

Matthias von Davier · Eugenio Gonzalez
Irwin Kirsch · Kentaro Yamamoto *Editors*

The Role of International Large-Scale Assessments: Perspectives from Technology, Economy, and Educational Research

 Springer

The Role of International Large-Scale Assessments: Perspectives from Technology, Economy, and Educational Research

Matthias von Davier • Eugenio Gonzalez
Irwin Kirsch • Kentaro Yamamoto
Editors

The Role of International Large-Scale Assessments: Perspectives from Technology, Economy, and Educational Research

 Springer

Editors

Matthias von Davier
Educational Testing Service
Princeton, NJ
USA

Irwin Kirsch
Educational Testing Service
Princeton, NJ
USA

Eugenio Gonzalez
Educational Testing Service
Princeton, NJ
USA

Kentaro Yamamoto
Educational Testing Service
Princeton, NJ
USA

ISBN 978-94-007-4628-2 ISBN 978-94-007-4629-9 (eBook)
DOI 10.1007/978-94-007-4629-9
Springer Dordrecht Heidelberg London New York

Library of Congress Control Number: 2012944183

© Springer Science+Business Media Dordrecht 2013
No part of this work may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, microfilming, recording or otherwise, without written permission from the Publisher, with the exception of any material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

Preface

This volume consists of the papers presented at the International Large-Scale Assessment conference held at Educational Testing Service (ETS) in March 2011. The conference was designed to present and discuss multidisciplinary issues related to the use and implementation of international large-scale assessments. It was geared towards funders, policymakers, managers, and technical staff of international large-scale assessment programs. The conference covered the following topics: large-scale assessments as change agents; technologies in large-scale assessments; the role of assessing cognitive skills in international growth and development; the utility and need for assessing noncognitive skills in large-scale assessments; the contributions of international large-scale studies in civic engagement and citizenship; and the role of large-scale assessments in research on educational effectiveness and school development.

The different perspectives brought together in this volume reflect the changing landscape of these surveys both in terms of the widening group of researchers and policymakers interested in these data and the issues that should and could be addressed. Among these new directions is a surge in the use of large-scale assessment data in the field of economics as well as an increased interest in how these assessments can inform and the use of technology in education and assessment. In addition, research in civics and citizenship studies as well as investigations focusing on motivation, interest, and self concept indicate great interest in the data collected in international comparisons of education and skills.

Bringing together expert authors to produce such a volume would not have been possible without the generous support of ETS. ETS provided funding for the speakers and sustenance for all conference participants. Also contributing to the success of the conference was the set of invited experts who agreed to provide reflection and discussion of the invited presentations. We would like to thank Esther Care, from Assessment & Teaching of 21st Century Skills Project, Melbourne, Australia; Guido Schwerdt, from the Program on Educational Policy and Governance, Harvard University, United States; Erik Amnå, Professor of Political Science, Örebro University, Sweden; Patrick Kyllonen, from ETS; and David Kaplan, Professor of Quantitative Methods, University of Wisconsin-Madison. These discussants helped

to sharpen and shape the questions and answers and, finally, to shape the final version of the chapters presented in this volume.

Finally, we want to acknowledge individuals at ETS who made the conference and this publication possible. The conference would not have been possible without the help of Judy Mendez and Judy Shahbazian, who helped with organization, arrangements, and general support of this endeavor, and Larry Hanover, who provided editorial reviews.

Matthias von Davier
Eugenio Gonzalez
Irwin Kirsch
Kentaro Yamamoto

Table of Contents

1 On the Growing Importance of International Large-Scale Assessments	1
Irwin Kirsch, Marylou Lennon, Matthias von Davier, Eugenio Gonzalez and Kentaro Yamamoto	
2 International Large-Scale Assessments as Change Agents	13
Jo Ritzen	
3 Technologies in Large-Scale Assessments: New Directions, Challenges, and Opportunities	25
Michal Beller	
4 The Role of International Assessments of Cognitive Skills in the Analysis of Growth and Development	47
Eric A. Hanushek and Ludger Woessmann	
5 The Utility and Need for Incorporating Noncognitive Skills Into Large-Scale Educational Assessments	67
Henry M. Levin	
6 The Contributions of International Large-Scale Studies in Civic Education and Engagement	87
Judith Torney-Purta and Jo-Ann Amadeo	
7 The Role of Large-Scale Assessments in Research on Educational Effectiveness and School Development	115
Eckhard Klieme	
8 Prospects for the Future: A Framework and Discussion of Directions for the Next Generation of International Large-Scale Assessments	149
Henry Braun	

Contributors

Jo-Ann Amadeo, Ph.D. Department of Psychology, Marymount University, Arlington, 22207 VA, USA
e-mail: jamadeo@marymount.edu

Michal Beller, Ph.D. RAMA—The National Authority for Measurement and Evaluation in Education, Ministry of Education, 125 Begin Blvd. 12th Floor, Tel Aviv 67012, Israel
e-mail: mbeller.rama@education.gov.il

Henry Braun, Ph.D. Center for the Study of Testing, Evaluation and Education Policy, Lynch School of Education Boston College, 140 Commonwealth Ave, 02467 Chestnut Hill, MA, USA
e-mail: braunh@bc.edu

Dr. Matthias von Davier Educational Testing Service, Rosedale Road, MS 13-E, 08541 Princeton, NJ, USA
e-mail: mvondavier@ets.org

Eugenio J. Gonzalez, Ph.D. Educational Testing Service, Rosedale Road, MS 13-E, 08541 Princeton, NJ, USA
e-mail: egonzalez@ets.org

Eric Hanushek, Ph.D. Hoover Institution, National Bureau of Economic Research and CESifo, Stanford University, Stanford, CA 94305-6010, USA
e-mail: hanushek@stanford.edu

Irwin Kirsch, Ph.D. Educational Testing Service, Rosedale Road, MS 13-E, 08541 Princeton, NJ, USA
e-mail: ikirsch@ets.org

Dr. Eckhard Klieme Center for Research on Educational Quality and Evaluation, German Institute for International Educational Research (DIPF), Goethe University, Schloßstraße 29, 60486 Frankfurt am Main, Germany
e-mail: klieme@dipf.de

Marylou Lennon Educational Testing Service, Rosedale Road, MS 13-E, 08541 Princeton, NJ, USA
e-mail: mlennon@ets.org

Henry M. Levin, Ph.D. Teachers College, Columbia University, 525 West 120 Street, New York, NY 10027, USA
e-mail: HL361@columbia.edu

Dr. Jo Ritzen Empower European Universities, International Economics of Science, Technology and Higher Education, Maastricht School of Governance, UNU-MERIT, Kloosterweg 54 Bunde, Netherlands
e-mail: jo.ritzen@empowereu.org

Judith Torney-Purta, Ph.D. Department of Human Development and Quantitative Methodology, University of Maryland, College Park, MD 20742, USA
e-mail: jtpurta@umd.edu

Dr. Ludger Woessmann Ifo Institute for Economic Research and CESifo, University of Munich, Poschingerstr. 5, 81679 Munich, Germany
e-mail: woessmann@ifo.de

Kentaro Yamamoto, Ph.D. Educational Testing Service, Rosedale Road, MS 13-E, 08541 Princeton, NJ, USA
e-mail: kyamamoto@ets.org

Chapter 1

On the Growing Importance of International Large-Scale Assessments

Irwin Kirsch, Marylou Lennon, Matthias von Davier, Eugenio Gonzalez and Kentaro Yamamoto

Large-scale assessments that compare the skills and knowledge demonstrated by populations across countries are relatively recent endeavors. These assessments have expanded in scope over time in response to increasing concern about the distribution of human capital and the growing recognition that skills contribute to the prosperity of nations and to better lives for individuals in those nations. Broadly defined, large-scale assessments are surveys of knowledge, skills, or behaviors in a given domain. The goal of large-scale assessments is to describe a population, or populations, of interest. As such, these assessments focus on group scores and can be distinguished from large-scale testing programs that focus on assessing individuals. The major themes laid out here—that these large-scale assessments have expanded over the past 50 years to include a greater number of surveys focusing on a broader range of populations and skill domains, that this work has led to new methodologies and modes of assessment, and that these assessments have grown to address the increasingly challenging questions posed by researchers and policymakers around the world—will be addressed in greater detail in each of the remaining chapters. We begin here by providing a general overview of the history of international large-scale assessments and the broadening role that these surveys have played in influencing policymakers around the world.

I. Kirsch (✉) · M. Lennon · E. Gonzalez · K. Yamamoto · M. von Davier
Educational Testing Service, Rosedale Road, MS 13-E, 08541 Princeton, NJ, USA
e-mail: ikirsch@ets.org

M. Lennon
e-mail: mlennon@ets.org

M. von Davier
e-mail: mvondavier@ets.org

E. Gonzalez
e-mail: egonzalez@ets.org

K. Yamamoto
e-mail: kyamamoto@ets.org

Large-Scale Assessments of Student Populations

Prior to the late 1950s, no systematic or standardized comparative data focusing on skills and knowledge had been collected at national or international levels. The foundational work in this area began with a focus on student skills. In 1958, a group of scholars met at the UNESCO Institute for Education in Hamburg to discuss issues associated with collecting systematic data about schools and education systems in a cross-country context. That meeting led to a study designed to investigate the feasibility of developing and conducting an assessment of 13-year-olds in 12 countries. The pilot 12-country study focused on five domains including mathematics, reading comprehension, geography, science, and non-verbal ability and was conducted between 1959 and 1962. The results of this pioneering study demonstrated the feasibility of conducting a large-scale international survey in which common cognitive instruments worked in a comparable manner across different cultures and languages (Naemi et al. [in press](#)).

A parallel effort in the United States began around this same time under the leadership of several prominent American scholars and policymakers. Francis Keppel, the US Commissioner of Education in the mid-1960s, was responsible for reporting to Congress about the condition of education in America. Keppel was concerned about the lack of systematic data on the educational attainment of students in the country. As he pointed out, most of the information that had been collected to date focused on the inputs of education—such as the number of classrooms, dollars spent, and school enrollment figures—rather than on the output of education in terms of skills and knowledge. This concern led Keppel to invite Ralph Tyler, Director of the Center for Advanced Study in the Behavioral Sciences at Stanford University, to develop a plan for the periodic national assessment of student learning. With Tyler as chair, the Carnegie Foundation funded two planning meetings for national student assessments in 1963 and 1964. A technical advisory group was formed in 1965 and chaired by John Tukey, head of the Department of Statistics at Princeton University and Associate Executive Director of Research Information Systems at AT&T Bell Laboratories. This work led to the National Assessment of Educational Progress (NAEP), which conducted its first assessment of in-school 17-year-olds in citizenship, science, and writing in 1969.

Rather than build an assessment around classical test theory models that focused primarily on measuring individual differences, Tyler's vision for NAEP was to focus on what groups of students knew and could do. In this scheme, groups were defined by educationally relevant variables such as gender, immigrant status and ethnic background. Tyler's idea was to convene panels of subject-matter experts, to have them identify key educational objectives within the domains to be assessed, and then to develop test items based on those objectives. Reports from these assessments would then focus on the performance of national populations or subgroups rather than individual students. Additionally, Tyler was adamant that assessment results not be based on any type of norm-referenced perspective such as grade-level norms.

As surveys such as NAEP progressed, one of the criticisms that arose was that interpretations were quite limited because they were fixed to the individual items used in the assessments. In the 1980s, Educational Testing Service (ETS) bid on and won the contract to conduct NAEP based on a monograph written by Samuel Messick, Albert Beaton, and Frederic Lord. In “National Assessment of Educational Progress reconsidered: a new design for a new era,” they introduced the idea of using Item Response Theory (IRT), an analytic approach with important advantages compared to the classical methods used previously in that it directly supports the creation of comparable scales across multiple forms of a test. In addition to incorporating IRT-based methodology, the work on NAEP led to developments of new methodologies including marginal estimation procedures that could optimize the reporting of proficiency scales based on very complex designs (von Davier et al. 2006).

NAEP and other surveys began by using a version of matrix sampling, an approach that is based on utilizing multiple, partially overlapping test forms. The introduction of balanced incomplete block (BIB) spiraling to large-scale assessment was another important innovation introduced in the 1980s. The goal of these developments was to broaden the item pool represented in the BIB-spiraled test forms in order to maximize the coverage of the constructs of interest. As an example, NAEP 8th grade mathematics assessments include a large number of test items across five subdomains of mathematics: number properties and operations; measurement; geometry; data analysis, statistics and probability; and algebra. Using BIB spiraling, each student is asked to respond to only a small subset of these items, reducing the burden on the test taker. Striking this balance of construct coverage and the reduction of test taker burden requires utilizing covariance information to create proficiency scales and the ability to generalize to populations of interest.

The use of IRT in combination with BIB-spiraling and covariance information among domains has made it possible to both broaden content coverage to include relevant facets of the cognitive constructs of interest and to extend inferences beyond individual items to the underlying construct. Just as we sample individuals and then make generalizations to populations, these scales, constructed with the help of IRT, represent a construct broadly and therefore make it possible to generalize beyond the specific items in the assessment to the construct domain that those items represent. These methodologies originally developed for NAEP are utilized in all the large-scale assessments covered in this volume, including the studies currently conducted by the International Association for the Evaluation of Educational Achievement (IEA) and the Organisation for Economic Co-operation and Development (OECD) that will be described next. Methodological innovations such as these have contributed to the growth and expansion of international large-scale assessments and allowed us to move beyond the questions raised by Tyler and others in the 1960s and 1970s and focus on increasingly complex questions raised by policymakers today.

Following the initial work that occurred from the 1960s through the 1980s, international large-scale assessments of student skills have expanded tremendously in terms of the number of assessments and participating countries. IEA continued to conduct important periodic large-scale international studies and, starting in 1995,

began to conduct continuous assessment cycles for the Trends in Mathematics and Science Study (TIMSS) followed by the Progress in Reading Literacy Study (PIRLS) in 2001. TIMSS is conducted every 4 years and focuses on achievement in mathematics and science at the fourth and eighth grades. PIRLS runs on a 5-year cycle and assesses how well children read after 4 years of primary school. By 2007, some 60 countries participated in TIMSS and over 40 countries participated in PIRLS. At the end of the 1990s, the OECD began the Programme for International Student Assessment (PISA) cycle of studies. PISA assesses the skills of 15-year-olds with the goal of gathering information about how well students have acquired the knowledge and skills essential for full participation in society. The first assessment was conducted in 2000 in over 30 countries and focused on the domains of reading, mathematics, and science. Since then, PISA has expanded in terms of the number of participating countries, with over 65 in the 2009 cycle, as well as the range of domains assessed, with cross-curricular areas such as problem solving and financial literacy being added to the assessment.

Large-Scale Assessments of Adults

In the 1990s, policy interest in the skills of adult workers and citizens led to the first international large-scale assessment focusing on adults ages 16–65. Working with Statistics Canada, ETS conducted the International Adult Literacy Survey (IALS) between 1994 and 1999, with 22 countries participating over three cycles. This assessment focused on prose, document, and quantitative literacy skills¹ and demonstrated the feasibility of conducting a household survey of adult literacy skills in an international context, maintaining comparability across countries and cultures. As such, IALS laid the foundation for subsequent surveys of adult skills and knowledge. The Adult Literacy and Life Skills Survey (ALL), which focused on a somewhat expanded set of adult skills including literacy, numeracy, and analytical problem solving, was conducted between 2003 and 2008 with some 11 countries participating.² The most recent adult survey, the OECD's Programme for the International Assessment of Adult Competencies (PIAAC), was conducting its first cycle in 2012 with 25 countries participating in 33 languages. PIAAC is a significant step forward in that it is the first computer-based household survey of adults, with interviewers taking laptops into people's homes and asking respondents to complete a background questionnaire and cognitive items on the computer. A parallel paper instrument is utilized for adults who are unable or unwilling to use the laptop equipment. For those adults taking the assessment on the computer, electronic reading tasks as well as scenario-based tasks assessing problem solving in technology envi-

¹ For definitions of these three literacy domains see Organisation for Economic Co-operation and Development and Statistics Canada (2000).

² For definitions of the ALL domains and more information about the survey see Statistics Canada and OECD (2005).

ronments complement the more traditional literacy and numeracy tasks that utilize texts, tables and static print-based stimulus material. PIAAC expands large-scale assessments by utilizing technology to administer the survey and, at the same time, embracing the fact that today's literacy-related tasks often take place in technology-based contexts such as web-based environments, spreadsheets and databases, or electronic mail.

The countries participating in today's student and adult large-scale surveys represent the overwhelming majority of GDP in the world and interest in the data these surveys yield continues to grow. For example, within the context of large-scale assessments, many countries now include special studies focusing on populations of particular interest such as the elderly, immigrants, and incarcerated adults. There is also interest in longitudinal studies as is the case in Canada, which is planning to use PIAAC to measure skills over time. Given that countries are looking more and more toward these assessments for data to drive and inform policy, it is likely that we will see international large-scale surveys continue to expand over time.

The Expanded Range of Large-Scale Assessments

As the aforementioned studies demonstrate, not only have we seen an expansion of who is assessed in terms of the range of participating countries and populations within those countries, but international large-scale assessments are also broadening the horizons in terms of what is being assessed. Earlier studies focused on in-school populations and measured typical academic domains such as mathematics, reading, and science. While these continue to be areas of interest, student assessments have expanded to measure a wider range of competencies and interests, reflecting a growing recognition of the need for lifelong learning as a tool to succeed in rapidly changing economies. Large-scale comparative surveys of adult populations began with a focus on literacy and quantitative skills and have expanded to include numeracy and problem solving in everyday adult contexts. With the growing importance of information technologies, measures of Information and Communication Technology (ICT) literacy skills, digital reading, and problem solving in technology environments have also been included in a number of studies.

The growing interest in assessing technology-related skills and knowledge has led to a growing interest in delivering assessments via computer. As has been mentioned, PIAAC is a household survey delivered on laptops. The call for tenders for PISA 2015 also focuses on moving that assessment more fully towards a computer-based platform. Computer-based assessments are making it possible to include new and innovative item types such as interactive scenario-based items and to collect a broader range of information including timing data and information about the processes test takers engage in when completing assessment tasks. This capability is, in turn, leading to a broadening of the cognitive constructs being measured. Additionally, computer-based assessments make it possible to take advantage of

psychometric advances such as the use of adaptive testing, which allows for more targeted and time efficient measures.

Another significant development in the history of international large-scale assessments has been the growing interest in broadening the information gained from cognitive measures through the use of extensive background questionnaires. Recent student and adult surveys typically include quite extensive background questionnaires. Student questionnaires address a range of topics including general attitudes and interests, day-to-day learning and leisure activities, and educational resources at home.

For adult assessments, questions about job requirements, literacy related activities at home and at work, and social outcomes such as engagement in civic activities have been included. Applying IRT scaling methodologies to these questionnaires has made it possible to create derived scales based on attitude and interest questions as well as on self-reported literacy-related activities and uses of technology. The use of IRT allows us to study differences across participating countries in terms of background characteristics along the same scales and in the same detail as is possible for the cognitive scales. Data from these questionnaires, in conjunction with the cognitive measures, are being used to inform increasingly complex policy questions about the relationships among learning, skills and outcomes.

Both the broadening of the cognitive constructs being addressed in large-scale comparative surveys and the interest in expanded coverage of policy relevant information collected in background questionnaires have driven the need to develop new methodologies for survey design and data analysis. What began as a basic desire to collect descriptive data in the 1960s and 1970s has now expanded to a much broader range of questions of policy interest. There is clearly growing interest on the part of stakeholders from different disciplines to address policy and research questions that are of interest both at the national and the international level.

Evidence-Based Policy Information

It is important to remember that the foundation of international large-scale assessments has always been some call for comparable information about the skills possessed by populations of interest and an understanding of how those skills are related to educational, economic and social outcomes. As such, the development of international large-scale assessments represents a cycle, as shown in Fig. 1.1. The initial work is motivated by policy questions which then drive the development of assessment frameworks and the design of instruments to address those questions. The desire to assess new aspects of existing constructs as well as to include new domains leads to advances in design and methodology that, finally, facilitate the analysis and interpretation of the survey data. This assessment data and the possibility of assessing new constructs as an outcome of more advanced methodologies leads, in turn, to new questions that then form the basis of the next cycle of assessment.