0.8) | 59 | 60 | 62 | 58 | 60 | 58 | 58 | 59 | 58 | 58 | 58 | 58 | 58 | 59 | 59 | 60 | 60 | 60 | 58 | 60 | 59 | 59 | 58 | 60 | 58 | 57 | 58 | 57 | 57 | 60 | 60 | 56 | 56 | 58 | 59 | 56 | 60 | 58 | 59 | 57 |
| **Norway** | 58 (0.4) | 59 | 60 | 62 | 59 | 62 | 57 | 58 | 59 | 58 | 59 | 58 | 58 | 60 | 59 | 59 | 60 | 58 | 59 | 58 | 60 | 60 | 60 | 56 | 62 | 58 | 59 | 57 | 55 | 55 | 60 | 61 | 57 | 56 | 57 | 59 | 56 | 58 | 57 | 58 | 60 |
| **Hong Kong** | 58 (1.0) | 59 | 58 | 61 | 58 | 60 | 59 | 57 | 58 | 57 | 57 | 57 | 57 | 59 | 59 | 58 | 61 | 58 | 61 | 58 | 60 | 58 | 59 | 58 | 61 | 57 | 58 | 60 | 58 | 59 | 59 | 62 | 56 | 55 | 56 | 59 | 57 | 57 | 57 | 60 | 60 |
| *Israel* | 57 (1.1) | 57 | 57 | 60 | 57 | 59 | 57 | 57 | 58 | 56 | 58 | 57 | 57 | 58 | 58 | 57 | 59 | 58 | 58 | 57 | 58 | 58 | 58 | 58 | 59 | 56 | 57 | 57 | 58 | 58 | 60 | 58 | 55 | 55 | 55 | 59 | 56 | 61 | 57 | 60 | 56 |
| **Switzerland** | 56 (0.5) | 57 | 56 | 60 | 56 | 58 | 56 | 56 | 57 | 56 | 57 | 57 | 56 | 58 | 58 | 56 | 59 | 57 | 58 | 57 | 58 | 58 | 58 | 56 | 60 | 56 | 56 | 55 | 56 | 56 | 59 | 60 | 55 | 55 | 55 | 58 | 56 | 61 | 56 | 58 | 56 |
| **Spain** | 56 (0.4) | 56 | 57 | 58 | 57 | 57 | 56 | 57 | 57 | 55 | 56 | 56 | 55 | 57 | 56 | 57 | 57 | 56 | 57 | 56 | 57 | 56 | 56 | 56 | 57 | 56 | 55 | 57 | 56 | 57 | 57 | 58 | 53 | 53 | 55 | 56 | 55 | 60 | 55 | 57 | 57 |
| *Scotland* | 55 (1.0) | 56 | 55 | 58 | 56 | 58 | 56 | 55 | 56 | 56 | 56 | 56 | 55 | 56 | 56 | 57 | 57 | 55 | 56 | 55 | 58 | 57 | 57 | 56 | 58 | 55 | 56 | 55 | 55 | 55 | 57 | 57 | 53 | 55 | 54 | 57 | 55 | 59 | 55 | 57 | 55 |
| *France* | 54 (0.6) | 54 | 54 | 57 | 54 | 56 | 53 | 54 | 54 | 54 | 55 | 54 | 53 | 55 | 53 | 56 | 55 | 54 | 56 | 54 | 55 | 55 | 54 | 55 | 57 | 54 | 54 | 57 | 54 | 54 | 55 | 59 | 53 | 53 | 53 | 57 | 52 | 56 | 53 | 54 | 52 |
| **Iceland** | 52 (0.9) | 51 | 51 | 55 | 53 | 55 | 53 | 52 | 53 | 43 | 52 | 53 | 52 | 53 | 53 | 53 | 53 | 52 | 54 | 53 | 53 | 55 | 53 | 53 | 55 | 52 | 52 | 52 | 52 | 52 | 55 | 54 | 50 | 51 | 52 | 54 | 51 | 55 | 54 | 53 | 51 |
| *Greece* | 52 (0.5) | 52 | 53 | 52 | 50 | 52 | 53 | 50 | 53 | 50 | 50 | 49 | 50 | 53 | 50 | 53 | 53 | 52 | 54 | 50 | 53 | 53 | 53 | 53 | 54 | 52 | 52 | 52 | 50 | 54 | 52 | 53 | 50 | 51 | 49 | 52 | 50 | 54 | 50 | 53 | 51 |
| *Denmark* | 51 (0.6) | 51 | 52 | 52 | 50 | 51 | 51 | 50 | 51 | 50 | 50 | 51 | 50 | 52 | 51 | 51 | 51 | 52 | 54 | 50 | 51 | 53 | 51 | 52 | 53 | 50 | 52 | 52 | 50 | 51 | 55 | 52 | 49 | 50 | 49 | 51 | 50 | 53 | 50 | 51 | 50 |
| **Belgium (Fr)** | 50 (0.7) | 50 | 50 | 54 | 50 | 53 | 51 | 50 | 51 | 49 | 50 | 49 | 50 | 51 | 50 | 48 | 51 | 50 | 53 | 51 | 52 | 52 | 51 | 51 | 54 | 50 | 50 | 52 | 50 | 50 | 52 | 56 | 49 | 49 | 50 | 53 | 49 | 53 | 50 | 51 | 49 |
| **Latvia (LSS)** | 50 (0.6) | 51 | 52 | 55 | 50 | 52 | 51 | 50 | 51 | 51 | 51 | 51 | 50 | 52 | 51 | 52 | 52 | 50 | 54 | 50 | 51 | 51 | 52 | 52 | 53 | 50 | 50 | 52 | 50 | 50 | 51 | 54 | 49 | 49 | 51 | 52 | 50 | 52 | 50 | 51 | 51 |
| **Portugal** | 50 (0.6) | 51 | 53 | 52 | 49 | 51 | 50 | 50 | 50 | 50 | 50 | 49 | 50 | 51 | 50 | 50 | 52 | 50 | 52 | 50 | 51 | 51 | 50 | 50 | 52 | 50 | 49 | 52 | 50 | 51 | 52 | 54 | 48 | 49 | 49 | 52 | 50 | 54 | 49 | 52 | 49 |
| **Romania** | 50 (0.8) | 51 | 52 | 52 | 50 | 51 | 50 | 50 | 50 | 49 | 50 | 51 | 50 | 52 | 50 | 51 | 52 | 50 | 54 | 50 | 50 | 51 | 49 | 50 | 51 | 50 | 50 | 52 | 50 | 49 | 52 | 52 | 47 | 49 | 49 | 51 | 50 | 54 | 49 | 52 | 51 |
| **Lithuania** | 49 (0.7) | 50 | 51 | 53 | 49 | 51 | 50 | 49 | 50 | 48 | 50 | 49 | 49 | 51 | 50 | 50 | 51 | 49 | 55 | 49 | 49 | 50 | 50 | 50 | 51 | 49 | 49 | 51 | 49 | 49 | 50 | 53 | 47 | 47 | 50 | 52 | 49 | 50 | 49 | 50 | 51 |
| *Iran, Islamic Rep.* | 47 (0.6) | 49 | 51 | 50 | 48 | 50 | 49 | 47 | 49 | 48 | 48 | 48 | 48 | 49 | 48 | 49 | 49 | 47 | 51 | 48 | 48 | 49 | 49 | 47 | 49 | 47 | 47 | 48 | 47 | 47 | 47 | 47 | 47 | 48 | 47 | 49 | 50 | 52 | 48 | 49 | 49 |
| *Cyprus* | 47 (0.4) | 48 | 47 | 50 | 47 | 48 | 47 | 47 | 48 | 47 | 47 | 47 | 46 | 48 | 43 | 47 | 50 | 47 | 49 | 47 | 48 | 48 | 48 | 48 | 50 | 47 | 47 | 48 | 47 | 48 | 48 | 50 | 45 | 47 | 47 | 49 | 47 | 51 | 43 | 49 | 48 |
| **Kuwait** | 43 (0.9) | 43 | 43 | 43 | 43 | 44 | 44 | 43 | 44 | 43 | 44 | 40 | 42 | 43 | 43 | 44 | 43 | 44 | 44 | 42 | 44 | 43 | 43 | 42 | 44 | 44 | 42 | 44 | 43 | 44 | 48 | 43 | 42 | 44 | 44 | 45 | 44 | 47 | 43 | 45 | 43 |
| **Colombia** | 39 (0.8) | 39 | 40 | 41 | 39 | 40 | 39 | 39 | 40 | 39 | 39 | 40 | 39 | 39 | 41 | 40 | 41 | 39 | 40 | 38 | 40 | 40 | 39 | 40 | 41 | 39 | 39 | 39 | 39 | 40 | 42 | 41 | 36 | 39 | 38 | 40 | 38 | 41 | 38 | 41 | 39 |
| **South Africa** | 27 (1.3) | 28 | 27 | 29 | 27 | 28 | 27 | 27 | 28 | 27 | 26 | 29 | 26 | 27 | 27 | 27 | 29 | 27 | 30 | 27 | 27 | 27 | 28 | 27 | 28 | 27 | 26 | 28 | 27 | 27 | 28 | 29 | 26 | 27 | 27 | 28 | 27 | 28 | 26 | 28 | 28 |
| **International Average** | 55 (0.7) | 56 | 56 | 59 | 56 | 58 | 56 | 56 | 56 | 56 | 56 | 56 | 55 | 57 | 57 | 57 | 57 | 55 | 58 | 56 | 57 | 57 | 57 | 56 | 58 | 55 | 56 | 57 | 55 | 56 | 58 | 58 | 54 | 55 | 56 | 57 | 55 | 58 | 55 | 57 | 56 |

*Eighth grade in most countries; see Table 2 for more information about the grades tested in each country.

**Of the 135 items in the science test, some items had two parts and some extended-response items were scored on a two- or three-point scale, resulting in 146 total score points.

( ) Standard errors for the average percent of correct responses on all items appear in parentheses. Standard errors for scores based on subsets of items are provided in Table B.3.

Because results are rounded to the nearest whole number, some totals may appear inconsistent.

Countries shown in italics did not satisfy one or more guidelines for sample participation rates, age/grade specifications, or classroom sampling procedures (see Figure A.3 for details).

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

SOURCE: IEA Third International Mathematics and Science Study (TIMSS), 1994-95.

## Table B.2 Test-Curriculum Matching Analysis Results - Science - Lower Grade (Seventh Grade)*

Average Percent Correct Based on Subsets of Items Specially Identified by Each Country as Addressing Its Curriculum (See Table B.4 for corresponding standard errors)

**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.

| Country | Average Percent Correct on All Items 146** | Singapore 83 | Korea 43 | Japan 45 | Czech Republic 108 | Slovenia 132 | Belgium (Fl) 46 | Bulgaria 105 | Netherlands 34 | England 104 | Hungary 98 | Austria 52 | Slovak Republic 111 | United States 146 | Canada 78 | Australia 94 | Hong Kong 32 | Germany 88 | Ireland 60 | Sweden 86 | New Zealand 110 | Norway 92 | Switzerland 36 | Russian Federation 49 | Spain 132 | Scotland 49 | Iceland 146 | France 26 | Belgium (Fr) 23 | Romania 91 | Greece 72 | Denmark 20 | Iran, Islamic Rep. 48 | Latvia (LSS) 46 | Portugal 81 | Cyprus 29 | Colombia 79 | Lithuania 108 | South Africa 26 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Singapore | 61 (1.2) | 66 | 65 | 65 | 62 | 62 | 62 | 63 | 63 | 63 | 61 | 65 | 62 | 61 | 65 | 66 | 62 | 63 | 65 | 67 | 65 | 64 | 66 | 65 | 62 | 66 | 61 | 68 | 61 | 64 | 63 | 68 | 63 | 62 | 62 | 66 | 64 | 64 | 66 |
| Korea | 61 (0.4) | 65 | 64 | 67 | 63 | 62 | 63 | 63 | 63 | 63 | 61 | 67 | 63 | 61 | 63 | 63 | 60 | 63 | 63 | 67 | 62 | 63 | 64 | 67 | 62 | 63 | 61 | 69 | 63 | 62 | 62 | 61 | 61 | 59 | 63 | 63 | 65 | 62 | 68 |
| Japan | 59 (0.3) | 60 | 61 | 67 | 60 | 60 | 60 | 61 | 63 | 61 | 58 | 63 | 62 | 59 | 62 | 61 | 59 | 64 | 59 | 66 | 62 | 61 | 64 | 66 | 62 | 63 | 59 | 65 | 63 | 61 | 62 | 66 | 61 | 59 | 60 | 59 | 62 | 59 | 68 |
| Czech Republic | 58 (0.8) | 60 | 61 | 64 | 61 | 60 | 60 | 61 | 63 | 61 | 59 | 64 | 62 | 58 | 60 | 60 | 55 | 63 | 62 | 65 | 62 | 61 | 63 | 64 | 59 | 59 | 58 | 65 | 65 | 62 | 61 | 62 | 62 | 57 | 59 | 59 | 63 | 61 | 61 |
| *Slovenia* | 57 (0.5) | 59 | 60 | 61 | 58 | 58 | 59 | 58 | 60 | 57 | 58 | 61 | 60 | 57 | 58 | 59 | 55 | 61 | 62 | 62 | 59 | 60 | 63 | 62 | 57 | 59 | 57 | 65 | 65 | 59 | 59 | 60 | 61 | 57 | 58 | 58 | 61 | 59 | 61 |
| Belgium (Fl) | 57 (0.5) | 58 | 58 | 60 | 60 | 58 | 60 | 58 | 65 | 57 | 59 | 61 | 60 | 57 | 61 | 59 | 59 | 62 | 60 | 66 | 61 | 63 | 66 | 63 | 59 | 61 | 57 | 61 | 64 | 59 | 62 | 66 | 63 | 57 | 59 | 58 | 64 | 59 | 58 |
| *Bulgaria* | 56 (1.0) | 57 | 60 | 60 | 57 | 57 | 58 | 58 | 60 | 57 | 56 | 57 | 60 | 56 | 59 | 57 | 54 | 59 | 59 | 66 | 58 | 59 | 65 | 63 | 57 | 61 | 56 | 62 | 61 | 58 | 58 | 61 | 57 | 57 | 57 | 58 | 62 | 58 | 59 |
| Netherlands | 56 (0.7) | 58 | 59 | 58 | 58 | 56 | 58 | 57 | 63 | 57 | 55 | 59 | 59 | 56 | 59 | 58 | 57 | 60 | 59 | 64 | 60 | 60 | 65 | 60 | 58 | 61 | 56 | 62 | 62 | 58 | 61 | 62 | 57 | 54 | 59 | 60 | 62 | 58 | 59 |
| England | 56 (0.6) | 57 | 57 | 57 | 56 | 56 | 56 | 57 | 60 | 58 | 55 | 56 | 57 | 56 | 58 | 58 | 55 | 58 | 57 | 62 | 60 | 59 | 62 | 58 | 57 | 60 | 56 | 56 | 59 | 57 | 58 | 58 | 59 | 53 | 57 | 55 | 59 | 58 | 59 |
| Hungary | 56 (0.6) | 57 | 57 | 57 | 57 | 56 | 56 | 57 | 61 | 55 | 56 | 59 | 58 | 56 | 58 | 58 | 52 | 59 | 58 | 60 | 58 | 58 | 60 | 62 | 56 | 53 | 56 | 60 | 59 | 58 | 57 | 61 | 59 | 52 | 56 | 51 | 60 | 58 | 57 |
| *Austria* | 55 (0.6) | 56 | 56 | 58 | 54 | 56 | 57 | 56 | 60 | 55 | 56 | 60 | 59 | 55 | 55 | 57 | 53 | 59 | 58 | 60 | 59 | 59 | 62 | 60 | 56 | 60 | 55 | 61 | 63 | 59 | 55 | 59 | 59 | 57 | 57 | 54 | 58 | 58 | 58 |
| Slovak Republic | 54 (0.6) | 56 | 58 | 61 | 57 | 54 | 55 | 56 | 58 | 55 | 55 | 54 | 58 | 54 | 56 | 56 | 53 | 55 | 59 | 60 | 57 | 58 | 60 | 61 | 55 | 56 | 54 | 61 | 58 | 57 | 56 | 59 | 59 | 54 | 55 | 54 | 58 | 56 | 57 |
| United States | 54 (1.1) | 55 | 54 | 54 | 55 | 54 | 56 | 55 | 59 | 55 | 54 | 56 | 56 | 54 | 56 | 56 | 54 | 55 | 57 | 57 | 58 | 57 | 61 | 56 | 55 | 58 | 54 | 59 | 59 | 56 | 58 | 57 | 56 | 54 | 57 | 54 | 58 | 56 | 58 |
| Canada | 54 (0.5) | 55 | 56 | 55 | 54 | 54 | 54 | 55 | 59 | 55 | 54 | 54 | 56 | 54 | 57 | 56 | 53 | 55 | 56 | 61 | 56 | 57 | 62 | 57 | 55 | 58 | 54 | 58 | 57 | 56 | 57 | 57 | 57 | 53 | 56 | 53 | 58 | 56 | 58 |
| Australia | 54 (0.7) | 52 | 57 | 57 | 54 | 54 | 54 | 56 | 58 | 55 | 53 | 55 | 56 | 54 | 55 | 56 | 52 | 57 | 56 | 60 | 55 | 56 | 60 | 61 | 55 | 57 | 57 | 57 | 57 | 56 | 57 | 54 | 55 | 50 | 55 | 53 | 57 | 56 | 58 |
| Hong Kong | 53 (1.2) | 56 | 57 | 57 | 54 | 54 | 56 | 56 | 60 | 55 | 56 | 58 | 54 | 53 | 55 | 56 | 54 | 56 | 55 | 60 | 57 | 55 | 61 | 59 | 54 | 58 | 53 | 60 | 60 | 55 | 55 | 54 | 55 | 52 | 54 | 51 | 58 | 54 | 61 |
| *Germany* | 53 (0.8) | 55 | 55 | 57 | 53 | 54 | 54 | 55 | 58 | 53 | 53 | 57 | 56 | 53 | 55 | 55 | 51 | 57 | 55 | 60 | 56 | 56 | 60 | 57 | 54 | 52 | 53 | 59 | 60 | 55 | 56 | 58 | 56 | 51 | 54 | 51 | 58 | 55 | 58 |
| *Ireland* | 52 (0.7) | 54 | 52 | 51 | 52 | 52 | 54 | 53 | 53 | 53 | 53 | 52 | 54 | 52 | 53 | 53 | 53 | 53 | 56 | 57 | 56 | 55 | 59 | 56 | 53 | 57 | 52 | 52 | 53 | 54 | 55 | 58 | 56 | 54 | 53 | 54 | 56 | 54 | 53 |
| Sweden | 51 (0.5) | 52 | 53 | 56 | 51 | 52 | 52 | 53 | 57 | 53 | 52 | 54 | 56 | 51 | 54 | 53 | 50 | 56 | 55 | 59 | 55 | 56 | 60 | 58 | 53 | 53 | 52 | 55 | 53 | 53 | 55 | 58 | 56 | 50 | 53 | 54 | 57 | 53 | 54 |
| New Zealand | 50 (0.7) | 52 | 50 | 51 | 50 | 50 | 52 | 52 | 54 | 50 | 50 | 51 | 52 | 50 | 52 | 53 | 50 | 55 | 51 | 56 | 55 | 54 | 58 | 51 | 52 | 53 | 50 | 55 | 54 | 52 | 55 | 52 | 53 | 52 | 52 | 53 | 52 | 52 | 52 |
| Norway | 50 (0.6) | 51 | 53 | 54 | 52 | 51 | 50 | 52 | 57 | 51 | 50 | 52 | 53 | 50 | 53 | 52 | 49 | 54 | 53 | 58 | 54 | 55 | 59 | 54 | 54 | 52 | 50 | 52 | 55 | 52 | 55 | 58 | 53 | 50 | 52 | 51 | 56 | 52 | 52 |
| Switzerland | 50 (0.4) | 53 | 54 | 54 | 50 | 51 | 52 | 51 | 55 | 51 | 51 | 55 | 53 | 50 | 52 | 52 | 49 | 55 | 53 | 57 | 54 | 54 | 60 | 55 | 52 | 52 | 50 | 56 | 56 | 57 | 53 | 55 | 54 | 50 | 52 | 54 | 56 | 51 | 52 |
| Russian Federation | 50 (0.8) | 52 | 51 | 55 | 50 | 51 | 51 | 50 | 57 | 51 | 50 | 53 | 54 | 50 | 52 | 51 | 50 | 51 | 55 | 52 | 52 | 53 | 59 | 61 | 51 | 53 | 49 | 52 | 56 | 56 | 54 | 52 | 53 | 54 | 51 | 54 | 57 | 53 | 53 |
| Spain | 49 (0.4) | 50 | 50 | 50 | 50 | 49 | 50 | 50 | 52 | 50 | 50 | 49 | 52 | 49 | 50 | 50 | 47 | 50 | 51 | 54 | 52 | 53 | 55 | 50 | 50 | 48 | 49 | 53 | 51 | 51 | 52 | 52 | 52 | 50 | 51 | 48 | 52 | 52 | 52 |
| Scotland | 48 (0.8) | 50 | 50 | 50 | 49 | 48 | 49 | 50 | 51 | 50 | 48 | 49 | 50 | 48 | 51 | 51 | 49 | 50 | 50 | 54 | 52 | 52 | 55 | 50 | 49 | 53 | 48 | 50 | 46 | 50 | 50 | 49 | 51 | 46 | 49 | 46 | 52 | 52 | 52 |
| Iceland | 46 (0.6) | 47 | 48 | 49 | 48 | 46 | 47 | 49 | 53 | 47 | 46 | 46 | 49 | 46 | 47 | 47 | 44 | 49 | 49 | 53 | 50 | 48 | 51 | 50 | 48 | 47 | 48 | 49 | 53 | 48 | 49 | 48 | 47 | 44 | 48 | 43 | 51 | 48 | 51 |
| France | 46 (0.6) | 48 | 49 | 50 | 46 | 47 | 47 | 47 | 48 | 48 | 47 | 51 | 48 | 46 | 48 | 47 | 45 | 45 | 46 | 53 | 49 | 48 | 55 | 50 | 47 | 50 | 46 | 57 | 50 | 48 | 48 | 51 | 47 | 44 | 47 | 44 | 50 | 47 | 48 |
| Belgium (Fr) | 45 (0.7) | 47 | 49 | 47 | 46 | 46 | 47 | 47 | 49 | 46 | 45 | 50 | 47 | 45 | 47 | 46 | 42 | 48 | 46 | 52 | 48 | 49 | 53 | 48 | 46 | 48 | 42 | 54 | 52 | 47 | 47 | 49 | 49 | 41 | 46 | 40 | 49 | 46 | 47 |
| *Romania* | 45 (0.7) | 43 | 45 | 42 | 41 | 45 | 42 | 47 | 50 | 41 | 45 | 50 | 44 | 42 | 42 | 46 | 37 | 47 | 48 | 50 | 46 | 47 | 48 | 45 | 45 | 44 | 45 | 52 | 51 | 47 | 47 | 46 | 48 | 43 | 46 | 42 | 48 | 44 | 42 |
| Greece | 45 (0.5) | 43 | 42 | 43 | 46 | 40 | 42 | 46 | 48 | 42 | 45 | 45 | 46 | 45 | 42 | 48 | 43 | 47 | 45 | 49 | 48 | 47 | 48 | 46 | 45 | 47 | 45 | 50 | 46 | 47 | 48 | 48 | 44 | 43 | 46 | 45 | 48 | 42 | 45 |
| Denmark | 44 (0.4) | 45 | 46 | 49 | 45 | 44 | 45 | 46 | 50 | 45 | 44 | 48 | 47 | 44 | 47 | 45 | 42 | 49 | 45 | 53 | 47 | 48 | 53 | 51 | 45 | 46 | 44 | 48 | 52 | 46 | 48 | 52 | 47 | 41 | 46 | 44 | 49 | 45 | 47 |
| Iran, Islamic Rep. | 42 (0.6) | 46 | 41 | 44 | 42 | 43 | 45 | 46 | 44 | 43 | 42 | 41 | 43 | 44 | 44 | 43 | 42 | 45 | 46 | 46 | 44 | 44 | 44 | 48 | 41 | 42 | 44 | 43 | 44 | 44 | 43 | 48 | 45 | 40 | 44 | 40 | 45 | 45 | 47 |
| Latvia (LSS) | 42 (0.5) | 42 | 43 | 43 | 42 | 42 | 45 | 43 | 47 | 42 | 42 | 44 | 44 | 42 | 42 | 43 | 41 | 43 | 46 | 48 | 45 | 45 | 48 | 45 | 42 | 43 | 42 | 47 | 42 | 44 | 44 | 48 | 45 | 42 | 45 | 41 | 47 | 43 | 44 |
| Portugal | 41 (0.5) | 43 | 41 | 42 | 41 | 41 | 42 | 42 | 43 | 41 | 41 | 44 | 44 | 41 | 42 | 42 | 37 | 43 | 45 | 47 | 44 | 44 | 45 | 43 | 42 | 43 | 41 | 47 | 47 | 42 | 44 | 47 | 45 | 42 | 45 | 42 | 46 | 44 | 42 |
| Cyprus | 40 (0.4) | 43 | 42 | 43 | 40 | 40 | 42 | 41 | 42 | 42 | 41 | 44 | 42 | 40 | 42 | 42 | 40 | 47 | 43 | 45 | 42 | 43 | 45 | 43 | 40 | 45 | 45 | 47 | 43 | 43 | 42 | 48 | 44 | 40 | 41 | 44 | 43 | 42 | 45 |
| *Lithuania* | 38 (0.7) | 40 | 39 | 44 | 37 | 38 | 41 | 39 | 42 | 39 | 38 | 42 | 42 | 38 | 40 | 40 | 36 | 41 | 42 | 44 | 40 | 40 | 45 | 46 | 38 | 40 | 38 | 41 | 44 | 40 | 40 | 46 | 40 | 38 | 40 | 38 | 44 | 39 | 43 |
| *Colombia* | 35 (0.7) | 37 | 36 | 35 | 37 | 35 | 39 | 37 | 39 | 39 | 35 | 37 | 37 | 35 | 34 | 36 | 33 | 38 | 38 | 46 | 37 | 37 | 34 | 36 | 35 | 38 | 35 | 39 | 44 | 37 | 38 | 36 | 39 | 38 | 37 | 34 | 44 | 38 | 39 |
| South Africa | 26 (1.0) | 27 | 27 | 28 | 26 | 25 | 27 | 26 | 28 | 26 | 25 | 27 | 27 | 26 | 25 | 26 | 26 | 27 | 29 | 29 | 27 | 27 | 27 | 30 | 26 | 26 | 26 | 28 | 27 | 27 | 27 | 26 | 27 | 25 | 26 | 25 | 27 | 27 | 30 |
| **International Average** | **50 (0.7)** | **52** | **52** | **53** | **51** | **50** | **51** | **51** | **54** | **51** | **50** | **53** | **52** | **50** | **52** | **52** | **49** | **53** | **52** | **56** | **53** | **53** | **56** | **54** | **51** | **52** | **50** | **55** | **54** | **52** | **52** | **54** | **53** | **49** | **51** | **50** | **54** | **52** | **53** |

*Seventh grade in most countries; see Table 2 for more information about the grades tested in each country.

**Of the 135 items in the science test, some items had two parts and some extended-response items were scored on a two- or three-point scale, resulting in 146 total score points.

( ) Standard errors for the average percent of correct responses on all items appear in parentheses.   Standard errors for scores based on subsets of items are provided in Table B.4.

Because results are rounded to the nearest whole number, some totals may appear inconsistent.

Countries shown in italics did not satisfy one or more guidelines for sample participation rates, age/grade specifications, or classroom sampling procedures (see Figure A.3 for details).

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

SOURCE: IEA Third International Mathematics and Science Study (TIMSS), 1994-95.

The international averages of each country's selected items presented across the last row of the tables show that the selection of items for the participating countries varied somewhat in average difficulty, ranging from 55% to 59% at the eighth grade and from 49% to 56% at seventh grade. Despite these differences, the overall picture provided by both Tables B.1 and B.2 reveals that different item selections do not make a major difference in how well countries perform relative to each other. The items selected by some countries were more difficult than those selected by others. The relative performance of countries on the various item selections did vary somewhat, but generally not in a statistically significant manner.[5]

Comparing the diagonal element for a country with the overall average percentage correct shows the difference between performance on this subset of items 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 eighth-grade students in Singapore was 70%. The diagonal element shows that Singaporean students had about the same average percent correct (72%) based on the smaller set of items selected as relevant to the curriculum in Singapore as they did overall. In the eighth grade, most countries had a difference of less than 5 percentage points between the two performance measures, with the largest difference of 7% for the Russian Federation (65% compared to 58%). Performance differences between the entire TIMSS test and the subset of items selected for the TCMA were, in general, somewhat larger for seventh-grade students, including a few countries with an average performance that was about 10 percentage points higher on the subsets of items selected for the TCMA for their own students – Switzerland, France, and the Russian Federation. Even these increases are not particularly large, however, considering that France and Switzerland both selected less than one-quarter of the items at the seventh grade.

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 performance on the overall test in comparison to each other also had higher or lower relative performance on the different sets of items selected for the TCMA. For example, at the eighth grade, Singapore had the highest average percent correct on the test as a whole and on all of the different item selections, with Japan, Korea, and the Czech Republic 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 the most extreme example, consider the 49 score points selected by the Russian Federation for the seventh grade. The Russian students did substantially better on these items than on the test as a whole, with 61% correct responses to these items, on average, compared to 50% average correct on the items on the test as a whole.

---

[5] Small differences in performance in these tables are not statistically significant. The standard errors for the estimated average percent correct statistics can found in Tables B.3 and B.4. We can say with 95% confidence that the value for the entire population will fall between the sample estimate plus or minus two standard errors.

However, all other countries also did better on these particular items, with an international average of 54% for the items selected by the Russian Federation compared with 50% on the test as a whole. Only 8 of the 22 countries that performed better than the Russian students on the overall test also did so on the items selected by the Russian Federation. However, 10 countries with the same or higher overall performance were within 5 percentage points of the Russian students on these items.

The TCMA results provide evidence that the TIMSS science test provides a reasonable basis for comparing achievement for 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 substantially on the same items. Insofar as countries rejected items that would be difficult for their own 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.

## Table B.3  Standard Errors for the Test-Curriculum Matching Analysis Results – Science – Upper Grade (Eighth Grade*)

See Table B.1 for the Test-Curriculum Matching Analysis Results

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

| Country | Number of Score Points Included | Average Percent Correct on All Items | Singapore | Korea | Japan | Czech Republic | Netherlands | Bulgaria | Slovenia | Austria | England | Hungary | Belgium (Fl) | Australia | Slovak Republic | Sweden | Canada | Ireland | United States | Russian Federation | Germany | New Zealand | Norway | Hong Kong | Israel | Switzerland | Spain | Scotland | France | Iceland | Greece | Denmark | Belgium (Fr) | Latvia (LSS) | Portugal | Romania | Lithuania | Iran, Islamic Rep. | Cyprus | Kuwait | Colombia | South Africa |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 146** | | 109 | 59 | 86 | 136 | 102 | 112 | 140 | 131 | 124 | 129 | 98 | 133 | 129 | 125 | 121 | 90 | 146 | 96 | 129 | 126 | 111 | 68 | 102 | 105 | 146 | 97 | 73 | 146 | 111 | 70 | 58 | 113 | 133 | 99 | 120 | 87 | 78 | 131 | 112 | 74 |
| Singapore | | 70 (1.0) | 0.9 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 0.9 | 0.9 | 1.0 | 0.9 | 0.9 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 1.1 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 0.9 | 1.0 |
| Korea | | 66 (0.4) | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.3 | 0.4 | 0.4 | 0.4 | 0.4 | 0.3 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.5 | 0.4 | 0.3 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 |
| Japan | | 65 (0.3) | 0.3 | 0.4 | 0.4 | 0.3 | 0.4 | 0.3 | 0.3 | 0.3 | 0.4 | 0.3 | 0.3 | 0.3 | 0.3 | 0.3 | 0.3 | 0.3 | 0.3 | 0.3 | 0.3 | 0.4 | 0.3 | 0.3 | 0.4 | 0.3 | 0.3 | 0.4 | 0.3 | 0.3 | 0.3 | 0.4 | 0.3 | 0.3 | 0.3 | 0.3 | 0.3 | 0.3 | 0.3 | 1.0 | 0.8 | 0.3 |
| Czech Republic | | 64 (0.8) | 0.9 | 0.9 | 0.8 | 0.8 | 0.8 | 0.8 | 0.9 | 0.8 | 0.9 | 0.8 | 0.9 | 0.9 | 0.9 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.9 | 0.9 | 0.8 | 0.8 | 0.9 | 0.8 | 1.0 | 0.9 | 0.8 | 0.9 | 0.9 | 0.9 | 0.8 | 0.8 | 0.9 |
| Netherlands | | 62 (1.1) | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.0 | 1.0 | 1.0 | 1.1 | 1.0 | 1.2 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 0.9 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.0 | 1.3 | 1.1 | 1.0 | 1.1 | 1.1 | 1.1 | 1.0 | 1.1 | 1.1 | 1.1 | 0.9 |
| Bulgaria | | 62 (1.0) | 1.0 | 1.3 | 0.9 | 0.9 | 0.9 | 1.0 | 1.0 | 0.9 | 1.0 | 0.9 | 0.9 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 0.9 | 0.9 | 0.9 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 0.9 | 1.1 | 1.0 | 0.9 | 0.9 | 0.9 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 0.9 | 1.0 | 1.0 | 1.1 |
| Slovenia | | 62 (0.5) | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.5 | 0.5 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 |
| Austria | | 61 (0.7) | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 |
| England | | 61 (0.6) | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 |
| Hungary | | 61 (0.6) | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 |
| Belgium (Fl) | | 60 (1.1) | 1.1 | 1.2 | 1.1 | 1.2 | 1.1 | 1.1 | 1.1 | 1.1 | 1.2 | 1.1 | 1.1 | 1.1 | 1.2 | 1.2 | 1.1 | 1.2 | 1.1 | 1.1 | 1.2 | 1.2 | 1.2 | 1.1 | 1.2 | 1.2 | 1.1 | 1.2 | 1.1 | 1.1 | 1.2 | 1.0 | 1.2 | 1.1 | 1.0 | 1.1 | 1.2 | 1.2 | 1.1 | 1.2 | 1.2 | 1.2 |
| Australia | | 60 (0.7) | 0.7 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 |
| Slovak Republic | | 59 (0.6) | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.8 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 |
| Sweden | | 59 (0.6) | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.5 | 0.5 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 |
| Canada | | 59 (0.5) | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.4 | 0.5 | 0.5 | 0.5 |
| Ireland | | 58 (0.9) | 0.9 | 1.0 | 1.0 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 1.0 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 1.0 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 1.0 | 1.0 | 0.9 | 0.9 | 1.0 | 0.9 | 0.9 | 1.0 | 1.0 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 |
| United States | | 58 (1.0) | 0.9 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 0.9 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 0.9 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 |
| Russian Federation | | 58 (0.8) | 0.8 | 0.9 | 0.8 | 0.7 | 0.7 | 0.8 | 0.8 | 0.8 | 0.7 | 0.8 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.8 | 0.7 | 0.8 | 0.8 | 0.8 | 0.7 | 0.7 | 0.8 | 0.8 | 0.8 | 0.7 | 0.8 | 0.8 | 0.8 | 0.7 | 0.8 |
| Germany | | 58 (1.0) | 1.0 | 1.1 | 1.1 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.1 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.1 | 1.0 | 1.0 | 0.9 | 1.1 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.1 | 0.9 | 1.0 | 0.9 | 1.0 | 0.9 | 1.0 | 1.0 |
| New Zealand | | 58 (0.8) | 0.9 | 0.9 | 0.9 | 0.8 | 0.9 | 0.8 | 0.8 | 0.8 | 0.9 | 0.8 | 0.8 | 0.9 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.9 | 0.8 | 0.9 | 0.9 | 0.9 | 0.8 | 0.8 | 0.9 | 0.8 | 0.9 | 0.8 | 0.9 | 0.9 | 0.8 | 0.8 | 0.9 | 0.8 | 0.9 | 0.8 | 0.9 | 0.9 |
| Norway | | 58 (0.4) | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 |
| Hong Kong | | 58 (1.0) | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 0.9 | 1.0 | 0.9 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 |
| Israel | | 57 (1.1) | 1.1 | 1.1 | 1.1 | 1.1 | 1.0 | 1.1 | 1.1 | 1.0 | 1.1 | 1.0 | 1.0 | 1.0 | 1.1 | 1.1 | 1.0 | 1.1 | 1.1 | 0.9 | 1.0 | 1.1 | 1.0 | 1.2 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.2 | 1.1 | 1.1 | 1.2 | 1.1 | 1.0 | 1.1 | 1.0 | 1.1 | 1.1 | 1.1 | 1.1 |
| Switzerland | | 56 (0.5) | 0.5 | 0.5 | 0.5 | 0.4 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.4 | 0.5 | 0.4 | 0.5 | 0.5 | 0.5 | 0.4 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 |
| Spain | | 56 (0.4) | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 |
| Scotland | | 55 (1.0) | 1.0 | 1.1 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 0.9 | 1.0 | 0.9 | 1.0 | 1.0 |
| France | | 54 (0.6) | 0.7 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 |
| Iceland | | 52 (0.9) | 0.8 | 1.0 | 1.0 | 0.9 | 1.0 | 0.9 | 0.9 | 0.9 | 1.0 | 0.9 | 1.0 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 1.0 | 1.0 | 1.0 | 0.9 | 0.9 | 0.8 | 0.9 | 1.0 | 0.9 | 0.9 | 0.9 | 0.9 | 1.0 | 0.9 | 0.8 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 1.0 |
| Greece | | 52 (0.6) | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 |
| Denmark | | 51 (0.6) | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.5 | 0.6 |
| Belgium (Fr) | | 50 (0.7) | 0.6 | 0.7 | 0.6 | 0.7 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.8 | 0.8 | 0.6 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.7 | 0.7 |
| Latvia (LSS) | | 50 (0.6) | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 |
| Portugal | | 50 (0.6) | 0.8 | 0.6 | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.8 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 |
| Romania | | 50 (0.8) | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.9 | 0.8 | 0.8 | 0.9 | 0.9 | 0.9 | 0.9 | 0.9 | 0.8 | 0.9 | 0.8 | 0.8 | 0.8 | 0.8 | 0.9 | 0.9 | 0.8 | 0.8 | 0.8 | 0.9 | 1.0 | 0.8 | 0.8 | 0.9 | 0.9 | 0.8 | 0.8 | 0.8 | 0.9 | 0.8 | 0.9 | 0.9 | 0.9 | 0.9 |
| Lithuania | | 49 (0.7) | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 |
| Iran, Islamic Rep. | | 47 (0.6) | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 |
| Cyprus | | 47 (0.4) | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.5 | 0.4 | 0.3 | 0.4 | 0.4 | 0.4 | 0.4 | 0.5 | 0.4 | 0.4 | 0.5 | 0.4 | 0.5 | 0.4 | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.5 | 0.5 | 0.4 | 0.4 | 0.6 |
| Kuwait | | 43 (0.9) | 1.0 | 1.0 | 1.1 | 0.8 | 1.0 | 0.9 | 0.9 | 0.9 | 1.0 | 0.9 | 0.8 | 0.9 | 0.8 | 0.9 | 0.9 | 0.9 | 0.9 | 0.8 | 0.8 | 1.0 | 1.0 | 0.8 | 0.9 | 0.9 | 0.9 | 0.9 | 1.0 | 1.0 | 0.9 | 0.8 | 0.9 | 1.0 | 0.9 | 0.9 | 1.0 | 0.9 | 1.0 | 0.9 | 1.0 | 1.0 |
| Colombia | | 39 (0.8) | 0.8 | 0.8 | 0.9 | 0.8 | 0.8 | 0.8 | 0.9 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.9 | 0.8 | 0.8 | 0.8 | 0.8 | 0.9 | 0.9 | 0.9 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.9 | 0.7 | 0.8 | 0.8 | 0.8 | 0.9 | 0.8 | 0.9 | 0.8 | 0.8 | 0.9 |
| South Africa | | 27 (1.3) | 1.3 | 1.2 | 1.4 | 1.2 | 1.3 | 1.3 | 1.3 | 1.3 | 1.3 | 1.3 | 1.3 | 1.3 | 1.3 | 1.3 | 1.3 | 1.3 | 1.3 | 1.2 | 1.3 | 1.4 | 1.3 | 1.4 | 1.2 | 1.3 | 1.3 | 1.3 | 1.2 | 1.3 | 1.4 | 1.3 | 1.2 | 1.2 | 1.3 | 1.3 | 1.3 | 1.3 | 1.4 | 1.3 | 1.3 | 1.3 |
| **International Average** | | **55 (0.7)** | **0.7** | **0.8** | **0.8** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.8** | **0.7** | **0.7** | **0.7** | **0.8** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.7** | **0.8** |

*Eighth grade in most countries; see Table 2 for more information about the grades tested in each country.
**Of the 135 items in the science test, some items had two parts and some extended-response items were scored on a two- or three-point scale, resulting in 146 total score points.
( ) Standard errors for the average percent of correct responses on all items 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 Table B.1. Because results are rounded to the nearest whole number, some totals may appear inconsistent.
Countries shown in italics did not satisfy one or more guidelines for sample participation rates, age/grade specifications, or classroom sampling procedures (see Figure A.3 for details).
Because population coverage falls below 65% Latvia is annotated LSS for Latvian Speaking Schools only.

SOURCE: IEA Third International Mathematics and Science Study (TIMSS), 1994-95.

# Table B.4 Standard Errors for the Test-Curriculum Matching Analysis Results - Science - Lower Grade (Seventh Grade*)

See Table B.3 for the Test-Curriculum Matching Analysis Results

**Instructions:** Read *across* the row for the standard error for the score based on the test items included by each of the countries across the top.

Read *down* the column under a country name for the standard error for the score of the country down the left on the items included by the country listed on the top.

Read along the *diagonal* for the standard error for the score for each different country based on its own decisions about the test items to include.

| Country | Average Percent Correct on All Items 146** | Singapore 83 | Korea 44 | Japan 46 | Czech Republic 108 | Slovenia 132 | Belgium (Fl) 45 | Bulgaria 106 | Netherlands 34 | England 105 | Hungary 98 | Austria 53 | Slovak Republic 111 | United States 146 | Canada 79 | Australia 94 | Hong Kong 33 | Germany 88 | Ireland 60 | Sweden 86 | New Zealand 110 | Norway 93 | Switzerland 36 | Russian Federation 49 | Spain 132 | Scotland 50 | Iceland 146 | France 26 | Belgium (Fr) 23 | Romania 92 | Greece 72 | Denmark 20 | Iran, Islamic Rep. 48 | Latvia (LSS) 46 | Portugal 81 | Cyprus 29 | Colombia 79 | Lithuania 109 | South Africa 26 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | *(Number of Score Points Included)* | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **Singapore** | 61 (1.2) | 1.3 | 0.4 | 0.3 | 0.7 | 0.5 | 0.5 | 1.0 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 1.0 | 0.5 | 0.7 | 1.3 | 0.8 | 0.7 | 0.6 | 0.7 | 0.6 | 0.4 | 0.6 | 0.4 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.5 | 0.4 | 0.6 | 0.6 | 0.5 | 0.4 | 0.8 | 0.6 | 1.0 |
| **Korea** | 61 (0.4) | 0.4 | 0.5 | 0.4 | 0.8 | 1.0 | 0.5 | 1.1 | 0.8 | 0.8 | 1.1 | 0.7 | 0.6 | 1.2 | 0.5 | 0.5 | 0.5 | 0.9 | 0.8 | 0.7 | 0.7 | 1.0 | 0.5 | 0.6 | 0.4 | 0.5 | 0.4 | 0.6 | 0.6 | 0.6 | 0.5 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 | 0.4 | 0.5 |
| **Japan** | 59 (0.3) | 0.3 | 0.4 | 0.4 | 0.8 | 0.8 | 0.3 | 0.4 | 0.4 | 0.4 | 0.9 | 0.3 | 0.3 | 0.4 | 0.5 | 0.3 | 0.9 | 0.8 | 0.7 | 0.6 | 0.7 | 0.8 | 0.4 | 0.4 | 0.3 | 0.5 | 0.3 | 0.4 | 0.5 | 0.4 | 0.3 | 0.4 | 0.3 | 0.4 | 0.3 | 0.4 | 0.8 | 0.3 | 0.4 |
| **Czech Republic** | 58 (0.8) | 0.7 | 1.0 | 0.7 | 0.8 | 0.8 | 0.7 | 1.0 | 0.8 | 0.7 | 0.9 | 0.7 | 0.8 | 0.8 | 0.7 | 0.8 | 0.9 | 0.7 | 0.7 | 0.6 | 0.7 | 0.8 | 0.8 | 0.7 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.7 | 0.9 | 0.6 | 0.9 | 0.8 | 0.7 | 0.8 | 0.8 | 0.8 |
| **Slovenia** | 57 (0.5) | 0.5 | 0.6 | 0.6 | 0.8 | 0.5 | 0.4 | 0.6 | 0.6 | 0.5 | 0.6 | 0.7 | 0.5 | 0.8 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.7 | 0.5 | 0.5 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.7 | 0.5 | 0.6 | 0.5 | 0.6 | 0.5 | 0.5 | 0.7 |
| **Belgium (Fl)** | 57 (0.5) | 0.5 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.7 | 0.6 | 0.7 | 0.7 | 0.5 | 0.5 | 0.6 | 0.6 | 0.7 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.5 | 0.8 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.8 | 0.6 | 0.5 | 0.6 |
| **Bulgaria** | 56 (1.0) | 1.0 | 1.1 | 1.1 | 1.0 | 1.0 | 1.1 | 1.1 | 1.1 | 1.0 | 1.0 | 1.2 | 0.9 | 1.0 | 1.0 | 1.0 | 1.1 | 1.0 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 0.8 | 1.0 | 1.2 | 1.0 | 1.2 | 1.4 | 0.9 | 1.0 | 1.1 | 1.0 | 1.0 | 0.9 | 1.0 | 1.0 | 1.0 | 1.1 |
| **Netherlands** | 56 (0.7) | 0.7 | 0.9 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.9 | 0.6 | 0.8 | 0.6 | 0.7 | 0.7 | 0.6 | 0.8 | 0.9 | 0.8 | 0.7 | 0.7 | 0.8 | 0.8 | 1.0 | 0.7 | 0.8 | 0.7 | 0.7 | 0.9 | 1.0 | 0.7 | 0.7 | 1.1 | 0.9 | 0.9 | 0.7 | 1.0 | 0.7 | 0.7 | 0.8 |
| **England** | 56 (0.6) | 0.6 | 0.7 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.8 | 0.6 | 0.7 | 0.7 | 0.6 | 0.8 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 |
| **Hungary** | 56 (0.6) | 0.6 | 0.7 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 | 0.6 | 0.7 | 0.7 | 0.8 | 0.7 | 0.6 | 0.7 | 0.6 | 0.8 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 |
| **Austria** | 55 (0.6) | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.8 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.8 | 0.7 | 0.6 | 0.7 | 0.8 | 0.6 | 0.6 | 0.6 | 0.8 | 0.6 | 0.6 | 0.7 |
| **Slovak Republic** | 54 (0.6) | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 | 0.5 | 0.6 | 0.7 |
| **United States** | 54 (1.1) | 1.0 | 1.1 | 1.1 | 1.0 | 1.0 | 1.1 | 1.0 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.2 | 1.1 | 1.1 | 1.0 | 1.1 | 1.1 | 1.2 | 1.1 | 1.1 | 1.2 | 1.1 | 1.2 | 1.0 | 1.0 | 1.1 | 1.3 | 1.1 | 1.2 | 1.1 | 1.3 | 1.1 | 1.0 | 1.1 |
| **Canada** | 54 (0.5) | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.6 | 0.5 | 0.7 | 0.5 | 0.5 | 0.6 | 0.5 | 0.5 | 0.6 | 0.5 | 0.7 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.7 | 0.6 | 0.5 | 0.6 | 0.5 | 0.6 | 0.8 | 0.5 | 0.5 | 0.6 | 0.5 | 0.6 | 0.5 | 0.6 | 0.5 | 0.5 | 0.6 |
| **Australia** | 54 (0.7) | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.8 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.8 |
| **Hong Kong** | 53 (1.2) | 1.3 | 1.3 | 1.3 | 1.3 | 1.2 | 1.1 | 1.2 | 1.3 | 1.2 | 1.2 | 1.3 | 1.2 | 1.2 | 1.3 | 1.3 | 1.3 | 1.2 | 1.1 | 1.2 | 1.3 | 1.2 | 1.4 | 1.2 | 1.2 | 1.4 | 1.2 | 1.4 | 1.3 | 1.2 | 1.3 | 1.3 | 1.2 | 1.2 | 1.2 | 1.3 | 1.2 | 1.2 | 1.3 |
| **Germany** | 53 (0.8) | 0.8 | 0.9 | 0.8 | 0.9 | 0.8 | 0.8 | 0.8 | 0.9 | 0.8 | 0.8 | 0.9 | 0.8 | 0.8 | 0.8 | 0.8 | 0.9 | 0.9 | 0.8 | 0.7 | 0.8 | 0.8 | 0.9 | 0.8 | 0.8 | 0.9 | 0.8 | 1.0 | 0.9 | 0.7 | 0.9 | 1.1 | 0.8 | 0.9 | 0.8 | 1.0 | 0.8 | 0.8 | 0.9 |
| **Ireland** | 52 (0.7) | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.9 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.9 | 0.9 | 0.7 | 0.8 | 0.8 | 0.7 | 0.7 | 0.7 | 0.9 | 0.7 | 0.6 | 0.8 |
| **Sweden** | 51 (0.5) | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.7 | 0.6 | 0.5 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.7 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.8 | 0.5 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.5 | 0.7 |
| **New Zealand** | 50 (0.7) | 0.7 | 0.8 | 0.8 | 0.7 | 0.7 | 0.6 | 0.7 | 0.8 | 0.7 | 0.8 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.8 | 0.7 | 0.6 | 0.7 | 0.7 | 0.7 | 1.0 | 0.7 | 0.8 | 0.8 | 0.7 | 1.0 | 1.0 | 0.7 | 0.8 | 0.9 | 0.7 | 0.8 | 0.7 | 0.9 | 0.7 | 0.7 | 0.8 |
| **Norway** | 50 (0.6) | 0.6 | 0.7 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.7 | 0.7 | 0.7 | 0.7 | 0.6 | 0.7 | 0.6 | 0.8 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.8 | 0.6 | 0.6 | 0.7 |
| **Switzerland** | 50 (0.6) | 0.6 | 0.7 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.8 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.8 | 0.8 | 0.6 | 0.7 | 0.8 | 0.6 | 0.7 | 0.6 | 0.8 | 0.6 | 0.6 | 0.7 |
| **Russian Federation** | 50 (0.8) | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.7 | 0.8 | 0.8 | 0.7 | 0.8 | 0.9 | 0.8 | 0.8 | 0.7 | 0.9 | 0.9 | 0.8 | 0.8 | 0.7 | 0.8 | 0.8 | 0.8 | 0.7 | 0.8 | 0.8 | 0.8 | 1.0 | 0.9 | 0.8 | 0.8 | 1.0 | 0.7 | 0.8 | 0.7 | 0.9 | 0.7 | 0.8 | 0.8 |
| **Spain** | 49 (0.4) | 0.4 | 0.5 | 0.5 | 0.4 | 0.4 | 0.4 | 0.5 | 0.5 | 0.4 | 0.5 | 0.5 | 0.4 | 0.4 | 0.4 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.4 | 0.4 | 0.5 | 0.4 | 0.5 | 0.6 | 0.4 | 0.5 | 0.6 | 0.5 | 0.5 | 0.5 | 0.6 | 0.4 | 0.4 | 0.6 |
| **Scotland** | 48 (0.8) | 0.8 | 0.9 | 0.8 | 0.8 | 0.7 | 0.7 | 0.7 | 0.9 | 0.7 | 0.8 | 0.8 | 0.8 | 0.7 | 0.8 | 0.8 | 0.9 | 0.8 | 0.7 | 0.8 | 0.8 | 0.8 | 1.0 | 0.8 | 0.8 | 0.8 | 0.7 | 1.0 | 1.0 | 0.7 | 0.8 | 0.9 | 0.7 | 0.8 | 0.7 | 0.9 | 0.7 | 0.7 | 0.8 |
| **Iceland** | 46 (0.6) | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.8 | 0.8 | 0.6 | 0.6 | 0.8 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.7 |
| **France** | 46 (0.6) | 0.6 | 0.7 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.8 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.8 | 0.6 | 0.6 | 0.7 | 0.6 | 0.8 | 0.8 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 | 0.6 | 0.9 | 0.6 | 0.6 | 0.7 |
| **Belgium (Fr)** | 45 (0.7) | 0.6 | 0.7 | 0.8 | 0.8 | 0.7 | 0.7 | 0.6 | 0.8 | 0.7 | 0.8 | 0.9 | 0.8 | 0.7 | 0.8 | 0.7 | 0.8 | 0.8 | 0.6 | 0.7 | 0.7 | 0.7 | 0.9 | 0.7 | 0.8 | 0.8 | 0.7 | 0.8 | 0.8 | 0.7 | 0.8 | 1.0 | 0.7 | 0.8 | 0.6 | 0.9 | 0.7 | 0.6 | 0.8 |
| **Romania** | 45 (0.8) | 0.7 | 0.8 | 0.8 | 0.8 | 0.7 | 0.7 | 0.8 | 0.8 | 0.7 | 0.8 | 0.9 | 0.8 | 0.8 | 0.7 | 0.8 | 0.8 | 0.8 | 0.8 | 0.7 | 0.8 | 0.7 | 0.8 | 0.7 | 0.8 | 0.8 | 0.7 | 0.9 | 0.8 | 0.8 | 0.8 | 1.0 | 0.8 | 0.8 | 0.6 | 0.9 | 0.7 | 0.6 | 0.8 |
| **Greece** | 45 (0.5) | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.5 | 0.5 | 0.6 | 0.5 | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.5 | 0.6 | 0.5 | 0.5 | 0.5 | 0.5 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.5 | 0.7 | 0.7 | 0.5 | 0.5 | 0.7 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.7 |
| **Denmark** | 44 (0.4) | 0.4 | 0.6 | 0.6 | 0.5 | 0.5 | 0.5 | 0.5 | 0.6 | 0.5 | 0.5 | 0.6 | 0.5 | 0.4 | 0.5 | 0.5 | 0.6 | 0.5 | 0.5 | 0.4 | 0.5 | 0.5 | 0.6 | 0.5 | 0.4 | 0.5 | 0.4 | 0.6 | 0.6 | 0.5 | 0.5 | 0.7 | 0.5 | 0.5 | 0.4 | 0.6 | 0.4 | 0.4 | 0.6 |
| **Iran, Islamic Rep.** | 42 (0.6) | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.5 | 0.6 | 0.7 | 0.6 | 0.6 | 0.7 | 0.6 | 0.7 | 0.7 | 0.6 | 0.7 | 0.8 | 0.6 | 0.6 | 0.5 | 0.7 | 0.7 | 0.6 | 0.7 |
| **Latvia (LSS)** | 42 (0.5) | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.7 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.5 | 0.6 | 0.6 | 0.7 | 0.5 | 0.5 | 0.6 | 0.5 | 0.7 | 0.7 | 0.5 | 0.6 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 | 0.6 | 0.6 | 0.8 |
| **Portugal** | 41 (0.5) | 0.5 | 0.6 | 0.6 | 0.5 | 0.5 | 0.5 | 0.6 | 0.5 | 0.5 | 0.5 | 0.6 | 0.5 | 0.5 | 0.5 | 0.5 | 0.6 | 0.5 | 0.5 | 0.4 | 0.5 | 0.4 | 0.6 | 0.5 | 0.5 | 0.5 | 0.4 | 0.6 | 0.7 | 0.5 | 0.5 | 0.8 | 0.5 | 0.6 | 0.5 | 0.7 | 0.5 | 0.5 | 0.7 |
| **Cyprus** | 40 (0.4) | 0.4 | 0.5 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 | 0.6 | 0.4 | 0.4 | 0.5 | 0.4 | 0.5 | 0.6 | 0.4 | 0.5 | 0.6 | 0.5 | 0.5 | 0.4 | 0.6 | 0.4 | 0.4 | 0.4 |
| **Lithuania** | 38 (0.7) | 0.8 | 0.8 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.8 | 0.7 | 0.8 | 0.8 | 0.7 | 0.7 | 0.7 | 0.8 | 0.8 | 0.8 | 0.8 | 0.7 | 0.8 | 0.8 | 0.8 | 0.8 | 0.7 | 0.8 | 0.7 | 0.9 | 0.9 | 0.8 | 0.8 | 1.0 | 0.8 | 0.7 | 0.7 | 0.9 | 0.8 | 0.7 | 0.8 |
| **Colombia** | 35 (0.7) | 0.6 | 0.9 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.6 | 0.8 | 0.6 | 0.6 | 0.6 | 0.6 | 0.7 | 0.8 | 0.8 | 0.7 | 0.8 | 0.7 | 0.9 | 1.0 | 0.8 | 0.7 | 1.0 | 1.0 | 0.7 | 0.6 | 1.0 | 0.8 | 0.8 | 0.6 |
| **South Africa** | 26 (1.0) | 1.0 | 1.1 | 1.0 | 1.0 | 0.7 | 0.9 | 1.0 | 1.2 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.1 | 1.1 | 1.1 | 1.0 | 0.9 | 1.1 | 1.1 | 1.1 | 1.2 | 0.9 | 1.1 | 1.1 | 1.0 | 1.1 | 1.2 | 1.1 | 1.1 | 1.1 | 1.0 | 1.0 | 1.1 | 1.1 | 1.1 | 1.0 | 1.1 |
| **International Average** | 50 (0.7) | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 | 0.8 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 | 0.8 | 0.7 | 0.7 | 0.7 |

*Seventh grade in most countries; see Table 2 for more information about the grades tested in each country.

**Of the 135 items in the science test, some items had two parts and some extended-response items were scored on a two- or three-point scale, resulting in 146 total score points.

( ) Standard errors for the average percent of correct responses on all items 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 Table B.2. Because results are rounded to the nearest whole number, some totals may appear inconsistent.

Countries shown in italics did not satisfy one or more guidelines for sample participation rates, age/grade specifications, or classroom sampling procedures (see Figure A.3 for details).

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

SOURCE: IEA Third International Mathematics and Science Study (TIMSS), 1994-95.