DOCUMENT RESUME

ED 413 218 SE 060 810

O'Sullivan, Christine Y.; Jerry, Laura; Ballator, Nada; AUTHOR

Herr, Fiona

TITLE NAEP 1996 Science State Report for Mississippi. Findings

from the National Assessment of Educational Progress.

INSTITUTION Educational Testing Service, Princeton, NJ.; National

Assessment of Educational Progress, Princeton, NJ.

SPONS AGENCY National Center for Education Statistics (ED), Washington,

DC.

REPORT NO NCES-97-499-MS

PUB DATE 1997-09-00

NOTE 132p.; For overall report, see ED 405 221. For other

> individual reports, see SE 060 786-831. "In collaboration with Audrey Champagne, Peggy Carr, Will Pfeiffenberger, and

Mistilina Sato."

AVAILABLE FROM National Library of Education, Office of Educational

Research and Improvement, U.S. Department of Education, 555

New Jersey Avenue, NW, Washington, DC 20208-5641;

1-800-424-1616 (limit one copy); also on NCES web site:

http://nces.ed.gov/naep/96state

PUB TYPE Information Analyses (070) -- Reports - Research (143)

EDRS PRICE MF01/PC06 Plus Postage.

DESCRIPTORS *Academic Achievement; Academic Standards; Educational

> Change; *Grade 8; Hands on Science; Junior High Schools; *National Competency Tests; Problem Solving; *Science Education; Science Process Skills; Sex Differences; *Standardized Tests; *Student Evaluation; Tables (Data)

IDENTIFIERS

*Mississippi; National Assessment of Educational Progress;

State Science Assessment (NAEP)

ABSTRACT

In 1990, the National Assessment of Educational Progress (NAEP) included a Trial State Assessment (TSA); for the first time in the NAEP's history, voluntary state-by-state assessments were made. The sample was designed to represent the 8th grade public school population in a state or territory. In 1996, 44 states, the District of Columbia, Guam, and the Department of Defense schools, took part in the NAEP state science assessment program. The NAEP 1996 state science assessment was at grade 8 only, although grades 4, 8, and 12 were assessed at the national level as usual. The 1996 state science assessment covered three major fields: earth, physical, and life sciences. In Mississippi, 2,469 students in 103 public schools were assessed. This report describes the science proficiency of Mississippi eighth-graders, compares their overall performance to students in the Southeast region of the United States and the entire United States (using data from the NAEP national assessment), presents the average proficiency for the three major fields, and summarizes the performance of subpopulations (gender, race/ethnicity, parents' educational level, Title I participation, and free/reduced lunch program eligibility). To provide a context for the assessment data, participating students, their science teachers, and principals completed questionnaires which focused on: instructional content (curriculum coverage, amount of homework); delivery of science instruction (availability of resources, type); use of computers in science instruction; educational background of teachers; and conditions facilitating science



+++++ ED413218 Has Multi-page SFR---Level=1 +++++
learning (e.g., hours of television watched, absenteeism). On the NAEP fields
of science scales that range from 0 to 300, Mississippi students had an
average proficiency of 133 compared to 148 throughout the United States. The
average science scale score of males did not differ significantly from that
of females in either Mississippi or the nation. At the eighth grade, White
students in Mississippi had an average science scale score that was higher

than those of Black and Hispanic students. (DDR/NB)



NAEP 1996 SCIENCE

State Report for Mississippi

U.S. DEPARTMENT OF EDUCATION Office of Educational Research and Improvement EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)
This document has been reproduced as received from the person or organization

- originating it.

 ☐ Minor changes have been made to improve reproduction quality.
- Points of view or opinions stated in this document do not necessarily represent official OERI position or policy.



BEST COPY AVAILABLE



U.S. DEPARTMENT OF EDUCATION OFFICE OF EDUCATIONAL RESEARCH AND IMPROVEMENT

What is The Nation's Report Card?

THE NATION'S REPORT CARD, the National Assessment of Educational Progress (NAEP), is the only nationally representative and continuing assessment of what America's students know and can do in various subject areas. Since 1969, assessments have been conducted periodically in reading, mathematics, science, writing, history/geography, and other fields. By making objective information on student performance available to policymakers at the national, state, and local levels, NAEP is an integral part of our nation's evaluation of the condition and progress of education. Only information related to academic achievement is collected under this program. NAEP guarantees the privacy of individual students and their families.

NAEP is a congressionally mandated project of the National Center for Education Statistics, the U.S. Department of Education. The Commissioner of Education Statistics is responsible, by law, for carrying out the NAEP project through competitive awards to qualified organizations. NAEP reports directly to the Commissioner, who is also responsible for providing continuing reviews, including validation studies and solicitation of public comment, on NAEP's conduct and usefulness.

In 1988, Congress established the National Assessment Governing Board (NAGB) to formulate policy guidelines for NAEP. The Board is responsible for selecting the subject areas to be assessed from among those included in the National Education Goals; for setting appropriate student performance levels; for developing assessment objectives and test specifications through a national consensus approach; for designing the assessment methodology; for developing guidelines for reporting and disseminating NAEP results; for developing standards and procedures for interstate, regional, and national comparisons; for determining the appropriateness of test items and ensuring they are free from bias; and for taking actions to improve the form and use of the National Assessment.

The National Assessment Governing Board

Honorable William T. Randall, Chair

Former Commissioner of Education State of Colorado Denver, Colorado

Mary R. Blanton, Vice Chair

Attorney Salisbury, North Carolina

Patsy Cavazos

Principal
W.G. Love Accelerated Elementary School
Houston, Texas

Catherine A. Davidson

Secondary Education Director Central Kitsap School District Silverdale, Washington

Edward Donley

Former Chairman
Air Products & Chemicals, Inc.
Allentown, Pennsylvania

Honorable James Edgar

Member Designate Governor of Illinois Springfield, Illinois

James E. Ellingson

Fourth-Grade Classroom Teacher Probstfield Elementary School Moorhead, Minnesota

Thomas H. Fisher

Director, Student Assessment Services Florida Department of Education Tallahassee, Florida

Michael J. Guerra

Executive Director Secondary Schools Department National Catholic Educational Association Washington, DC

Edward H. Haertel Professor of Education Stanford University Stanford, California Jan B. Loveless

President Loveless and Associates Midland, Michigan

Marilyn McConachie

Vice-Chairperson
Illinois State Board of Education
Northbrook, Illinois

William J. Moloney

Superintendent of Schools Calvert County Public Schools Prince Frederick, Maryland

Honorable Annette Morgan

Former Member Missouri House of Representatives Jefferson City, Missouri

Mark D. Musick

President Southern Regional Education Board Atlanta, Georgia

Mitsugi Nakashima

First Vice-Chairperson Hawaii State Board of Education Honolulu, Hawaii

Michael T. Nettles

Professor of Education & Public Policy University of Michigan Ann Arbor, Michigan and Director Frederick D. Patterson Research Institute United Negro College Fund

Honorable Norma Paulus

Superintendent of Public Instruction Oregon State Department of Education Salem, Oregon

Honorable Roy Romer Governor of Colorado Denver, Colorado Honorable Edgar D. Ross

Judge
Territorial Court of the Virgin Islands
Christiansted, St. Croix
U.S. Virgin Islands

Fannie L. Simmons

Mathematics Coordinator
District 5 of Lexington/Richland County
Ballentine, South Carolina

Adam Urbanski

President Rochester Teachers Association Rochester, New York

Deborah VoltzAssistant Professor

Assistant Professor
Department of Special Education
University of Louisville
Louisville, Kentucky

Marilyn A. Whirry

Twelfth-Grade English Teacher Mira Costa High School Manhattan Beach, California

Dennie Palmer Wolf

Senior Research Associate Harvard Graduate School of Education Cambridge, Massachusetts

Ramon C. Cortines (Ex-Officio)

Acting Assistant Secretary
Office of Educational Research
and Improvement
U.S. Department of Education
Washington, DC

Roy Truby
Executive Director, NAGB
Washington, DC



NATIONAL CENTER FOR EDUCATION STATISTICS

NAEP 1996 SCIENCE STATE REPORT

for

MISSISSIPPI

Christine Y. O'Sullivan
Laura Jerry
Nada Ballator
Fiona Herr

In collaboration with
Audrey Champagne, Peggy Carr,
Will Pfeiffenberger, and Mistilina Sato

September 1997

U.S. Department of Education
Office of Educational Research and Improvement

Prepared by Educational Testing Service under a cooperative agreement with the National Center for Education Statistics.



U.S. Department of Education

Richard W. Riley Secretary

Office of Educational Research and Improvement

Ramon C. Cortines
Acting Assistant Secretary

National Center for Education Statistics

Pascal D. Forgione, Jr. Commissioner

Education Assessment Group

Gary W. Phillips
Associate Commissioner

September 1997

SUGGESTED CITATION

O'Sullivan, C.Y., Jerry, L., Ballator, N., and Herr, H. NAEP 1996 Science State Report for Mississippi, Washington, DC: National Center for Education Statistics, 1997.

FOR MORE INFORMATION

Contact: Arnold A. Goldstein 202-219-1741

For ordering information on this report, write:

National Library of Education
Office of Educational Research and Improvement
U.S. Department of Education
555 New Jersey Avenue, NW
Washington, D.C. 20208-5641

or call 1-800-424-1616 (in the Washington, DC, metropolitan area call 202-219-1651).

This report also is available on the World Wide Web: http://www.ed.gov/NCES/naep

The work upon which this publication is based was performed for the National Center for Education Statistics, Office of Educational Research and Improvement, by Educational Testing Service.

Educational Testing Service is an equal opportunity, affirmative action employer.

Educational Testing Service, ETS, and the ETS logo are registered trademarks of Educational Testing Service.



Table of Contents

HIGHLIGHTS	1
INTRODUCTION	7
What Was Assessed?	
Who Was Assessed?	
Reporting NAEP Science Results	
Interpreting NAEP Results	
How Is This Report Organized?	
Other Reports of NAEP 1996 Science Results	
PART ONE Science Scale Score Results	
Item Maps	20
CHAPTER 1 Science Scale Score Results for Eighth-Grade Students	23
Comparisons Between Mississippi and Other Participating Jurisdi Performance in the NAEP Fields of Science	
CHAPTER 2 Science Scale Score Results for Eighth-Grade Students b Subpopulations	
Gender	32
Race/Ethnicity	
Students' Reports of Parents' Highest Education Level	
Title I Participation	
Free/Reduced-Price Lunch Program Eligibility	
PART TWO Finding a Context for Understanding Students' Scien	ice
Performance in Public Schools	39
CHAPTER 3 School Science Education Policies and Practices	41
Emphasis on Science in the School	4
Resource Availability to Teachers	
Parents as Classroom Aides	
Student Absenteeism	



CHAPTER 4	Science Classroom Practices	49
Curriculu	m Coverage	50
Eighth-G	rade Students' Course Taking	52
Instructio	nal Emphasis	54
Science I	łomework	57
Computer	r Use in Science Instruction	60
CHAPTER 5	Student Performance on Hands-On Science Tasks	65
NAEP H	ands-On Science Tasks	66
Sample (Questions from a Task	67
Instructio	n Related to Scientific Investigation	69
CHAPTER 6	nfluences Beyond School that Facilitate Learning Science	75
Discussin	g Studies at Home	76
Literacy 1	Materials in the Home	77
Televisio	n Viewing Habits	78
Parental S	Support	80
Student M	Mobility	81
Students'	Views About Science	82
APPENDIX A	Reporting NAEP 1996 Science Results	83
APPENDIX B	The NAEP 1996 Science Assessment	99
APPENDIX C	Technical Appendix: The Design, Implementation,	
and Analysis of	f the 1996 State Assessment Program in Science	107
APPENDIX D	Teacher Preparation	121
ACKNOWI FD	CMENTS	127



HIGHLIGHTS

Monitoring the performance of students in subjects such as science is a key concern of the citizens, policy makers, and educators who direct educational reform efforts. The 1996 National Assessment of Educational Progress (NAEP) in science assesses the current level of science performance as a mechanism for informing education reform. This science assessment is the first to be constructed on a new framework, and it is also the first to be given at the state level. This report contains results for public school students at grade 8.

What Is NAEP?

The National Assessment of Educational Progress (NAEP), the "Nation's Report Card," is the only ongoing nationally representative assessment of what America's students know and can do in various academic subjects. Since 1969, NAEP assessments have been conducted with national samples of students in the areas of reading, mathematics, science, writing, and other fields. By making information on student performance available to policy makers, educators, and the general public, NAEP is an integral part of our nation's evaluation of the conditions and progress of education.

NAEP is a congressionally mandated project of the National Center for Education Statistics, U.S. Department of Education. Results are provided only for group performance. NAEP is forbidden by law to report results at an individual or school level.

In 1990 Congress authorized a voluntary state-by-state NAEP assessment. The 1990 Trial State Assessment in mathematics at grade 8 was the first state-level NAEP assessment. Since then, state-level assessments have taken place in 1992 and 1994 in reading (grade 4), in 1992 and 1996 in mathematics (grades 4 and 8), and in 1996 in science (grade 8). In 1996, 44 states, the District of Columbia, Guam, and the Department of Defense Schools took part in the NAEP state assessment program. The NAEP 1996 state science assessment was at grade 8 only, although grades 4, 8, and 12 were assessed at the national level as usual.



NAEP 1996 Science Assessment

The NAEP 1996 science assessment was developed using a new framework. This framework was produced by educators, administrators, assessment experts, and curriculum specialists using a national consensus process. The framework was designed to reflect current practices in science teaching. It called for the use of multiple-choice questions and constructed-response questions that required both short and extended responses. The constructed-response questions served as indicators of students' ability to know and integrate facts and scientific concepts, their ability to reason, and their ability to communicate scientific information. In the 1996 assessment, these constructed-response questions constituted nearly 80 percent of the total student response time. The NAEP 1996 assessment in science also included hands-on tasks that enabled students to demonstrate directly their knowledge and skills related to scientific investigation.

The 1996 science framework was structured according to a matrix that consisted of the three traditional fields of science (earth, physical, and life) crossed with three processes of knowing and doing science (conceptual understanding, scientific investigation, and practical reasoning). A central category encompassing the nature of science and the nature of technology was woven throughout the assessment, as was a themes category representing major ideas or key concepts that transcend scientific disciplines.¹

Students' science performance is summarized on the NAEP science scales, which range from 0 to 300 at each grade. While the scale score ranges are identical for grades 4, 8, and 12, the scales were derived independently at each grade. For example, scale scores on the grade 8 scale cannot imply anything about performance at grade 12 in the national assessment. The science scale is discussed in Appendix C of this report, the NAEP 1996 Science State Report for Mississippi (see C.9). Note that the national average for the combined public and nonpublic school population is 150; the average for public schools only (appropriate for most tables in this report) is 148.

Comparison of Mississippi to the Nation

Table H.1 shows the distribution of science scale scores for eighth-grade students attending public schools in Mississippi, the Southeast region, and the nation in 1996.

• The average science scale score for eighth graders in public schools in Mississippi was 133. This average was lower than that for public school students across the nation (148).²

² Differences reported as significant are statistically different at the 95 percent confidence level. This means that with 95 percent confidence there is a real difference in the average science scale score between the two populations of interest.



More details about the NAEP 1996 science assessment can be found in Appendix B of this report, the NAEP 1996 Science State Report for Mississippi.



TABLE H.1

Distribution of Science Scale Scores for Public School Students at Grade 8

	Average	10th	25th	50th	75th	90th
	Scale Score	Percentile	Percentile	Percentile	Percentile	Percentile
Mississippi	133 (1.4)	91 (3.0)	111 (1.6)	134 (1.5)	155 (1.3)	174 (1.5)
Southeast	141 (1.9)	96 (2.9)	118 (2.7)	.143 (2.1)	165 (1.9)	183 (1.2)
Nation	148 (0.9)	102 (1.6)	126 (1.3)	151 (0.9)	172 (1.1)	191 (1.3)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

Major Findings for Student Subpopulations

The preceding section provided a view of the overall science performance of eighth-grade students in Mississippi. It is also important to examine the average science scale scores of subgroups within the population. Typically, NAEP presents results for demographic subgroups defined by gender, race/ethnicity, and parental education. In addition, in 1996 NAEP collected information on student participation in two federally funded programs: Title I programs and the free/reduced-price lunch component of the National School Lunch Program.

The reader is cautioned against using NAEP results to make simple or causal inferences related to subgroup membership. Differences among groups of students are almost certainly associated with a broad range of socioeconomic and educational factors not discussed in NAEP reports and possibly not addressed by the NAEP assessment program.



Results related to gender and race/ethnicity for public school students are highlighted below. More complete results for the various demographic subgroups examined by the NAEP science assessment can be found in Chapter 2 of this report, the NAEP 1996 Science State Report for Mississippi.

- The average science scale score of males did not differ significantly from that of females in either Mississippi or the nation.
- At the eighth grade, White students in Mississippi had an average science scale score that was higher than those of Black and Hispanic students.

Finding a Context for Understanding Students' Science Performance in Public Schools

The science performance of students in Mississippi may be better understood when viewed in the context of the environment in which students are learning. This educational environment is largely determined by school policies and practices, by characteristics of science instruction in the school, by home support for academics and other home influences, and by students' own views about science. Information about this environment is gathered by means of questionnaires completed by principals and teachers as well as questions answered by students as part of the assessment.

Because NAEP is administered to a sample of students that is representative of all eighth-grade students in Mississippi schools, NAEP results provide a view of the educational practices in Mississippi that may be useful for improving instruction and setting policy. However, despite the richness of context provided by the NAEP results, it is very important to note that NAEP data cannot establish a cause-and-effect relationship between educational environment and students' scores on the NAEP science assessment.



The following results are for public school students:

School Science Education Policies and Practices³

- In Mississippi, the percentage of eighth-grade students attending public schools that reported science was a priority (39 percent) was not significantly different from the percentage of eighth-grade students nationwide (43 percent).
- The percentage of eighth-grade students in Mississippi who attended schools that were expected to follow a district or state curriculum (95 percent) was not significantly different from the national percentage (94 percent).
- In Mississippi, 93 percent of eighth graders attended schools that reported providing instruction in science every day. This percentage did not differ significantly from that of eighth graders across the nation (92 percent).
- Relatively few of the students in Mississippi had teachers who reported receiving all of the resources they needed for classroom instruction (10 percent). This was not significantly different from the corresponding percentage of eighth-grade students nationwide (11 percent).
- In Mississippi, 43 percent of the eighth-grade students were taught by teachers who reported that there was a curriculum specialist available to help or advise them in science. This figure did not differ significantly from that of students across the nation (43 percent).

Science Classroom Practices⁴

- About one fifth of the eighth-grade students in Mississippi had science teachers who reported spending a lot of time on earth science (20 percent), about half reported spending a lot of time on physical science (53 percent), and about one quarter reported spending a lot of time on life science (23 percent).
- Less than one fifth of the students in Mississippi (14 percent) had teachers who reported they planned to place moderate emphasis on the understanding of key science concepts by their students. This percentage was smaller than that of students whose teachers planned heavy emphasis on conceptual understanding (86 percent).



More detailed results related to school policies and practices can be found in Chapter 3 of this report, the NAEP 1996 Science State Report for Mississippi.

⁴ More detailed results related to classroom practices can be found in Chapter 4 of this report, the NAEP 1996 Science State Report for Mississippi.

- In Mississippi, the percentage of eighth-grade students whose teachers reported they planned to give moderate emphasis to developing science problem-solving skills (39 percent) was smaller than that of students whose teachers planned heavy emphasis on this topic (59 percent).
- Teachers of 54 percent of the students in Mississippi reported that they
 planned to place moderate emphasis on knowing how to communicate
 ideas in science effectively, not significantly different from* the
 percentage of students whose teachers reported giving this topic heavy
 emphasis (40 percent).
- In Mississippi, 23 percent of eighth graders reported not spending any time on science homework in a typical week while 31 percent spent one hour or more on their science homework each week.

Scientific Investigations⁵

- Of the eighth-grade students in Mississippi, 81 percent had teachers who reported giving moderate to heavy emphasis on the development of data analysis skills. This percentage was not significantly different from* that of students nationwide (89 percent).
- More than half of the eighth graders in Mississippi had teachers who reported their students performed hands-on activities or investigations in science once a week or more (61 percent).

Influences Beyond School That Facilitate Learning Science⁶

- The percentage of eighth graders in Mississippi who reported watching six or more hours of television a day (29 percent) was greater than the percentage for the nation (17 percent).
- In Mississippi, 41 percent of eighth graders agreed that science is useful for solving everyday problems.

⁶ More detailed results related to influences beyonds school that facilitate learning science can be found in Chapter 6 of this report, the NAEP 1996 Science State Report for Mississippi.



^{*} Although the difference may appear large, recall that "significance" here refers to "statistical significance."

More detailed results related to scientific investigations can be found in Chapter 5 of this report, the NAEP 1996 Science State Report for Mississippi.

INTRODUCTION

Improving education is often seen as an important first step as the United States attempts to remain competitive in an increasingly technical global economy. At the 1996 Governors' Summit in Palisades, New Jersey, the President and the Governors reaffirmed the need to strengthen our schools and strive for world-class standards. Furthermore, in his 1997 State of the Union Address, President Clinton placed education center stage and called for states to commit to national standards that represent what all students must know to succeed in the knowledge-based economy of the twenty-first century.

In 1983, the National Commission on Excellence in Education issued a report entitled A Nation at Risk: The Imperative for Educational Reform that was critical of education in the United States.⁷ Interest in reform was also fueled by the publication of other reports and analyses that pointed out the deficiencies of the educational system and noted how these could be rectified.⁸ Since then, organizations from the public and private sectors have assumed pivotal roles in providing support to state and local educational establishments as they seek to reform their educational systems in areas such as the development of standards, revision of curricula, development of appropriate assessment techniques, and professional development.⁹ In addition to these activities, organizations such as the National Science Teachers Association and the American Association for the Advancement of Science have worked closely with the National Research Council to produce documents that help teachers interpret the National Science Education Standards that were published in 1995.¹⁰ As the new century approaches, commitment to science reform continues.



A Nation at Risk: The Imperative for Educational Reform. (Washington, DC: National Commission on Excellence in Education, 1983).

Educating Americans for the 21st Century: A Report to the American People and the National Science Board. (Washington, DC: National Science Board, Commission on Precollege Education in Mathematics, Science, and Technology, 1983).

Statewide Systemic Initiatives in Science, Mathematics, and Engineering. (Arlington, VA: The National Science Foundation, 1995-1996); Scope, Sequence, and Coordination of Secondary School Science. Volume I: The Content Core; Volume II: Relevant Research. (Washington, DC: National Science Teachers Association, 1992); Benchmarks for Science Literacy. (Washington, DC: Project 2061, American Association for the Advancement of Science, 1993); New Standards Project. (Washington, DC: National Research Council, 1995).

¹⁰ National Science Education Standards. (Washington, DC: National Research Council, 1996).

Monitoring the performance of students in science is a key concern of the state and national policy makers and educators who direct educational reform efforts. To this end, the 1996 National Assessment of Educational Progress (NAEP) is an important source of information on what the nation's students know and can do in science.

What Was Assessed?

The science assessment was crafted to measure the content and skills specified in the science framework for the 1996 NAEP. Two organizing concepts underlie the science framework. First, scientific knowledge should be structured so as to make factual information meaningful. The way in which knowledge is structured should be influenced by the context in which the knowledge is being presented. Second, science performance depends on knowledge of facts, the ability to integrate this knowledge into larger constructs, and the capacity to use the tools, procedures, and reasoning processes of science to develop an increased understanding of the natural world. Thus, the framework called for the NAEP 1996 science assessment to include the following:

- Multiple-choice questions that assess students' knowledge of important facts and concepts and that probe their analytical reasoning skills;
- Constructed-response questions that explore students' abilities to explain, integrate, apply, reason about, plan, design, evaluate, and communicate scientific information; and
- Hands-on tasks that probe students' abilities to use materials to make observations, perform investigations, evaluate experimental results, and apply problem-solving skills.

The core of the science framework is organized along two dimensions. The first dimension divides science into three major fields: earth, physical, and life sciences. The second dimension defines characteristic elements of knowing and doing science: conceptual understanding, scientific investigation, and practical reasoning. Each question in the assessment is categorized as measuring one of the elements of knowing and doing within one of the fields of science (e.g., scientific investigation in the context of earth science). The framework also contains two overarching domains — the nature of science and the organizing themes of science. The nature of science encompasses the historical development of science and technology, the habits of mind that characterize science, and the methods of inquiry and problem solving. It also includes the nature of technology — specifically, design issues involving the application of science to real-world problems and associated trade-offs or compromises. The themes of science include the notions of systems and their application in the scientific disciplines, models and their functioning in the development of scientific understanding, and patterns of change as they are exemplified in natural phenomena. A fuller description of the framework is provided in Appendix B.



Who Was Assessed?

School and Student Characteristics

Table I.1 provides demographic profiles of the eighth-grade students in Mississippi, the Southeast region, and the nation. These profiles are based on data collected from the students and schools participating in the 1996 state and national science assessments at grade 8. As described in Appendix A, the state data and the regional and national data are drawn from separate samples.

To ensure comparability across jurisdictions, NCES has established guidelines for school and student participation rates. Appendix A highlights these guidelines, and jurisdictions failing to meet these guidelines are noted in tables and figures in NAEP reports containing state-by-state results. For jurisdictions failing to meet the initial school participation rate of 70 percent, results are not reported.

Schools and Students Assessed

Table I.2 summarizes participation data for schools and students sampled in Mississippi for the 1996 state assessment program in science.¹¹

In Mississippi, 103 public schools participated in the 1996 eighth-grade science assessment. These numbers include participating substitute schools that were selected to replace some of the nonparticipating schools from the original sample. The weighted school participation rate after substitution in 1996 was 95 percent for public schools, which means that the eighth-grade students in this sample were directly representative of 95 percent of all the eighth-grade public school students in Mississippi.

In each school, a random sample of students was selected to participate in the assessment. In Mississippi in 1996, on the basis of sample estimates, 0 percent of the eighth-grade public school population were classified as students with limited English proficiency (LEP). In addition, 10 percent of eighth graders in public schools had an Individual Education Plan (IEP). An IEP is a plan written for a student who has been determined to be eligible for special education. The IEP typically sets forth goals and objectives for the student and describes a program of activities and/or related services necessary to achieve the goals and objectives.

]



For a detailed discussion of the NCES guidelines for sample participation, see Appendix A of this report or the Technical Report of the NAEP 1996 State Assessment Program in Science. (Washington, DC: National Center for Education Statistics, 1997).



TABLE I.1

Profile of Students in Mississippi, the Southeast Region, and the Nation at Grade 8

Barrania Cuberana		Public		
Demographic Su	bgroups	Percentage		
RACE/ETHNICITY				
Mississippi	White Black Hispanic Asian/Pacific Islander American Indian	50 (2.1) 44 (1.9) 6 (0.6) 0 (0.1) 1 (0.2)		
Southeast	White Black Hispanic Asian/Pacific Islander American Indian	65 (3.8) 26 (3.3) 8 (1.3) 1 (0.3) 1 (0.4)		
Nation	White Black Hispanic Asian/Pacific Islander American Indian	68 (0.4) 15 (0.3) 12 (0.3) 2 (0.3) 2 (0.3)		
PARENTS' EDUCA		· ·		
Mississippi	Did not finish high school Graduated from high school Some education after high school Graduated from college I don't know.	8 (0.6) 24 (0.9) 16 (0.7) 42 (1.3) 10 (0.5)		
Southeast	Did not finish high school Graduated from high school Some education after high school Graduated from college I don't know.	10 (0.6) 26 (2.0) 19 (1.5) 37 (2.1) 8 (0.7)		
Nation	Did not finish high school Graduated from high school Some education after high school Graduated from college I don't know.	7 (0.5) 21 (1.0) 20 (0.7) 42 (1.3) 10 (0.6)		
GENDER				
Mississippi	Male Female	50 (1.1) 50 (1.1)		
Southeast	Male Female	49(0.8) 51(0.8)		
Nation	Male Female	51 (1.2) 49 (1.2)		

(continued on next page)





TABLE I.1 (continued)

Profile of Students in Mississippi, the Southeast Region, and the Nation at Grade 8

Demographic Subgroups		Public
		Percentage
TITLE I		
Mississippi	Participated Did not participate	33 (3.0) 67 (3.0)
Southeast	Participated Did not participate	13 (3.0) 87 (3.0)
Nation	Participated Did not participate	13 (2.3) 87 (2.3)
FREE/REDUCED-I	PRICE LUNCH	
Mississippi	Eligible Not eligible Information not available	52 (1.9) 42 (2.0) 6 (2.5)
Southeast	Eligible Not eligible Information not available	32 (4.3) 41 (7.5) 26 (9.9)
Nation	Eligible Not eligible Information not available	29 (1.6) 51 (3.6) 20 (4.4)

The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). The percentages for Race/Ethnicity may not add to 100 percent because some students categorized themselves as "Other." **** Standard error estimates cannot be accurately determined.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

Schools were permitted to exclude certain students from the assessment, provided that the following criteria were met. To be excluded, a student had to be categorized as LEP or had to have an IEP and (in either case) be judged incapable of participating in the assessment. The intent was to assess all selected students; therefore, all selected students who were capable of participating in the assessment should have been assessed. However, schools were allowed to exclude those students who, in the judgment of school staff, could not meaningfully participate. The NAEP guidelines for inclusion are intended to assure uniformity of inclusion criteria from school to school. Note that some students classified as LEP and some students having an IEP were deemed eligible to participate and were included in the assessment. In Mississippi, the students who were excluded from the assessment because they were categorized as LEP or had an IEP represented 6 percent of the public school population in grade 8.



In Mississippi, 2,469 public school eighth-grade students were assessed in 1996. The weighted student participation rate was 92 percent for public schools. This means that the sample of eighth-grade students who took part in the assessment was directly representative of 92 percent of the eligible public school student population in participating schools in Mississippi (that is, all students from the population represented by the participating schools, minus those students excluded from the assessment). The overall weighted response rate (school rate times student rate) was 87 percent for public schools. This means that the sample of students who participated in the assessment was directly representative of 87 percent of the eligible eighth-grade public school population in Mississippi.

In accordance with standard practice in survey research, the results presented in this report were based on calculations that incorporate adjustments for the nonparticipating schools and students. Hence, the final results derived from the sample provide estimates of the science performance for the full population of eligible public school eighth-grade students in Mississippi. However, in instances where nonparticipation rates are large, these nonparticipation adjustments may not adequately compensate for the missing sample schools and students.

In order to guard against potential nonparticipation bias in published results, the National Center for Education Statistics (NCES) has established minimum participation levels as a condition for the publication of 1996 state assessment program results. NCES also established additional guidelines addressing four ways in which nonparticipation bias could be introduced into a jurisdiction's published results (see Appendix A). In 1996 Mississippi met minimum participation levels for public schools at grade 8. However, Mississippi's nonpublic schools did not participate in the 1996 state assessment. Hence, results are included in this report only for public schools. Mississippi met all other established NCES participation guidelines.

In the analysis of student data and reporting of results, nonresponse weighting adjustments have been made at both the school and student level, with the aim of making the sample of participating students as representative as possible of the entire eligible eighth-grade population. For details of the nonresponse weighting adjustment procedures, see the Technical Report of the NAEP 1996 State Assessment Program in Science.





TABLE 1.2

School and Student Participation at Grade 8 in Mississippi

	Public
SCHOOL PARTICIPATION	
Weighted school participation rate before substitution	89%
Weighted school participation rate after substitution	95%
Number of schools originally sampled	109
Number of schools not eligible	3
Number of schools in original sample participating	96
Number of substitute schools provided	9
Number of substitute schools participating	7
Total number of participating schools	103
STUDENT PARTICIPATION	
Weighted student participation rate after makeups	92%
Number of students selected to participate in the assessment	2,914
Number of students withdrawn from the assessment	133
Percentage of students who were of Limited English Proficiency	0%
Percentage of students excluded from the assessment due to Limited English Proficiency	0%
Percentage of students who had an Individualized Education Plan	10%
Percentage of students excluded from the assessment due to Individualized Education Plan status	6%
Number of students to be assessed	2,693
Number of students assessed	2,469
Overall weighted response rate	87%



Reporting NAEP Science Results

The NAEP Science Scale

The NAEP 1996 science assessment spans the broad field of science in each of the grades assessed. Because of the survey nature of the assessment and the breadth of the domain, each student participating cannot be expected to answer all the questions in the assessment since this would impose an unreasonable burden on students and their schools. Thus, each student was administered a portion of the assessment, and data were combined across students to report on the achievement of eighth graders and on the achievement of subgroups of students (e.g., subgroups defined by gender or parental education).

Student responses to the assessment questions were analyzed to determine the percentage of students responding correctly to each multiple-choice question and the percentage of students achieving each of the score categories for constructed-response questions. Item response theory (IRT) methods were used to produce scales that summarized results for each of the three fields of science (i.e., earth, physical, and life) at each grade level. An overall composite scale also was developed at each of grades 4, 8, and 12 by weighting the separate scales based on the relative importance of each field of science in the NAEP science framework. Results presented in this report are based on this overall composite scale, which ranges from 0 to 300.

The use of separate grade-specific reporting scales for the science assessment is consistent with the National Assessment Governing Board's 1993 policy that future NAEP assessments be developed using within-grade frameworks and that scaling be carried out within grade. Because this science assessment was based on a new framework, and no comparisons with previous NAEP science assessments were possible, a new scale was developed. The ranges of the science scales (from 0 to 300) differ by design from the 0-to-500 reporting scales used in other NAEP subject areas and were chosen to minimize confusion with other common test scales and to discourage inappropriate cross-grade comparisons.

The national average on the science scale is 150, including both public and nonpublic school students. The average for the nation's public school students appears most frequently in this report, and it is slightly lower. (Additional details of the scaling procedures can be found in Appendix C of this report, in the NAEP 1996 Technical Report, and in the Technical Report of the NAEP 1996 State Assessment Program in Science.)



14

Science Achievement Levels

A companion report, being issued by the National Assessment Governing Board, will present the NAEP 1996 science results in terms of achievement levels. As authorized by the NAEP legislation and adopted by the National Assessment Governing Board, the achievement levels are based on the Board's judgments about what are reasonable performance expectations for students on the NAEP 1996 science assessment. The achievement levels for the NAEP 1996 science assessment were adopted on an interim basis, indicating that they may be revised when other information becomes available, such as the fourth- and twelfth-grade results from the Third International Mathematics and Science Study (TIMSS).

Interpreting NAEP Results

This report describes science performance for eighth graders and compares the results for various groups of students within that population — for example, those who have certain demographic characteristics or who responded to a specific background question in a particular way. The report examines the results for individual demographic groups and for individual background questions. It does not include an analysis of the relationships among combinations of these subpopulations or background questions.

Because the percentages of students in these subpopulations and their average science scale scores are based on samples, rather than on the entire population of eighth graders in a jurisdiction, the numbers reported are necessarily estimates. As such, they are subject to a measure of uncertainty, reflected in the standard error of the estimate. When the percentages or average scale scores of certain groups are compared, it is essential to take the standard error into account, rather than to rely solely on observed similarities or differences. Therefore, the comparisons discussed in this report are based on statistical tests that consider both the magnitude of the difference between the means or percentages and the standard errors of those statistics.

The statistical tests determine whether the evidence, based on the data from the groups in the sample, is strong enough to conclude that the averages or percentages are really different for those groups in the population. If the evidence is strong (i.e., the difference is statistically significant), the report describes the group averages or percentages as being different (e.g., one group performed higher than or lower than another group) — regardless of whether the sample averages or sample percentages appear to be about the same or not. If the evidence is not sufficiently strong (i.e., the difference is not significant), the averages or percentages are described as being not significantly different — again, regardless of whether the sample averages or sample percentages appear to be about the same or widely discrepant. Rather than relying on the apparent magnitude of the difference between sample averages or percentages, the reader is cautioned to rely on the results of the statistical tests to determine whether those sample differences are likely to represent actual differences between the groups in the population. The statistical tests and the Bonferroni procedure, which is used when more than two groups are being compared, are discussed in greater detail in Appendix A.



In addition, some of the percentages reported in the text of the report are given qualitative descriptions (e.g., relatively few, about half, etc.). The descriptive phrases used and the rules used to select them are also described in Appendix A.

The tables in the Highlights and in Part 1 (Chapters 1 and 2) show not only the average scale scores for students but also the distribution of their scores at five selected percentiles. The distribution of the scores through these percentiles encourages the reader to consider the performance of the students in the various groupings (whether by state, region, gender, participation in federal programs, etc.) as overlapping ranges of heterogeneous performance, rather than as a simple monolithic average. As an example, consider Table 2.5 which shows that, for the nation, the 75th percentile for students eligible for free or reduced-price lunch is 157 while the average scale score for students who were not eligible for this service is 155. This means that at least 25 percent of the students eligible for free or reduced-price lunch performed above the average for students who were not eligible.

How Is This Report Organized?

The NAEP 1996 Science State Report for Mississippi is a computer-generated report that describes the science performance of eighth-grade students in Mississippi, the Southeast region, and the nation. The system to generate the state reports was developed because reports customized with each jurisdiction's data would otherwise have been impossible to produce in a timely fashion. Because the process is automated, the variables reported were chosen as those most likely to be of interest to most jurisdictions. Unfortunately, this means that some variables of particular interest may not be reported here; however, each jurisdiction will receive all reportable data on CD ROM, and all data will be available on the NCES Web site (http://www.ed.gov/NCES/naep). Also because of the process, the language in the bullets and in parts of the text sometimes seem awkward. It is hoped that understanding the reason for these awkwardnesses will enable the reader to overlook them.

A separate report describes additional eighth-grade science assessment results for the nation and the states, as well as the national results for grades 4 and 12.¹² This State Report consists of four sections:

- This **Introduction** provides background information about what was assessed, who was sampled, and how the results are reported.
- Part One shows the distribution of science scale score results for eighth-grade students in Mississippi, the Southeast region, and the nation.
- Part Two relates eighth-grade public school students' science scale scores to contextual information about school characteristics, instruction, and home support for science in Mississippi, the Southeast region, and the nation. In addition, Chapter 5 discusses student results of the hands-on tasks.

O'Sullivan, C.Y., C.M. Reese, and J. Mazzeo. NAEP 1996 Science Report Card for the Nation and the States. (Washington, DC: National Center for Education Statistics, 1997).



• Several Appendices are presented to support the results discussed in the report:

Appendix A Reporting NAEP 1996 Science Results
Appendix B The NAEP 1996 Science Assessment

Appendix C Technical Appendix Appendix D Teacher Preparation

Other Reports of NAEP 1996 Science Results

Related reports may be of interest to the reader:

- Cross-State Data Compendium for the 1996 Grade 8 Science Assessment
- Technical Report of the NAEP 1996 State Assessment Program in Science
- NAEP 1996 Science Report Card for the Nation and the States

As presently planned, there will be three additional reports appearing in late 1997 and early 1998. One report will contain sample items and examples of student work on these questions. A second report will cover policy and practices in the schools and classrooms in the United States. A third report will cover special components of the NAEP science assessment, including the advanced science assessment and the hands-on exercises.



PART ONE

Science Scale Score Results

The following chapters describe the average science scale scores of eighth-grade students in Mississippi. As described in the Introduction, the NAEP science scale is a composite of the three major fields of science: earth, physical, and life. Student performance is generally reported on this composite scale and so reflects average student scores across the three fields. Student performance may also be summarized on separate NAEP fields of science scales that range from 0 to 300.

This part of the report contains two chapters. Chapter 1 compares the overall science performance of public school students in Mississippi to the nation. (Results for the Southeast region are also presented.) It also contains a U.S. map comparing the average scale scores in Mississippi with other states, and a table showing students' scale score distributions for the three fields of science. Chapter 2 summarizes science performance for subpopulations of public school students as defined by gender, race/ethnicity, parental education, participation in Title I services and programs, and eligibility for the free/reduced-price lunch component of the National School Lunch Program (NSLP).

The NAEP 1996 assessment in science is the first developed using a new framework, described in Appendix B. The scale developed to report results from the 1996 science assessment is a within-grade scale comprised of three fields of science scales. Appendix A describes reporting on the scale, and Appendix C describes the construction of the scale.



Item Maps

Students' performance is summarized on the NAEP science scale which ranges from 0 to 300. Nationally, public school students' scale scores ranged from about 102 for those scoring at the 10th percentile to about 191 for those performing at the 90th percentile. Sample questions are shown in Figure 1.1 illustrating the range of performance on the NAEP science scale for grade 8. Each question is one that is likely to be answered correctly by a student whose score is at or near the given percentile.

To illustrate the range of performance in more detail, questions from the assessment were "mapped" onto a 0 to 300 scale, as in Figure 1.2. The item map is a visual representation of the scale showing selected questions in positions corresponding to their difficulty. The item map shows which questions a student of any particular ability is likely to answer correctly. The position of the question on the scale represents a dividing line. Students who attained scores greater than the score corresponding to the question's difficulty are likely to answer it correctly, while students with scores below that degree of difficulty are less likely to answer it correctly.

More specifically, students who scored below the scale score associated with a particular question had less than a 65 percent probability of earning a given amount of credit on a constructed-response question or less than a 74 percent probability of correctly answering a multiple-choice question. A small proportion of these students — those near but below the question's position on the scale — may be more likely than not to answer the question correctly (between 50 and 65 or 74 percent). Such students are not considered "able" to answer the question, since they have not achieved sufficient consistency in their responses.

This discussion and the item map illustrations refer to eighth-grade students in the national assessment, whose scores may not resemble those of eighth-grade students in Mississippi.



FIGURE 1.1

Sample Questions Likely to Be Answered Correctly by Grade 8 Students At or Near Selected Percentiles

Percentile	Question				
10th	Find typical yearly rainfall from a graph. (104)				
25th	Explain the impact of fish death on an ecosystem. (127)				
50th	Identify the effect of acid rain. (150)				
75th	Understand where earthquakes occur. (172)				
90th	Explain why lightning is seen before thunder is heard. (194)				

The value in parentheses represents the scale score attained by students who had a 65 percent probability of reaching a given level on a constructed-response question (in italic type) or a 74 percent probability of correctly answering a 4-option multiple-choice question (in regular type).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



Figure 1.2 is an item map for grade 8.¹³ Multiple-choice questions are shown in regular type; constructed-response questions are in italic type.¹⁴ An example of how to interpret the item map may be helpful. In this figure, a multiple-choice question involving interpreting a graph maps at the 136 point on the scale. This means that eighth-grade students with science scale scores at or above 136 are likely to answer this question correctly — that is, they have at least a 74 percent chance of doing so.¹⁵ Put slightly differently, this question is answered correctly by at least 74 of every 100 students scoring at or above the 136 scale-score level. Note that this does not mean that students at or above the 136 scale score always answer the question correctly or that students below the 136 scale score always answer it incorrectly.

As another example, consider the constructed-response question that maps at a scale score of 194. This question concerns the differing speeds of light and sound. Scoring of responses to this question allows for partial credit by using a three-level scoring guide. Mapping a question at the 194 scale score indicates that at least 65 percent of the students performing at or above this point achieved a score of 3 ("Complete") on the question. Among students with lower scores, less than 65 percent gave complete responses to the question.



Details on the procedures used to develop the item map are provided in the forthcoming NAEP 1996 Technical Report. The procedures are similar to those used in past NAEP assessments.

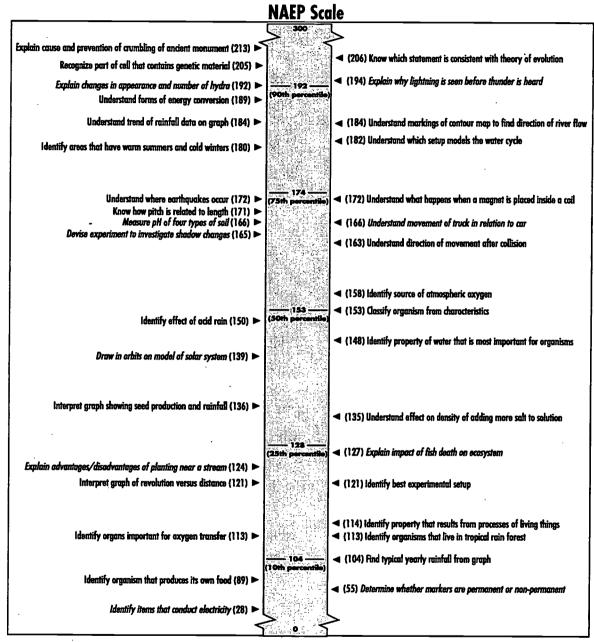
The placement of constructed-response questions is based on (1) the "mapping" of a score of 3 on a 3-point scoring guide for short constructed-response questions and (2) the "mapping" of a score of at least 3 on a 4-point scoring guide and a score of at least 4 on a 5-point scoring guide for extended constructed-response questions.

¹⁵ For constructed-response questions, a criterion of 65 percent was used. For multiple-choice questions, the criterion was 74 percent. The use of a higher criterion for multiple-choice questions reflected the students' ability to "guess" the correct answer from among the alternatives.



FIGURE 1.2 — GRADE 8

Map of Selected Questions on the NAEP Science Scale for Grade 8



NOTE: Position of questions is approximate and an appropriate scale range is disployed for grade 8. Italic type indicates a constructed-response question. Regular type denotes a multiple-choice question.

Each grade 8 science question was mapped onto the NAEP 0-to-300 science scale. The position of the question on the scale represents the scale score attained by students who had a 65 percent probability of reaching a given score level on a constructed-response question or a 74 percent probability of correctly answering a 4-option multiple-choice question. Only selected questions are presented. Percentiles of scale score distribution are referenced on the map.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



CHAPTER 1

Science Scale Score Results for Eighth-Grade Students

To remain competitive in the global economy, a technologically and scientifically literate citizenry is required. As a result, reform in science and mathematics education in the United States has gained increasing attention. The 1983 publication A Nation At Risk: The Imperative for Educational Reform called for overall reform of the United States educational system, with heavy emphasis placed on mathematics and science.¹⁶ The National Goals Panel was convened in 1989 to further focus attention on education reform. In 1991 the National Science Foundation's Statewide Systemic Initiative began awarding grants to support state reform in K-12 mathematics and science instruction.¹⁷ During the 1990s many states have been developing standards for science curriculum, teaching, and assessment using guidance from reform efforts such as the American Association for the Advancement of Science's Project 2061, the National Science Teachers Association's Scope, Sequence, and Coordination of High School Science, and the recently published National Research Council's National Science Education Standards. 18 A reaffirmation of the goal for world-class standards in education was made at the 1996 Governors' Summit in Palisades, NJ. All these efforts address ways to produce innovative science curricula aimed at improving national scientific literacy. As a means of informing the progress of such reform, the U.S. Department of Education supports programs geared toward assessing the current level of science knowledge and skills including the Third International Mathematics and Science Study (TIMSS), 19 administered in 1995, and the 1996 National Assessment of Educational Progress (NAEP) in science.



¹⁶ A Nation at Risk: The Imperative for Educational Reform. (Washington, DC: National Commission on Excellence in Education, 1983).

¹⁷ Statewide Systemic Initiative. (Washington, DC: National Science Foundation, 1990).

¹⁸ Science for All Americans: A Project 2061 Report on Literacy Goals in Science, Mathematics and Technology.
(Washington, DC: American Association for the Advancement of Science, 1989); Scope, Sequence, and Coordination of High School Science. (Washington, DC: National Science Teachers Association, 1995); National Science Education Standards. (Washington DC: National Research Council, 1996).

¹⁹ The Third International Mathematics and Science Study was conducted in 1994 in the Southern Hemisphere and in 1995 in the Northern Hemisphere.

The NAEP 1996 state science assessment at grade 8 was the first time science has been assessed at the state level. It continues the state-level component begun in 1990 with the NAEP Trial State Assessment (TSA). The NAEP 1996 assessment in science had 47 participating jurisdictions.²⁰ Results for 46 jurisdictions were reported for the science assessment.²¹

The science framework for the 1996 National Assessment of Educational Progress²² was developed through a consensus process involving educators, policy makers, business people, assessment experts and curriculum specialists. The 1996 NAEP science assessment included multiple-choice questions, constructed-response exercises, and (for the first time) hands-on tasks. Because the 1996 assessment was based on an essentially new framework, it is not possible to compare results from the 1996 assessment with those from the previous NAEP science assessment in 1990.

Table 1.1 shows the distribution of science scale scores for eighth-grade students attending public schools in Mississippi, the Southeast region, and the nation.

• The average science scale score for eighth-grade public school students in Mississippi was 133. This average was lower than that for public school students across the nation (148).²³



TABLE 1.1

Distribution of Science Scale Scores for Public School Students

	Average Scale Score	10th Percentile	25th Percentile	50th Percentile	75th Percentile	90th Percentile
Mississippi	133 (1.4)	91 (3.0)	111 (1.6)	134 (1.5)	155 (1.3)	174 (1.5)
Southeast	141 (1.9)	96 (2.9)	118 (2.7)	143 (2.1)	165 (1.9)	183 (1.2)
Nation	148 (0.9)	102 (1.6)	126 (1.3)	151 (0.9)	172 (1.1)	191 (1.3)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science

Assessment.



Jurisdiction refers to states, territories, the District of Columbia, and the Department of Defense Education Activities (DoDEA) domestic and international schools. The DoDEA schools also made special arrangements to assess their fourth-grade students in science.

²¹ One jurisdiction did not meet minimum participation levels for public or nonpublic schools and did not have any results reported.

²² Science Framework for the 1996 National Assessment of Educational Progress. (Washington, DC: National Assessment Governing Board, 1993).

Differences reported as significant are statistically different at the 95 percent confidence level. This means that with 95 percent confidence there is a real difference in the average science scale score between the two populations of interest.

Comparisons Between Mississippi and Other Participating Jurisdictions

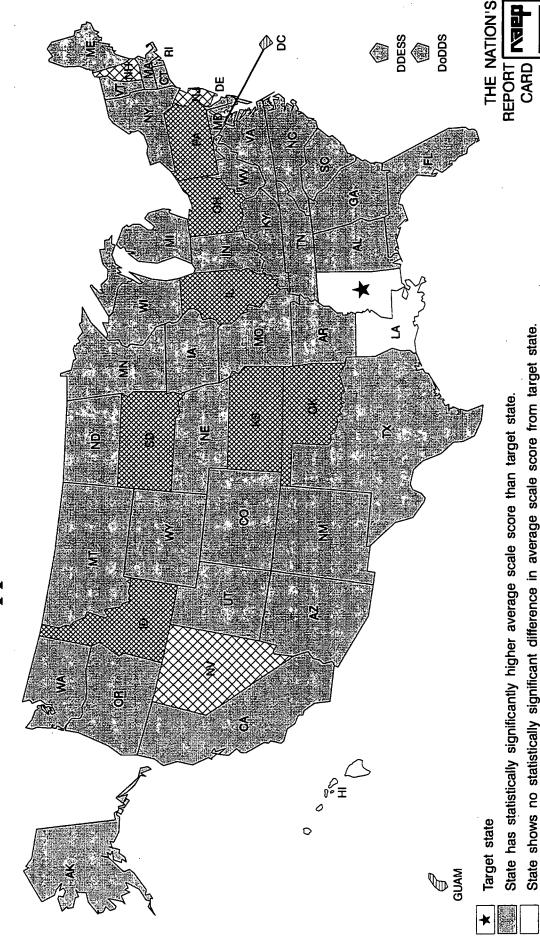
The map on the following page shows how the average science scale score for eighth-grade public school students in Mississippi compares with those of other jurisdictions participating in the NAEP 1996 science assessment. The different shadings on the map indicate whether or not the average scale scores of public school students in the other jurisdictions were statistically different from that of public school students in Mississippi ("Target State"). States with horizontal lines have a significantly lower average science scale score than Mississippi while states with gray shading have a significantly higher average scale score. Unshaded states have average scale scores that did not differ significantly from the average for Mississippi. States with large crosshatching did not meet minimum participation rate guidelines established by NCES for the NAEP assessments. A description of the statistical procedures used to produce this map is contained in Appendix A.



The NAEP 1996 State Assessment

Comparisons of Overall Science Scale Scores at Grade 8

Mississippi Public School Students





BEST COPY AVAILABLE

State has statistically significantly lower average scale score than target state.

State did not meet minimum participation rate guidelines.

State did not participate.

State Assessment

Performance in the NAEP Fields of Science

The core of the science framework is organized along two dimensions. The first divides science into three major fields: earth, physical, and life. The second dimension defines characteristic elements of knowing and doing science: conceptual understanding, scientific investigation, and practical reasoning. Each question is categorized as measuring one of the elements of knowing and doing within one of the fields of science.

Table 1.2 shows the distribution of scale scores for each of the three fields of science for Mississippi, the Southeast region, and the nation. Appendix B describes the three fields of science in more detail, and Appendix C contains a discussion of the scaling procedures used to develop the three fields of science scales and the composite NAEP science scale.

 Students in Mississippi performed lower than students nationwide in the physical science, earth science, and life science fields described in the science framework.



TABLE 1.2

Distribution of Science Scale Scores for Public School Students by Fields of Science

	Average Scale Score	10th Percentile	25th Percentile	50th Percentile	75th Percentile	90th Percentile
Physical Science	<u> </u>					_
Mississippi	132 (1.6)	87 (2.9)	109 (1.8)	133 (1.0)	156 (1.5)	175 (1.8)
Southeast	141 (2.2)	94 (2.2)	117 (2.9)	143 (3.0)	166 (2.3)	184 (2.1)
Nation	149 (1.0)	101 (2.0)	126 (1.3)	151 (1.2)	173 (1.2)	192 (1.6)
Earth Science	1 ' 1					
Mississippi	134 (1.6)	88 (2.3)	110 (1.3)	135 (2.0)	158 (2.2)	178 (2.2)
Southeast	142 (2.2)	95 (3.3)	118 (3.7)	143 (2.3)	167 (1.4)	186 (2.1)
Nation	149 (1.0)	101 (1.9)	126 (1.5)	150 (1.2)	173 (1.3)	192 (1.9)
Life Science	` '	• ,				
Mississippi	133 (1.6)	90 (2.1)	111 (1.4)	134 (1.4)	156 (1.9)	175 (1.7)
Southeast	141 (2.0)	94 (3.3)	117 (3.8)	144 (1.8)	166 (1.1)	183 (1.0)
Nation	148 (1.1)	100 (2.2)	126 (1.3)	151 (1.0)	173 (1.1)	191 (1.7)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



CHAPTER 2

Science Scale Score Results for Eighth-Grade Students by Subpopulations

The previous chapter provided a view of the overall science performance of eighth-grade students in Mississippi and the nation. It is also important to examine the average performance of subgroups since past NAEP assessments in science, as well as in other academic subjects, have shown substantial differences among groups defined by gender, racial/ethnic background, parental education, and other demographic characteristics.²⁴ A key contribution of NAEP to the ongoing conversations concerning education reform is the ability to monitor the performance of subgroups of students in academic achievement.

The NAEP 1996 state assessment in science provides performance information for subgroups of eighth graders in Mississippi, the Southeast region, and the nation. In addition to the more typical demographic subgroups defined by gender, race/ethnicity, and parental education, the 1996 assessment also collected information on two federally funded programs — student participation in Title I programs and services, and student eligibility for the free/reduced-price school lunch program.

²⁴ Jones, L.R., I.V.S. Mullis, S.A. Raizen, I.R. Weiss, and E.A. Weston. The 1990 Science Report Card: NAEP's Assessment of Fourth, Eighth, and Twelfth Graders. (Washington, DC: National Center for Education Statistics, 1992); Campbell, J.R., C.M. Reese, C. O'Sullivan, and J.A. Dossey. NAEP 1994 Trends in Academic Progress. (Washington, DC: National Center for Education Statistics, 1996).



A description of the subgroups and how they are defined is presented in Appendix A. The reader is cautioned against making simple or causal inferences related to the performance of various subgroups of students or about the effectiveness of Title I programs. Average performance differences between two groups of students may in part be due to socioeconomic or other factors. For example, differences observed among racial/ethnic subgroups are almost certainly associated with a broad range of socioeconomic and educational factors not discussed in this report and possibly not addressed by the NAEP assessment program. Similarly, differences in performance between students eligible for Title I programs and those not eligible does not account for the initial performance level of the students prior to placement in Title I programs or differences in course content and emphasis between the two groups.

Gender

Previous NAEP results for science have shown a significant difference in the average scale scores of male and female eighth graders, with males having consistently higher scale scores.²⁵ As shown in Table 2.1, the NAEP 1996 state science assessment results for eighth graders in Mississippi are not consistent with those general findings.

The average science scale score of males did not differ significantly from that of females in either Mississippi or the nation.

THE NATION'S TABLE 2.1	TABLE 2.1
1996 State Assessment	Distribution of Science Scale Scores for Public School Students by Gender

	Average Scale Score	10th Percentile	25th Percentile	50th Percentile	75th Percentile	90th Percentile
Male						
Mississippi	134 (1.8)	90 (2.7)	112 (2.2)	135 (2.4)	157 (2.0)	175 (1.2)
Southeast	142 (2.1)	96 (3.0)	118 (4.6)	144 (3.8)	167 (2.3)	185 (1.7)
Nation	149 (1.1)	101 (1.8)	126 (2.0)	153 (1.1)	174 (1.2)	192 (1.2)
Female						
Mississippi	132 (1.3)	92 (1.9)	110 (2.4)	133 (2.6)	154 (1.2)	172 (2.3)
Southeast	140 (1.9)	96 (3.6)	118 (2.7)	143 (1.8)	163 (1.9)	182 (1.8)
Nation	148 (1.2)	103 (1.3)	127 (1.4)	150 (1.3)	170 (1.7)	189 (3.4)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

²⁵ Campbell, J.R., K.E. Voelkl, and P.L. Donahue. NAEP 1996 Trends in Academic Progress. (Washington, DC: National Center for Education Statistics, 1997); Jones, L.R., I.V.S. Mullis, S.A. Raizen, I.R. Weiss, and E.A. Weston. The 1990 Science Report Card: NAEP's Assessment of Fourth, Eighth, and Twelfth Graders. (Washington, DC: National Center for Education Statistics, 1992).



Race/Ethnicity

As part of the background questions administered with the NAEP 1996 science assessment, students were asked to identify the racial/ethnic subgroup that best describes them. The five mutually exclusive categories were White, Black, Hispanic, Asian or Pacific Islander, and American Indian or Alaskan Native.

Findings from previous NAEP science assessments have shown that racial/ethnic differences exist in science performance. However, when interpreting differences in subgroup performance, confounding factors related to socioeconomic status, home environment, and educational opportunities available to students need to be considered. The distribution of eighth-grade science scale scores for Mississippi, the Southeast region, and the nation by race/ethnicity are shown in Table 2.2.28

 White students in Mississippi demonstrated an average science scale score that was higher than those of Black and Hispanic students.



TABLE 2.2

Distribution of Science Scale Scores for Public School Students by Race/Ethnicity

	Average Scale Score	10th Percentile	25th Percentile	50th Percentile	75th Percentile	90th Percentile
Mississippi	149 (1.2)	113 (2.8)	131 (1.0)	150 (1.7)	167 (1.3)	183 (2.1)
Southeast	153 (1.2)	114 (3.5)	135 (1.2)	155 (1.5)	174 (1.7)	188 (2.4)
Nation	159 (1.1)	120 (1.3)	140 (1.2)	160 (1.2)	179 (1.2)	196 (1.8)
Black						
Mississippi	119 (1.4)	83 (1.1)	101 (2.1)	119 (1.4)	137 (1.7)	154 (1.6)
Southeast	116 (1.8)	79 (2.6)	96 (2.0)	116 (1.4)	136 (2.0)	153 (3.5)
Nation	120 (1.2)	81 (1.8)	99 (1.1)	120 (1.1)	140 (1.6)	158 (1.8)
Hispanic	. 					
Mississippi	105 (3.8)	64 (12.7)	81 (5.3)	103 (8.9)	126 (2.5)	146 (6.8)
Southeast	126 (4.2)	83 (9.7)	105 (3.0)	125 (8.7)	149 (3.1)	167 (4.3)
Nation	127 (1.8)	83 (3.3)	104 (2.6)	129 (1.6)	152 (2.7)	170 (2.8)

The NAEP science scale ranges from 0 to 300. Results are reported for racial/ethnic subgroups meeting established sample size requirements (see Appendix A). The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



²⁶ Campbell, J.R., K.E. Voelkl, and P.L. Donahue. NAEP 1996 Trends in Academic Progress. (Washington, DC: National Center for Education Statistics, 1997); Jones, L.R., I.V.S. Mullis, S.A. Raizen, I.R. Weiss, and E.A. Weston. The 1990 Science Report Card: NAEP's Assessment of Fourth, Eighth, and Twelfth Graders. (Washington, DC: National Center for Education Statistics, 1992).

²⁷ McKenzie, F.D. "Educational Strategies for the 1990s." The State of Black America 1991. (New York: National Urban League, 1991).

²⁸ Results are reported for racial/ethnic subgroups meeting established sample size requirements (see Appendix A).

Students' Reports of Parents' Highest Education Level

Students were asked to indicate the highest level of education completed by each parent. Four levels of education were identified: did not finish high school, graduated from high school, some education after high school, and graduated from college. A choice of "I don't know" was also available. For this analysis, the highest education level reported for either parent was used.

In general, results show that increasing parental education is associated with increases in student performance. In reviewing these results, it is important to note that, nationally, approximately 10 percent of eighth graders did not know the level of education that either of their parents had completed. For public school students in Mississippi, this percentage was 10 percent. Despite the fact that some research has questioned the accuracy of student-reported data from similar groups of students,²⁹ past NAEP assessments in science, as well as other subject areas, have found that student-reported level of parental education exhibits a consistent positive relationship with student performance on the assessments.³⁰ Other research has corroborated NAEP findings.³¹

Table 2.3 shows the results for eighth-grade public school students reporting that neither parent graduated from high school, at least one parent graduated from high school, at least one parent received some education after high school, at least one parent graduated from college, or that they did not know their parents' highest education level. The following pertains to those students who reported knowing the educational level of one or both parents.

• The average science scale score of students in Mississippi who reported that neither parent graduated from high school did not differ significantly from that of students who reported that at least one parent graduated from high school but was lower than that of students who reported that at least one parent received some education after high school or at least one parent graduated from college.

³¹ National Education Longitudinal Study. National Education Longitudinal Study of 1988: Base Year Student Survey. (Washington, DC: National Center for Education Statistics, 1995).



²⁹ Looker, E.D. "Accuracy of Proxy Reports of Parental Status Characteristics." Sociology of Education, 62(4), pp. 257-276, 1989.

³⁰ Jones, L.R., I.V.S. Mullis, S.A. Raizen, I.R. Weiss, and E.A. Weston. The 1990 Science Report Card: NAEP's Assessment of Fourth, Eighth, and Twelfth Graders. (Washington, DC: National Center for Education Statistics, 1992); Campbell, J.R., K.E. Voelkl, and P.L. Donahue. NAEP 1996 Trends in Academic Progress. (Washington, DC: National Center for Education Statistics, 1997); Reese, C.M., K.E. Miller, J. Mazzeo, and J.A. Dossey. NAEP 1996 Mathematics Report Card. (Washington, DC: National Center for Education Statistics, 1997).



TABLE 2.3

Distribution of Science Scale Scores by Public School Students' Reports of Parents' Highest Education Level

	Average Scale Score	10th Percentile	25th Percentile	50th Percentile	75th Percentile	90th Percentile
Did not finish high school Mississippi	125 (2.5)	91 (7.2)	106 (6.0)	125 (2.4)	145 (3.8)	160 (4.8)
Southeast Nation	133 (2.6) 131 (2.0)	98 (7.5) 86 (3.0)	113 (7.2) 108 (2.6)	135 (3.2) 134 (4.0)	151 (3.5) 153 (5.6)	166 (8.0) 170 (3.7)
Graduated from high school Mississippi Southeast Nation	126 (1.9) 134 (2.9) 140 (1.5)	85 (3.1) 89 (5.9) 98 (2.0)	105 (4.3) 111 (5.8) 119 (2.1)	128 (2.5) 136 (1.9) 142 (1.6)	148 (2.6) 158 (2.7) 163 (1.4)	165 (3.8) 177 (1.5) 181 (1.2)
Some education after HS Mississippi Southeast Nation	142 (1.8) 147 (2.3) 155 (1.2)	105 (5.4) 105 (4.9) 113 (1.0)	124 (2.9) 127 (3.3) 137 (1.5)	144 (2.7) 149 (5.1) 158 (2.8)	162 (2.5) 169 (2.3) 176 (2.2)	178 (3.2) 184 (1.7) 191 (1.4)
Graduated from college Mississippi Southeast Nation	138 (1.9) 150 (2.0) 157 (1.3)	95 (1.8) 104 (3.9) 112 (2.1)	115 (2.2) 127 (3.9) 137 (1.0)	139 (1.8) 154 (2.5) 160 (1.4)	162 (2.5) 173 (2.6) 180 (1.5)	181 (3.0) 189 (4.0) 198 (1.3)
I don't know. Mississippi Southeast Nation	119 (2.6) 124 (3.1) 133 (2.6)	79 (2.7) 83 (10.9) 88 (3.6)	98 (5.5) 101 (2.5) 109 (3.6)	119 (5.3) 122 (2.3) 134 (6.5)	142 (6.6) 147 (6.7) 157 (3.8)	161 (3.6) 167 (7.1) 174 (4.4)

The NAEP science scale ranges from 0 to 300. Results are reported for parental education subgroups meeting established sample size requirements (see Appendix A). The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

BEST COPY AVAILABLE



Title I Participation

The Improving America's Schools Act of 1994 (P.L. 103-382) reauthorized the Elementary and Secondary Education Act of 1965 (ESEA). Title I Part A of the ESEA provides financial assistance to local educational agencies to meet the educational needs of children who are failing or most at risk of failing.³² Title I programs are designed to help disadvantaged students meet challenging academic performance standards. Through Title I, schools are assisted in improving teaching and learning and in providing students with opportunities to acquire the knowledge and skills outlined in their state's content and performance standards. For high poverty Title I schools, all children in the school may benefit through participation in schoolwide programs. Title I funding supports state and local education reform efforts and promotes coordinating of resources to improve education for all students.

NAEP first collected student-level information on participation in Title I programs in 1994. The NAEP program will continue to monitor the performance of Title I program participants in future assessments. The Title I information collected by NAEP refers to current participation in Title I services. Students who participated in such services in the past but do not currently receive services are not identified as Title I participants. Differences between students who receive Title I services and those who do not should not be viewed as an evaluation of Title I programs. Typically, Title I services are intended for students who score poorly on assessments. To properly evaluate Title I programs, the performance of students participating in such programs must be monitored over time and their progress must be assessed.³³

Table 2.4 presents results for eighth-grade students by Title I participation.

- For students receiving Title I services, the average science scale score of students in Mississippi (120) was not significantly different from* that of students nationwide (127). The average scale score of Mississippi students who were not receiving Title I services (139) was lower than that of their national counterparts (152).
- The average scale score of Mississippi students who were receiving Title
 I services was lower than that of students who were not.

³³ For a study of mathematics performance of Title I students in 1991-1992, see U.S. Department of Education, PROSPECTS: The Congressionally Mandated Study of Educational Growth and Opportunity, Interim Report: Language Minority and Limited English Proficient Students. (Washington, DC: U.S. Department of Education, 1995).



^{*} Although the difference may appear large, recall that "significance" here refers to "statistical significance."

³² U.S. Department of Education, Office of Elementary and Secondary Compensatory Education Programs. Improving Basic Programs Operated by Local Education Agencies. (Washington, DC: U.S. Department of Education, 1996).



TABLE 2.4

Distribution of Science Scale Scores for Public School Students by Title I Participation

	Average Scale Score	10th Percentile	25th Percentile	50th Percentile	75th Percentile	90th Percentile
Participating						
Mississippi	120 (2.1)	82 (2.3)	100 (2.3)	120 (5.7)	140 (2.5)	157 (2.5)
Southeast	113 (3.3)!	74 (4.2)!	91 (3.3)!	112 (3.1)!	134 (5.5)!	156 (6.5)!
Nation	127 (4.9)	82 (4.1)	102 (5.1)	126 (5.5)	152 (6.2)	170 (7.4)
Not participating						
Mississippi	139 (1.6)	98 (3.0)	118 (2.5)	141 (1.3)	161 (1.7)	178 (1.5)
Southeast	145 (2.0)	103 (1.9)	124 (3.0)	147 (1.7)	168 (1.3)	185 (1.2)
Nation	152 (1.2)	107 (1.8)	131 (1.6)	154 (1.3)	174 (1.3)	192 (2.3)

The NAEP science scale ranges from 0 to 300. Results are reported for students participating in Title I programs only if established sample size requirements are met (see Appendix A). The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

Free/Reduced-Price Lunch Program Eligibility

The free/reduced-price lunch component of the National School Lunch Program (NSLP), offered through the U.S. Department of Agriculture (USDA), is designed to ensure that children near or below the poverty line receive nourishing meals.³⁴ Eligibility for free or reduced-price meals is determined through the USDA's Income Eligibility Guidelines; it is included in this report as an indicator of poverty. The program is available to public schools, nonprofit private schools, and residential child care institutions.

NAEP first collected information on student-level eligibility for the federally funded NSLP in 1996. The NAEP program will continue to monitor the performance of these students in future assessments.



³⁴ U.S. General Services Administration. Catalog of Federal Domestic Assistance. (Washington, DC: Executive Office of the President, Office of Management and Budget, 1995).

Table 2.5 shows the results for eighth graders based on their participation in this program.

- For students who were eligible for free or reduced-price lunch, the average science scale score of students in Mississippi (121) was lower than that of students nationwide (133). Similarly, the average scale score of students who were not eligible for this service was lower for Mississippi (148) than for the nation (155).
- The average scale score of Mississippi students who were eligible for free or reduced-price lunch was lower than that of students who were not.



TABLE 2.5

Distribution of Science Scale Scores for Public School Students by Free/Reduced-Price Lunch Eligibility

	Average Scale Score	10th Percentile	25th Percentile	50th Percentile	75th Percentile	90th Percentile
Eligible						
Mississippi	121 (1.5)	82 (2.2)	101 (2.7)	121 (1.6)	141 (1.9)	160 (1.6)
Southeast	122 (2.0)	82 (1.9)	100 (3.7)	121 (1.9)	145 (3.2)	161 (1.9)
Nation	133 (1.7)	87 (3.0)	108 (2.0)	133 (2.0)	157 (1.7)	176 (2.5)
Not eligible	1					
Mississippi	148 (1.5)	110 (2.8)	129 (2.7)	149 (2.2)	167 (1.3)	183 (1.5)
Southeast	150 (1.6)!	108 (4.2)!	130 (1.5)!	151 (2.7)!	171 (1.8)!	187 (2.4)!
Nation	155 (1.3)	114 (2.6)	136 (1.4)	157 (1.7)	176 (1.2)	194 (2.8)
Information not available	1					
Mississippi	134 (5.6)!	95 (9.6)!	115 (7.6)!	136 (6.9)!	151 (7.6)!	169 (8.9)!
Southeast	152 (3.4)!	111 (7.1)!	132 (4.1)!	155 (2.4)!	174 (3.3)!	186 (1.5)!
Nation	154 (3.6)!	109 (6.2) !	134 (5.0)!	157 (2.8)!	178 (2.9)!	196 (4.7)!

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



PART TWO

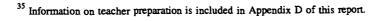
Finding a Context for Understanding Students' Science Performance in Public Schools

The science performance of public school students in Mississippi can be better understood when viewed in the context of the environment in which the students are learning. This educational environment is largely determined by school characteristics, by characteristics of science instruction in the school, by home support for academics and other home influences, and by the students' own views about science. NAEP gathers information about this environment by means of the questionnaires administered to principals, teachers, and students.

Because NAEP is administered to a sample of students that is representative of the eighth-grade student population in the schools of Mississippi, NAEP results provide a view of the educational practices in Mississippi, useful for improving instruction and setting policy. However, despite the richness of the NAEP results, it is very important to note that NAEP data cannot establish a cause-and-effect relationship between educational environment and student scores on the NAEP science assessment.

The variables contained in Part Two are from the school characteristics and policies questionnaire, teacher questionnaires, and student background questionnaires. Part Two consists of four chapters: Chapter 3 discusses school characteristics related to science instruction;³⁵ Chapter 4 describes classroom practices related to science instruction, including curriculum, instructional emphases, coursework, and computer use; Chapter 5 describes portions of a hands-on task and explores student exposure to these experiences; and Chapter 6 covers some potential influences from the home and from the students' own views about science.

To provide additional information, the bullets below sometimes contain combined results from one or more categories (i.e., collapsed categories). When this is the case, the summed numbers reported in the bullets may be slightly different from the sums of the rounded numbers presented in the tables for each of the categories.





CHAPTER 3

School Science Education Policies and Practices

School programs and conditions, instructional practices, and resource availability vary from state to state and even among schools within a locality. The information in this chapter is intended to give insight into those policies or practices that are associated with students' success in science.

The variables reported here reflect information from the questionnaires completed by principals and teachers of the public school students in the NAEP 1996 science assessment. In all cases, analyses are done at the student level. School and teacher-reported results are given in terms of the percentage of students who attend schools or who have teachers reporting particular practices.³⁶

Emphasis on Science in the School

In the school characteristics and policies questionnaire, principals or other head administrators were asked several questions relating to the priority placed on science within their schools. Table 3.1 presents their responses.

- The percentage of eighth-grade students in Mississippi who attended schools with a special focus on science (3 percent) was not significantly different from the national percentage (8 percent).
- The percentage of eighth-grade students in Mississippi attending schools that reported science was a priority (39 percent) was not significantly different from the national percentage (43 percent). The average scale score for students in these schools (133) was lower than that of students in schools nationwide reporting that science was a priority (147).
- The average scale score of students in Mississippi schools that reported that science was a priority (133) was not significantly different from that of students in schools where science was not a priority (135).
- The percentage of eighth-grade students in Mississippi who attended schools that reported having a district or state curriculum that the school was expected to follow (95 percent) was not significantly different from the national percentage (94 percent).



³⁶ Appendix A provides more details on the units of analysis used to derive the results presented in this report.



TABLE 3.1

Public Schools' Reports on Science as a Priority

Mississippi	Southeast	Nation		
Percentage and Average Scale Score				

Is this a school with a special			
focus on science?*			
Yes	3 (1.7)	18 (8.4)	8 (2.7)
	145 (13.2)!	*** (**.*)	137 (5.0)!
Has your school identified science		·	Ì
as a priority in the last two years?			
Yes	39 (5.0)	49 (11.2)	43 (6.8)
	133 (2.4)	140 (3.5)!	147 (3.3)
No	61 (5.0)	51 (11.2)	57 (6.8)
	135 (1.9)	142 (3.2)!	151 (1.7)
Does your district or state have a curriculum in science that your school is expected to follow?*	,	, ,	
Yes	95 (2.0)	100 (****)	94 (2.0)
	134 (1.5)	141 (1.8)	149 (1.0)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). * The response category "No" was inappropriate here because the question permitted several options to be selected; consequently, only "Yes" responses were tallied. ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. **** Standard error estimates cannot be accurately determined. SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

Principals were also asked how often students received science instruction. Schools using block scheduling (i.e., extended periods of instruction on fewer days) were not separately identified. Consequently, students in schools with block scheduling who receive science instruction two or three times weekly may receive as many *hours* of instruction as students under traditional scheduling who receive instruction every day. Table 3.2 shows the following:

- In Mississippi, 93 percent of eighth graders attended schools that reported providing instruction in science every day. This percentage did not differ significantly from that of eighth graders across the nation (92 percent).
- The average scale score for students receiving science instruction every day (134) was lower than that of students nationwide receiving this much instruction (150).





TABLE 3.2

Public Schools' Reports on Time Spent in Science Instruction

How often does a typical	Mississippi	Southeast	Nation
eighth-grade student in your school receive instruction in science?	Percent	age and Average Scal	e Score

Twice a week or less/Not taught	0 (****)	0 (****)	0 (****)
	*** (**.*)	*** (**.*)	*** (**.*)
Three or four times a week	7 (3.2)	1 (****)	8 (2.7)
	129 (4.9)!	*** (**.*)	147 (4.8)!
Every day	93 (3.2)	99 (****)	92 (2.7)
	134 (1.6)	142 (1.8)	150 (1.2)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. **** Standard error estimates cannot be accurately determined.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment

Resource Availability to Teachers

Resources available to teachers and schools vary. Past surveys have shown that teachers' perceptions of the availability of resources (e.g., materials, staff, and time) are variable across the country.³⁷ Previous NAEP assessments in other subject areas have shown an overall positive relationship in most states between teachers' reports of resource availability and their students' performance.³⁸

Availability of Instructional Materials

Teachers often see the lack of resources and materials as a key problem for science instruction. In 1993 a national survey of elementary and secondary school educators reported that deficiencies related to instructional resources were the most serious problems for science instruction in their schools.³⁹ In that survey, schools reported spending a total of \$0.51 per elementary student per year and \$0.88 per middle grade student per year on science supplies, and \$50 per year on science software. (The average price for one piece of software is \$100.)



³⁷ U.S. Department of Education. Schools and Staffing in the United States: A Statistical Profile, 1993-94. (Washington, DC: National Center for Education Statistics, 1996).

³⁸ For example, see Miller, K.E., J.E. Nelson, and M. Naifeh. Cross-State Data Compendium for the NAEP 1994 Grade 4 Reading Assessment. (Washington, DC: National Center for Education Statistics, 1995); National Center for Education Statistics. State-by-State Background Questionnaire Data Appendix: NAEP 1992 Mathematics Assessment, Grades 4 and 8. (Washington, DC: Office of Educational Research and Improvement, 1994).

³⁹ Weiss, I.R. A Profile of Science and Mathematics Education in the United States: 1993. (Chapel Hill, NC: Horizon Research, 1994).

Teachers whose students participated in the NAEP 1996 science assessment were asked to categorize how well their school systems provided them with the classroom instructional materials they needed. The results are shown in Table 3.3.

- Relatively few of the students in Mississippi had teachers who reported receiving all the resources they needed (10 percent). This percentage was not significantly different from that of students across the nation (11 percent).
- The average science scale score of students in Mississippi whose teachers reported receiving all the resources they needed (138) was not significantly different from that of students whose teachers received some or none of the resources they needed (134).



TABLE 3.3

Public School Teachers' Reports on Resource Availability

Which of the following statements is true about how well your school system provides	Mississippi	Southeast	Nation	
you with the instructional materials and other resources you need to teach your class?	Percentage and Average Scale Score			
I get some or none of the resources I need.	47 (4.7)	47 (9.5)	37 (4.1)	
	134 (2.4)	143 (3.6)!	144 (2.0)	
I get most of the resources I need.	43 (4.5)	39 (7.0)	52 (4.1)	
	135 (2.4)	143 (3.1)!	153 (2.1)	
I get all the resources I need.	10 (2.4)	14 (5.2)	11 (3.1)	
	138 (2.6)!	139 (4.9)!	154 (5.4)!	

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



Availability of Curriculum Specialist in the School

Table 3.4 shows the percentages and average scale scores of eighth-grade students in public schools whose teachers indicated they had a curriculum specialist available to help or advise them in science.

• In Mississippi, less than half of the students were taught by teachers who reported that there was a curriculum specialist available to help or advise them in science (43 percent). This figure did not differ significantly from that of students across the nation (43 percent).



TABLE 3.4

Public School Teachers' Reports on Curriculum Specialists

Is there a curriculum specialist	Mississippi	Southeast	Nation	
available to help or advise you in science?	Percentage and Average Scale Score			
Yes	43 (4.5)	58 (8.0)	43 (3.9)	
	134 (2.1)	142 (3.1)	148 (2.7)	
No	57 (4.5)	42 (8.0)	57 (3.9)	
	135 (2.0)	144 (2.6)!	152 (1.5)	

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



Parents as Classroom Aides

When school personnel and parents develop a positive line of communication, they strengthen the learning environment for the students both at school and at home. One of the most frequent reasons cited by school personnel for contacting parents is to request parent volunteer time at school.⁴⁰ The principals of the participating public schools were asked if parents were used as classroom aides. As shown in Table 3.5, principals for eighth graders reported the following:

A small percentage of the students in Mississippi (5 percent) were in schools that reported routinely using parents as aides in classrooms while 55 percent of students in Mississippi attended schools where parents were not used as classroom aides.



TABLE 3.5

Public Schools' Reports on Parents as Aides in Classrooms

Does your school use parents as	Mississippi	Southeast	Nation	
aides in classrooms?	Percentage and Average Scale Score			
No ·	55 (5.3)	44 (9.1)	43 (6.0)	
	133 (1.9)	138 (3.2)!	146 (2.4)	
Yes, occasionally	41 (4.9)	42 (10.7)	46 (6.3)	
	131 (2.5)	138 (3.6)!	150 (2.7)	
Yes, routinely	5 (2.9)	14 (7.8)	11 (3.6)	
	*** (**.*)	**** (**.*)	152 (6.9)!	

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



⁴⁰ U.S. Department of Education. *The Condition of Education 1995*. (Washington, DC: National Center for Education Statistics, 1995). 46



Student Absenteeism

School principals were asked if student absenteeism was a serious, moderate, or minor problem, or not a problem. Table 3.6 shows results for eighth graders based on principals' reports.

- In Mississippi, 41 percent of the eighth-grade public school students attended schools that reported that absenteeism was a moderate to serious problem. This percentage was greater than that for the nation (22 percent).
- The average scale score of students in Mississippi attending schools that reported that absenteeism was not a problem (138) was not significantly different from* that of students in schools where absenteeism was a moderate to serious problem (131).

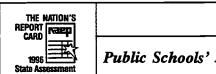


TABLE 3.6

Public Schools' Reports on Student Absenteeism

To what degree is student	Mississippi	Southeast	Nation	
absenteeism a problem in your school?	Percentage and Average Scale Score			
Not a problem	7 (.2.5)	7 (3.9)	28 (4.8)	
	138 (8.4)!	**** (**.*)	156 (3.1)	
Minor	52 (5.1)	56 (6.4)	50 (4.9)	
	135 (2.0)	144 (3.0)	149 (1.5)	
Moderate to serious	41 (4.9)	37 (6.0)	22 (3.7)	
	131 (2.6)	· 138 (2.5)l	140 (3.0)	

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

BEST COPY AVAILABLE



^{*} Although the difference may appear large, recall that "significance" here refers to "statistical significance."

CHAPTER 4

Science Classroom Practices

Science education in the nation's schools has received considerable attention at the national, state, district, school, and classroom levels. In recent years, a number of national and international programs have measured student performance in science. The latest national trend report indicates that although eighth graders' scores have shown recent increases, there is no significant difference in average scores between 1970 and 1996.⁴¹ A recent international study, the Third International Mathematics and Science Study (TIMSS), demonstrated that eighth-grade students' performance in the United States was slightly above average compared with that of students in 40 other countries.⁴²

Using guidance from such programs as the Statewide Systemic Initiative, Project Scope, Sequence, and Coordination, Benchmarks for Science Literacy, and the National Science Education Standards,⁴³ many states are currently involved in re-evaluating their existing standards and developing new frameworks and criteria for science instruction in their state. TIMSS has also pointed out some differences between classroom practices in the United States and in the 40 other participating nations that may guide development of more effective science instruction.⁴⁴ This chapter focuses on curricular and instructional content issues in Mississippi public schools and their relationship to students' science performance.

For some of the issues discussed in this chapter, student- and teacher-reported results for similar questions are presented. In these situations, some discrepancies may exist between student- and teacher-reported percentages. It is not possible to offer conclusive reasons for these discrepancies or to determine whose reports more accurately reflect eighth-grade classroom activities. The results merely present students' and teachers' impressions of the science classroom.



⁴¹ Campbell, J.R., K.E. Voelkl, and P.L. Donahue. NAEP 1996 Trends in Academic Progress. (Washington, DC: National Center for Education Statistics, 1997).

⁴² Beaton, A.E., M.O. Martin, I.V.S. Mullis, E.J. Gonzalez, T.A. Smith, and D.L. Kelly. Science Achievement in the Middle School Years: IEA's Third International Mathematics and Science Study (TIMSS). (Chestnut Hill, MA: TIMSS International Study Center, 1996).

⁴³ National Science Foundation, 1990, Statewide Systemic Initiative, provided grants to further research and initiatives in science reform; Scope, Sequence and Coordination of High School Science. Vol.1. The Content Core: A Guide for Curriculum Developers. (Washington, DC: National Science Teachers Association, 1992); American Association for the Advancement of Science. Benchmarks for Science Literacy. (New York: Oxford University Press, 1993); National Research Council. National Science Education Standards. (Washington, DC: National Academy Press, 1996).

⁴⁴ National Center for Education Statistics. *Pursuing Excellence*. (Washington, DC: U.S. Government Printing Office, 1996).

Curriculum Coverage

The NAEP 1996 science assessment examines three fields of science: earth, physical, and life. In grades 4 and 12, the 1996 NAEP framework emphasized the three fields of science more or less equally; however, the framework specified a heavier emphasis on life science at grade 8, consistent with the increasingly recognized importance of human biology for this age group.⁴⁵ Eighth-grade public school teachers were asked how much time was spent on the three traditional fields of science in their classes and the results are presented in Table 4.1.

- In Mississippi, 20 percent of the eighth-grade public school students had teachers who reported spending a lot of time on earth science. This percentage was smaller than that for the nation (41 percent). Students in Mississippi in classrooms where a lot of time was spent on earth science had an average scale score (131) that was lower than that of similar students nationwide (149).
- In Mississippi, 53 percent of the public school students had teachers who reported spending a lot of time on physical science. This figure was not significantly different from that of their national counterparts (49 percent). The average science scale score in classrooms where physical sciences was covered a lot was lower in Mississippi (134) than nationwide (151).
- In Mississippi, 23 percent of the students had teachers who reported spending a lot of time on life science. This was not significantly different from the percentage nationwide (19 percent). The average scale score for students in these classrooms (135) was lower than that of students across the nation spending a lot of time on life science (147).



52

[.]

⁴⁵ National Research Council. National Science Education Standards. (Washington, DC: National Academy Press, 1996).



TABLE 4.1

Public School Teachers' Reports on Curriculum Coverage

How much time do you spend on	Mississippi	Southeast	Nation
each of the following areas of			
science in this class?	Percent	age and Average Scal	e Score

Earth science	None	6 (1.9)† 136 (3.0)!	6 (2.8)† **** (**.*)	7 (1.8)† 153 (4.4)!
	A little	15 (2.5)† 129 (4.7)	9 (4.6)† *** (**.*)	11 (3.1)† 153 (5.6)!
	Some	59 (4.0)† 138 (1.9)	45 (9.1)† 147 (2.6)i	41 (5.0)† 151 (2.1)
	A lot	20 (3.3)† 131 (2.6)	· 40 (6.4)† 138 (3.4)	41 (5.6)† 149 (2.9)
Physical science	None	0 (****)† **** (**.*)	0 (****)† *** (**.*)	3 (1.2)† 141 (9.5)!
	A little	2 (0.9)† **** (**.*)	9 (2.7)† *** (**.*)	12 (3.6)† 152 (4.4)!
	Some	45 (4.5)† 136 (2.4)	41 (7.9)† 149 (2.9)!	36 (4.9)† 152 (2.8)
	A lot	53 (4.5)† 134 (2.1)	50 (7.8)† 142 (2.5)!	49 (4.9)† 151 (1.8)
Life science	None	7 (2.2)† 136 (2.6)!	12 (5.7)† 151 (2.9)!	17 (5.1)† 155 (5.0)!
	A little	12 (2.1)† 130 (4.5)	23 (7.5)† 144 (2.7)!	22 (4.1)† 152 (3.5)
	Some	59 (4.0)† 136 (2.2)	51 (11.9)† 147 (3.6)!	41 (6.1)† 149 (2.5)
	A lot	23 (3.4)† 135 (2.9)	14 (4.2)† 140 (4.0)!	19 (4.7)† 147 (2.6)!

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. **** Standard error estimates cannot be accurately determined. † Interpret with caution — more than 15% of the respondents did not answer this question.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

BEST COPY AVAILABLE



Eighth-Grade Students' Course Taking

Exposure to science and the opportunity to learn science have a positive effect on the science performance of students.⁴⁶ To investigate whether there is a relationship between science performance of students on the 1996 NAEP assessment and their study of science in school, information on the types of science classes in which eighth-grade students were enrolled and the amount of time spent each week on science instruction was collected. As noted for Table 3.2, in which school principals answered a similar question concerning the frequency of science instruction, students in schools with block scheduling were not identified separately. Consequently, students under block scheduling who receive science instruction two or three times weekly may be receiving as much instruction as students in traditional settings who have science every day.

Based on students' responses shown in Table 4.2:

- In eighth grade, 2 percent of the students in Mississippi reported not taking a science course this year. This did not differ significantly from the national percentage (3 percent).
- In Mississippi, the average scale score for students taking life science (126) was not significantly different from that of students taking physical science (130).
- The average scale score for Mississippi students taking life science (126) was not significantly different from that of students taking earth science (122).
- In Mississippi, 83 percent of the students reported studying science three or more times a week. The average scale score for students who reported studying science three or more times a week in Mississippi (136) was lower than that of students studying at this level nationwide (152).



Council of Chief State School Officers. State Indicators of Science and Mathematics Education. (Washington, DC: CCSSO, 1995).





TABLE 4.2

Public School Students' Reports on Their Science Classes

Mississippi	Southeast	Nation
Percenta	age and Average Sca	ale Score

Which best describes the science			
course you are taking?			
I am not taking science this year.	2 (0.4) *** (**.*)	2 (0.4)	3 (0.9) 120 (3.0)!
Life science	7 (0.9)	11 (1.5)	12 (1.5)
	126 (4.0)	. 122 (2.9)!	133 (3.5)
Physical science	36 (2.8)	27 (4.8)	25 (2.6)
	130 (2.0)	146 (3.0)!	154 (1.6)
Earth science	13 (1.2)	23 (4.2)	23 (3.1)
	122 (2.8)	137 (3.5)!	148 (3.6)
General science	12 (1.0)	16 (3.2)	19 (1.5)
	143 (2.7)	145 (3.7)!	156 (1.7)
Integrated science	30 (2.4)	21 (4.9)	17 (1.8)
	141 (2.0)	153 (2.2)!	156 (1.6)
About how often do you study		1	
science in school?			
Never	3 (0.4)	4 (0.7)	4 (0.5)
	115 (3.5)	116 (5.7)	126 (3.2)
Less than once a week	5 (0.5)	4 (0.4)	4 (0.3)
	117 (4.3)	128 (4.8)	136 (3.0)
1 or 2 times a week	10 (0.8)	. 8 (0.9)	7 (0.8)
	121 (2.6)	133 (4.8)	138 (2.6)
3 or 4 times a week	13 (1.3)	10 (0.8)	13 (1.9)
	133 (2.5)	138 (2.5)	146 (2.2)
Every day	70 (1.7)	75 (1.9)	71 (2.7)
	137 (1.5)	145 (2.1)	153 (1.3)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



Instructional Emphasis

The framework that guided the development of the NAEP 1996 science assessment identified three ways of knowing and doing science — conceptual understanding, scientific investigation, and practical reasoning.⁴⁷ In addition, the science education reform effort has focused heavily on students' ability to communicate their understanding of science to others.⁴⁸ To assess students' opportunities to learn and communicate the knowledge and skills outlined in the framework, teachers were asked about their plans for science instruction during the entire year. Their responses are shown in Table 4.3.

- In Mississippi, the percentage of eighth-grade students whose teachers reported they planned to give moderate emphasis to knowing science facts and terminology (37 percent) was smaller than that of students whose teachers planned heavy emphasis on knowing facts and terminology (62 percent).
- The average scale score of students whose teachers gave moderate emphasis to knowing facts and terminology (137) was not significantly different from that of students whose teachers heavily emphasized this topic (133).
- Less than one fifth of the students in Mississippi (14 percent) had teachers who reported they planned to place moderate emphasis on the understanding of key science concepts by their students. This percentage was smaller than that of students whose teachers planned heavy emphasis on conceptual understanding (86 percent).
- The average scale score of students whose teachers planned moderate emphasis on the understanding of science concepts (130) was not significantly different from that of students whose teachers placed heavy emphasis on this topic (135).
- In Mississippi, the percentage of eighth-grade students whose teachers reported they planned to give moderate emphasis to developing science problem-solving skills (39 percent) was smaller than that of students whose teachers planned heavy emphasis on this topic (59 percent).
- Teachers of 54 percent of the students in Mississippi reported that they
 planned to place moderate emphasis on knowing how to communicate
 ideas in science effectively, not significantly different from* the
 percentage of students whose teachers reported giving this topic heavy
 emphasis (40 percent).

⁴⁸ American Association for the Advancement of Science. Benchmarks for Science Literacy. (New York: Oxford University Press, 1993); National Research Council. National Science Education Standards. (Washington, DC: National Academy Press, 1996).



^{*} Although the difference may appear large, recall that "significance" here refers to "statistical significance."

⁴⁷ Science Framework for the 1996 National Assessment of Educational Progress. (Washington, DC: National Assessment Governing Board, 1993).



TABLE 4.3

Public School Teachers' Reports on Instructional Emphasis

Think about your plans for your science instruction during the entire year. About how much emphasis will you give to the following as an objective for your students?

Mississippi Southeast Nation

Percentage and Average Scale Score

Knowing science facts and terminology			
Little or no emphasis	2 (0.9)	0 (****)	5 (2.3)
	*** (**.*)	*** (**.*)	154 (4.0)!
Moderate emphasis	37 (4.3)	47 (6.8)	57 (3.4)
	137 (2.3)	144 (2.5)!	153 (1.4)
Heavy emphasis	62 (4.3)	52 (6.8)	38 (3.9)
	133 (2.0)	141 (2.9)	145 (2.6)
Understanding key science concepts			
Little or no emphasis	0 (****) *** (**.*)	0 (****) *** (**.*)	0 (****)
Moderate emphasis	14 (2.9)	18 (8.3)	11 (2.4)
	130 (4.1)!	144 (3.3)!	143 (2.4)!
Heavy emphasis	86 (2.9)	82 (8.3)	89 (2.5)
	135 (1.6)	142 (2.4)	151 (1.2)
Developing science problem-solving Skills			
Little or no emphasis	2 (1.0)	1 (****)	3 (.1.6)
	*** (**.*)	*** (**.*)	140 (20.9)!
Moderate emphasis	39 (4.0)	34 (6.0)	28 (3.7)
	132 (2.5)	139 (4.6)!	148 (3.4)
Heavy emphasis	59 (3.9)	65 (6.1)	69 (4.3)
	135 (2.0)	145 (2.2)	152 (1.3)
Knowing how to communicate ideas in science effectively			
Little or no emphasis	6 (1.9)	4 (1.4)	16 (3.3)
	135 (4.5)!	*** (**.*)	151 (2.7)!
Moderate emphasis	54 (4.7)	56 (8.2)	42 (4.3)
	134 (1.9)	142 (2.6)!	149 (2.3)
Heavy emphasis	40 (4.3)	40 (8.3)	42 (4.4)
	134 (2.5)	144 (3.4)!	151 (1.5)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. **** Standard error estimates cannot be accurately determined.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

BEST COPY AVAILABLE



With the explosion of the information age, mainstream news and the Internet afford opportunities to access up-to-date scientific information. Science instruction could benefit by taking advantage of such opportunities. To determine if these opportunities were being explored, eighth-grade teachers and students were asked how often they have classroom discussions about science stories that appear in the news. The results are presented in Table 4.4.

- In Mississippi, 53 percent of eighth-grade students were taught by teachers who reported frequent (once a week or more) classroom discussions of science in the news. A small percentage of the students (7 percent) had teachers who reported never or hardly ever discussing science in the news.
- When students were asked how often they discussed science in the news,
 30 percent reported frequent discussions while 52 percent reported never or hardly ever discussing it.



TABLE 4.4

Public School Teachers' and Students' Reports on Discussions of Science in the News

	Mississippi	ssippi	Southeast		Nation	
How often do your students (do you) discuss science in the news?	Teacher	Student	Teacher	Student	Teacher	Student
		Percenta	ige and Av	erage Sca	le Score	

Never or hardly ever	7 (1.9)	52 (1.4)	0 (0.1)	44 (2.0)	8 (2.6)	44 (1.3)
	125 (5.1)!	131 (1.6)	(**.*)	135 (2.2)	155 (7.6)!	144 (1.2)
Once or twice a month	40 (4.1)	18 (0.9)	56 (9.5)	24 (1.3)	44 (4.9)	22 (1.1)
	135 (2.0)	140 (2.1)	142 (2.5)!	151 (2.5)	150 (2.1)	155 (1.9)
Once or twice a week	43 (4.7)	19 (0.9)	27 (5.0)	22 (1.3)	33 (2.9)	22 (0.9)
	137 (2.2)	138 (2.1)	139 (3.7)	146 (2.5)	149 (2.0)	154 (1.8)
Almost every day	11 (2.9)	11 (0.7)	16 (8.3)	10 (0.8)	16 (4.9)	11 (1.1)
	131 (5.0)!	126 (2.7)	153 (1.9)!	138 (3.7)	153 (3.8)!	147 (2.8)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



Science Homework

Past NAEP science assessments have shown a positive relationship between science homework and performance.⁴⁹ To examine the relationship between homework and science scale scores in Mississippi, the teachers of the assessed students were asked to report the amount of science homework they assigned each week, and students were asked to report the amount of time they spent on science homework each week.

Tables 4.5 and 4.6 show the teachers' and students' responses for eighth-grade public school students in Mississippi. (Students had an additional response choice "I am not taking a science course this year," but no analogous option was available to teachers.) According to the teachers' responses:

- In Mississippi, 2 percent of the eighth graders were not assigned science homework each week. In addition, 84 percent of the students were assigned an hour or more of homework each week.
- The percentage of students in Mississippi whose teachers assigned an hour or more of homework each week (84 percent) was not significantly different from the corresponding national percentage (86 percent).



TABLE 4.5

Public School Teachers' Reports on Homework in Science

About how much time do you	Mississippi	Southeast	Nation		
expect a student in this class to spend doing homework each week?	Percentage and Average Scale Score				
None	2 (0.7)	1 (****) *** (**.*)	2 (0.8) 134 (4.5)!		
1/2 hour	14(3.0)	12 (4.0)	12 (2.3)		
	138(5.4)!	137 (4.0)i	142 (3.3)!		
1 hour	36 (3.6)	46 (6.1)	42 (4.1)		
	137 (1.7)	142 (2.8)!	152 (2.1)		
2 hours	35 (3.5)	24 (5.6)	28 (4.4)		
	135 (2.2)	144 (5.9)!	152 (3.0)		
More than 2 hours	13 (2.9)	17 (5.7)	15 (4.8)		
	126 (5.1)!	150 (4.9)!	156 (3.9)!		

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. **** Standard error estimates cannot be accurately determined.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



Jones, L.R., I.V.S. Mullis, S.A. Raizen, I.R. Weiss, and E.A. Weston. The 1990 Science Report Card: NAEP's Assessment of Fourth, Eighth, and Twelfth Graders. (Washington, DC: National Center for Education Statistics, 1992).

6

The eighth-grade students' reports indicated that:

- About one quarter of the eighth graders did not spend any time on science homework in a typical week (23 percent) while 31 percent spent one hour or more on their science homework each week.
- The percentage of students in Mississippi who spent an hour or more on homework each week (31 percent) was not significantly different from the percentage of students nationwide spending this much time on homework each week (34 percent).



TABLE 4.6

Public School Students' Reports on Homework in Science

If you are taking science this year,	Mississippi	Southeast	Nation		
about how much time do you spend doing science homework each week?	Percentage and Average Scale Score				
I am not taking a science	-				
course this year.	2 (0.5)	3 (0.4)	4 (0.9)		
	*** (**.*)	*** (**.*)	127 (3.1)i		
None	23 (1.4)	23 (2.3)	22 (1.5)		
	135 (2.0)	140 (2.7)	147 (1.6)		
1/2 hour	43 (1.2)	44 (1.5)	40 (1.4)		
	135 (1.7)	144 (2.4)	151 (1.1)		
1 hour	20 (0.9)	19 (1.2)	19 (0.7)		
	128 (2.2)	139 (2.9)	148 (1.6)		
2 hours	6 (0.6)	6 (0.6)	8 (0.5)		
	136 (3.8)	144 (3.9)	156 (2.7)		
3 hours	3 (0.4)	2 (0.3)	3 (0.4)		
	135 (3.5)	••• (••••)	157 (3.1)		
More than 3 hours	3 (0.3)	3 (0.3)	4 (0.4)		
	137 (4.6)	••• (••.•)	152 (3.5)		

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



60

In addition to being asked about science homework in general, students were asked how often they use a computer at home for schoolwork. Because the question was not restricted to science homework, students' reports most likely included homework for other academic areas such as English and mathematics. Given the trend that home computers are steadily assuming more importance for completing homework assignments,⁵⁰ it seems useful that NAEP monitor the prevalence of this practice and its relationship to performance.

Based on the reports of eighth graders in Mississippi, as shown in Table 4.7:

- About half of the students reported that there was no computer at home (50 percent) and another 18 percent reported never or hardly ever using their home computer to do homework.
- About one fifth of the eighth graders reported using their home computer to do homework almost every day (10 percent) or once or twice a week (11 percent).
- The average scale score for students who used a computer almost every day for homework (132) was not significantly different from that of students who never or hardly ever did so (136).
- The average scale score for students who used a computer almost every day for homework (132) was lower than that of students who used a computer at home once or twice a month (146).



TABLE 4.7

Public School Students' Reports on Using Computers at Home

How often do you use a computer	Mississippi	Southeast	Nation
at home for schoolwork?	Percent	Percentage and Average Scale Score	
There is no computer at home.	50 (1.3)	43 (1.5)	, 36 (1.2)
	130 (1.4)	138 (2.3)	143 (1.0)
Never or hardly ever	18 (0.9)	16 (0.9)	17 (0.9)
	136 (2.0)	137 (3.0)	144 (1.6)
Once or twice a month	11 (0.7)	13 (0.7)	15 (0.5)
	146 (3.0)	153 (3.2)	160 (1.8)
Once or twice a week	11 (0.7)	16 (0.8)	17 (1.1)
	137 (2.3)	149 (2.5)	157 (1.9)
Almost every day	10 (0.8)	12 (1.4)	15 (0.7)
	132 (3.1)	143 (2.8)	154 (1.9)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



⁵⁰ U.S. Department of Education. Digest of Education Statistics 1995. (Washington, DC: National Center for Education Statistics, 1995).

Computer Use in Science Instruction

The use of computers in the collection of data, interpretation of results, and communication of findings is part of the *Benchmarks for Science Literacy* and the recently published *National Science Education Standards.*⁵¹ Recommendations for facilitating science instruction in the nation's schools often include more use of computers. Computers can be used to demonstrate scientific concepts, simulate scientific phenomena, deliver instruction, and collect and analyze data. Of course, effective computer use may depend on many factors other than availability, such as teachers' training or whether computers have been incorporated into the curriculum effectively.

Computers are increasingly important in students' homes, where they are used for homework as well as for other pursuits. Since 1984 the percentage of students in grades 7 through 12 who use a computer at school or at home has increased over two-fold, to approximately 60 percent of students using a computer at school and 30 percent using one at home.⁵²

Given the potential role of computers in science instruction, NAEP asked teachers in Mississippi about the availability and use of computers in science instruction. As presented in Table 4.8, when eighth-grade science teachers in Mississippi were asked about the availability of computers, their responses indicated the following:

- In Mississippi, 50 percent of the students were in science classes where computers were not available. This percentage was greater than that for the nation (17 percent).
- The average scale score of Mississippi students whose teachers reported not having any computers available (136) was not significantly different from that of students whose teachers reported having one computer in the classroom (138).

⁵² U.S. Department of Education. Digest of Education Statistics 1995. (Washington, DC: National Center for Education Statistics, 1995).



American Association for the Advancement of Science. Benchmarks for Science Literacy. (New York: Oxford University Press, 1993); National Research Council. National Science Education Standards. (Washington, DC: National Academy Press, 1996).



TABLE 4.8

Public School Teachers' Reports on the Availability of Computers

Which best describes the	Mississippi	Southeast	Nation		
availability of computers for use by your science students?	Percentage and Average Scale Score				
None available	50 (4.3) 136 (1.8)	19 (5.4) 135 (4.9)!	17 (3.4) 149 (5.6)!		
One within the classroom	17 (3.2) 138 (4.0)	31 (8.1) 140 (2.6)!	22 (4.8) 149 (3.2)!		
Two or three within the classroom	2 (1.1)	5 (2.3) 146 (4.5)!	9 (4.6) 156 (7.2)!		
Four or more within the classroom	0 (****) **** (**.*)	0 (****) *** (**.*)	7 (3.0) 159 (2.8)!		
Available in a computer laboratory		1			
but difficult to access or schedule	20 (3.5) 134 (3.9)	39 (10.1) 148 (2.5)!	32 (4.9) 149 (2.1)		
Available in a computer laboratory			<u> </u>		
and easy to access or schedule	9 (2.8)	6 (3.0)	13 (2.6)		
-	124 (3.4)!	138 (5.7)!	148 (2.4)		

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. **** Standard error estimates cannot be accurately determined.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

BEST COPY AVAILABLE



The availability of computers varies from school to school, and the uses for computers can vary widely from class to class. Computers can be used in many ways to help students learn science, including simulating scientific phenomena or illustrating models. Also, the frequency of use can vary, regardless of the primary use in the classroom. Teachers in Mississippi were asked how they used computers and how often they were used in their science classrooms. Also, students were asked how often they used computers when doing science in school. The responses of eighth-grade public school teachers to the purpose of use for science instruction, as shown in Table 4.9, indicate the following:

- The percentage of Mississippi students whose teachers reported that they used computers for simulations and modeling (6 percent) was smaller than the corresponding national percentage (26 percent).
- The percentage of students in Mississippi whose teachers reported that their use of computers for instruction in science was for data analysis and other applications (11 percent) was smaller than that of students nationwide (20 percent).
- About three quarters of the eighth graders had teachers who reported not using a computer for science instruction (73 percent). This percentage was greater than the percentage for the nation (46 percent).

Table 4.10 presents teacher and student reports on the frequency of use of computers for science.

- In Mississippi, 77 percent of the students had teachers who reported never or hardly ever using a computer with their classes, while a small percentage reported doing so almost every day (1 percent) or once or twice a week (4 percent).
- In Mississippi, 80 percent of the students reported never or hardly ever using computers to do science in school. Furthermore, 4 percent of the students reported using computers almost every day and 7 percent used them once or twice a week.





TABLE 4.9

Public School Teachers' Reports on the Use of Computers for Instruction in Science

How do you use computers for	Mississippi	Southeast	Nation	
instruction in science?	Percentage and Average Scale Score			
Drill and practice	6 (2.3)	6 (2.2)	8 (4.4)	
	127 (6.2)!	134 (7.4)I	155 (6.8)!	
Playing science/learning games	8 (2.3)	11 (4.8)	20 (3.8)	
	141 (6.8)!	132 (4.6)!	150 (3.9)	
Simulations and modeling	6 (2.2)	13 (5.4)	26 (5.5)	
	136 (5.2)!	137 (2.9)!	153 (2.4)!	
Data analysis and other applications	11 (2.5)	20 (5.2)	20 (3.5)	
	132 (3.6)!	144 (3.5)	149 (1.6)	
Word processing	9 (2.8)	23 (4.7)	22 (3.5)	
	139 (6.7)!	146 (3.0)l	152 (2.2)	
I do not use computers for				
science instruction.	73 (4.1)	63 (5.4)	46 (4.2)	
	135 (1.6)	142 (3.1)	149 (2.1)	

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



TABLE 4.10

Public School Teachers' and Students' Reports on the Frequency of Computer Use

	Mississippi		Southeast		Nation	
How often do your students (do you) use a computer for science?	Teacher	Student	Teacher	Student	Teacher	Student
		Percent	le Score			
Never or hardly ever	77 (3.6) 136 (1.7)	80 (1.3) 136 (1.4)	74 (6.5) 144 (2.6)	66 (3.2) 142 (2.2)	62 (4.3) 150 (1.8)	67 (1.8) 150 (1.1)
Once or twice a month	17 (3.1) 133 (3.2)	9 (0.9) 130 (3.1)	20 (5.4) 136 (3.8)!	19 (1.9) 147 (3.1)	31 (4.0) 151 (2.2)	18 (1.1) 154 (1.9)
Once or twice a week	4 (1.9) 117 (6.5)!	7 (0.6) 122 (3.8)	6 (3.1) 142 (3.9)!	11 (1.4) 138 (4.5)	7 (2.4) 156 (4.0)!	10 (1.0) 145 (2.9)
Almost every day	1 (****) *** (**.*)	4 (0.4) 114 (4.0)	0 (****)	5 (0.7) 127 (4.0)	0 (0.3)	5 (0.5) 135 (3.6)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. **** Standard error estimates cannot be accurately determined.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



BEST COPY AVAILABLE

CHAPTER 5

Student Performance on Hands-On Science **Tasks**

A number of goals for science education have been put forward in a series of reports authored by government agencies and professional societies over the last 15 years.⁵³ These goals include acquisition of a core of scientific understanding, ability to apply science knowledge in practical ways, familiarity with experimental design, and the ability to carry out scientific experiments. The reports also offered recommendations for the science curricula and instruction needed to achieve these goals, such as encouraging active student participation in hands-on science, incorporating cooperative group learning, and assignment of sustained projects to students.54

A 1993 national survey indicated that science teachers devote 21 to 26 percent of class time to hands-on or manipulative activities.⁵⁵ While research on the relationship between exposure to hands-on science tasks and overall science performance is sparse and inconclusive, a recent study has demonstrated a positive relationship for eighth-grade students between the frequency of hands-on activities and their performance on a standardized assessment.56



.65

⁵³ National Science Board Commission on Precollege Education in Mathematics, Science, and Technology. Educating America for the 21st Century. (Washington, DC: National Science Foundation, 1983); Science for All Americans: A Project 2061 Report on Literacy Goals in Science, Mathematics, and Technology. (Washington, DC: American Association for the Advancement of Science, 1989); Aldridge, B.G. Essential Changes in Secondary School Science: Scope, Sequence, and Coordination. (Washington, DC: National Science Teachers Association, 1989); National Research Council. Fulfilling the Promise: Biology Education in the Nation's Schools. (Washington, DC: National Academy Press, 1990).

⁵⁴ Science Framework for the 1996 National Assessment of Educational Progress. (Washington, DC: National Assessment Governing Board, 1993).

⁵⁵ Blank, R.K. and D. Gruebel. State Indicators of Science and Mathematics Education. (Washington, DC: Council of Chief State School Officers, 1995).

⁵⁶ Stohr-Hunt, P.M. "An Analysis of Frequency of Hands-On Experience and Science Achievement." Journal of Research on Science Teaching, 33. (1996, pp. 101-109). 66

NAEP included assessments of higher-order thinking skills in science and mathematics as early as 1986 through a pilot assessment that required students to work on various hands-on tasks. Although the NAEP 1990 science assessment measured skills that were integral to scientific investigation,⁵⁷ hands-on tasks were not included. When the 1996 science framework⁵⁸ was developed in the early 1990s, it took into account the current reforms in science education by specifying three question types that probed understanding of conceptual and reasoning skills: performance exercises, constructed-response questions, and multiple-choice questions. It was envisaged that in the performance exercises, students would manipulate selected physical objects and try to solve a scientific problem using the objects before them. Hands-on tasks that met these criteria were developed for the 1996 science assessment, and each student who participated in the assessment was given an opportunity to conduct one of them.

NAEP Hands-On Science Tasks

Four different hands-on tasks were administered in the NAEP 1996 science assessment. Each task was designed to use materials to perform an investigation, make observations, evaluate experimental results, and apply problem-solving skills. In addition, tasks shared the following characteristics:

- Diagrams were included to guide students through the procedures;
- Multiple-choice and constructed-response questions were embedded throughout the tasks; and
- Scientific investigation was integrated with conceptual understanding and practical reasoning.

The creation of the hands-on tasks presented special challenges. Since the assessment was administered in a variety of settings, ranging from laboratories to cafeterias, all of the required equipment necessary to conduct each task had to be provided in a self-contained kit produced according to standard specifications to ensure uniformity. There were some limitations on materials and equipment. For example, live materials (with the exception of seeds) and equipment that required an electric outlet were not used. Safety was also an important concern and was addressed in a number of ways. The state's safety regulations were considered; no toxic or corrosive chemicals were used; assessment administrators were trained in appropriate laboratory safety; and students were provided with goggles for some tasks.

A brief summary of one of the four hands-on tasks is described in this chapter. Several questions from the hands-on task are also shown with their scoring criteria.

⁵⁸ Science Framework for the 1996 National Assessment of Educational Progress. (Washington, DC: National Assessment Governing Board, 1993).



⁵⁷ Science Objectives: 1990 Assessment. (Princeton, NJ: The National Assessment of Educational Progress, 1989).

Sample Questions from a Task

A brief summary of one of the four tasks given to grade 8 students in Mississippi is presented below with sample questions in Figures 5.1 and 5.2.

Salt Solutions: Estimating the Salt Concentration of an Unknown Salt Solution Using the "Floating Pencil Test"

An instrument constructed from a pencil and thumbtack served as a hydrometer in this task. Students were asked to observe, measure, and compare the lengths of a portion of the pencil, marked with calibrations for ease of measurement, that floated above the surface in distilled water and in a 25 percent salt solution. Based on these observations, students were asked to predict how the addition of more salt to the salt solution would affect the floating pencil. Students then measured the length of the pencil that floated above the surface of a solution of unknown salt concentration and used the results of their previous observations to estimate the salt concentration of the unknown solution. The task assessed students' ability to make simple observations, measure length using a ruler, apply observations to an unknown, draw a graph, interpolate from graphical data, and make a generalized inference from observations. The task also assessed students' understanding of the value of performing multiple trials of the same procedure.

Figure 5.1 shows a data table that was presented in the first stage of the task. Questions 3, 4, and 5 are also presented in this figure. Students were asked to measure the length of pencil floating above the surface in three solutions: distilled water, a 25 percent salt solution, and a solution containing an unknown concentration of salt. The students recorded two measurements for each of the 3 solutions in Table 1 and calculated the average of each pair of readings. The scoring rubrics for **Complete** responses are shown in Figure 5.1.



67



FIGURE 5.1

Salt Solutions Task: Questions 3, 4, and 5

3. Now take the pencil out of the water and dry it with a paper towel.

Use the ruler to measure the length of the pencil that was above the water. Record the length in Table 1 below under Measurement 1.

TABLE 1

	Length of Pencil Above Water Surface (cm)					
Type of Solution	Measurement 1	Measurement 2	Average			
Distilled Water						
Salt Solution						
Unknown Salt Solution						

- 4. Now place the pencil back in the distilled water and repeat steps 2 and 3. Record your measurement in Table 1 under **Measurement 2**.
- 5. Calculate the average of Measurements 1 and 2 and record the result in the data table.

(You can calculate the average by adding Measurement 1 + Measurement 2 and then dividing by two.)

SCORING RUBRIC

Measurement: A Complete response has three pairs of measurements that agree within a given tolerance and also are in the correct relative order.

Average: A Complete response correctly calculates the average for each set of data.



Students were then presented with graph paper and asked to plot the average of the measurements for distilled water and 25 percent salt solution against salt concentration. They were told to assume a linear relationship between the height of the pencil above the solution and the salt concentration, and then asked to use the graph to determine the salt concentration of the unknown solution (Figure 5.2). The scoring rubric for a Complete response is also shown.

THE NATION'S REPORT NEED	FIGURE 5.2	
1995 State Assessment	Salt Solutions Task: Question 14	

the unl	cnown solu	tion?	·	
Explain	n how you	determined ye	our answer.	
				<u></u>
	-			

Unknown Solution: A Complete response gave a salt concentration consistent with the graph and correctly explained how the graph was used to obtain the answer.

Instruction Related to Scientific Investigation

Research devoted to the effectiveness of hands-on tasks is ongoing, although there is evidence that eighth graders who are exposed to hands-on activities more frequently perform better on standardized assessments.⁵⁹ Eighth-grade science teachers in Mississippi were asked about the emphasis they placed on laboratory skills and data analysis in their science classes and about the frequency and nature of hands-on activities or investigations assigned by them. Students were asked about the frequency and nature of hands-on activities or investigations conducted by them.

As mentioned before, a direct cause-and-effect relationship between educational environment and student scores on the NAEP science assessment is not implied. For instance, the motivation and expectations of teachers or students reporting hands-on investigations hardly ever or once or twice a week may be a factor in the average score differences. However, responses to teacher (and school) questionnaires provide a broad view of educational practices that should prove useful for improving instruction and setting policy.



⁵⁹ Stohr-Hunt, P.M. "An Analysis of Frequency of Hands-On Experience and Science Achievement." Journal of Research on Science Teaching, 33. (1996, pp. 101-109).

Teachers' and students' responses regarding scientific investigation are presented in Tables 5.1 through 5.5.

- The percentage of eighth-grade students in Mississippi whose teachers reported placing heavy emphasis on the development of laboratory skills and techniques (24 percent) was smaller than the percentage nationwide (42 percent). Students whose teachers reported heavy emphasis on laboratory skills and techniques in Mississippi had an average scale score (135) which was lower than that of students nationwide whose teachers reported this (153).
- The percentage of eighth-grade students in Mississippi whose teachers reported moderate to heavy emphasis on the development of data analysis skills (81 percent) was not significantly different from* that of students nationwide (89 percent). Eighth-grade students whose teachers reported moderate to heavy emphasis on data analysis skills had an average science scale score (134) which did not differ significantly from that of students whose teachers reported little or no emphasis on the development of data analysis skills (137).



TABLE 5.1

Public School Teachers' Reports on Science Instruction Related to Performance Tasks

Think about your plans for your science instruction during the entire year. About how much emphasis will you give to each of the following?	Mississippi Southeast Nation				
	Percentage and Average Scale Score				

21 (3.0) 129 (2.9)	12 (4.1) 124 (5.3)!	13 (2.5) 135 (3.6)!
55 (3.8)	48 (7.6) 145 (2.7)	44 (4.7) 152 (2.0)
24 (3.4)	40 (6.5)	42 (4.5) 153 (2.1)
100 (2.0)	().	.55 (=,
19 (3.0) 137 (2.7)	18 (3.7) 133 (5.0)!	11 (2.7) 139 (5.5)!
60 (3.9) 134 (1.9)	51 (5.0) 144 (2.7)	65 (5.3) 151 (1.6)
22 (3.4) 132 (3.1)	30 (4.0) 145 (3.1)	24 (4.3) 153 (3.0)
	129 (2.9) 55 (3.8) 136 (2.1) 24 (3.4) 135 (2.9) 19 (3.0) 137 (2.7) 60 (3.9) 134 (1.9) 22 (3.4)	129 (2.9) 124 (5.3)! 55 (3.8) 48 (7.6) 136 (2.1) 145 (2.7) 24 (3.4) 40 (6.5) 135 (2.9) 144 (1.9)! 19 (3.0) 18 (3.7) 137 (2.7) 133 (5.0)! 60 (3.9) 51 (5.0) 134 (1.9) 144 (2.7) 22 (3.4) 30 (4.0)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science

^{*} Although the difference may appear large, recall that "significance" here refers to "statistical significance."



- About two thirds of the eighth-grade students in Mississippi (67 percent)
 had teachers who reported doing a science demonstration at least once
 a week, not significantly different from* the percentage of students
 nationwide (59 percent). Less than half of eighth-grade students in
 Mississippi (37 percent) reported that their teacher performed science
 demonstrations at least once a week.
- The eighth-grade students in Mississippi whose teachers reported doing a science demonstration at least once a week had an average scale score (135) which was lower than that of their national counterparts (151).



TABLE 5.2

Public School Teachers' and Students' Reports on the Frequency of Science Demonstrations

	Mississippi		Southeast		Nati <i>o</i> n		
How often do you (does your teacher) do a science demonstration?	Teacher	Student	Teacher	Student	Teacher	Student	
		Percent	age and A	verage Scale Score			
Never or hardly ever	5 (1.7) 126 (4.9)!	32 (1.1) 129 (1.7)	2 (1.0)	31 (2.8) 133 (2.7)	2 (0.8) 149 (11.6)!	30 (1.3) 141 (1.5)	
Once or twice a month	28 (3.5) 134 (2.6)	31 (1.0) 138 (1.7)	41 (7.3) 142 (3.8)!	31 (1.7) 144 (2.3)	39 (4.1) 150 (2.0)	29 (1.1) 151 (1.3)	
Once or twice a week	58 (3.8) 135 (2.0)	25 (0.9) 135 (1.6)	43 (5.2) 142 (2.1)	26 (1.7) 149 (2.4)	49 (3.5) 152 (1.9)	28 (1.2) 156 (1.4)	
Almost every day	9 (2.1) 137 (3.0)!	12 (0.9) 128 (2.8)	15 (3.9) 148 (2.6)!	12 (1.6) 145 (2.2)	10 (2.3) 144 (2.0)!	14 (0.9) 153 (2.0)	

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

BEST COPY AVAILABLE



^{*} Although the difference may appear large, recall that "significance" here refers to "statistical significance."

- The percentage of eighth-grade students in Mississippi whose teachers reported their science students performed hands-on tasks once a week or more (61 percent) was smaller than the national percentage (83 percent). The percentage of students in Mississippi whose teachers reported their students never or hardly ever did hands-on tasks (9 percent) was greater than nationwide levels (1 percent).
- The eighth-grade students in Mississippi whose teachers reported their students did hands-on tasks at least once a week had an average science scale score (134) which was lower than that of students nationwide whose teachers reported this same level of hands-on task experience (153).
- The eighth-grade students in Mississippi whose teachers reported their students did hands-on tasks almost every day had an average scale score (136) which did not differ significantly from that of students whose teachers reported doing hands-on activities once or twice a month (136).



TABLE 5.3

Public School Teachers' and Students' Reports on the Frequency of Hands-on Activities or Investigations

How often do your students (do you)	Mississippi		Southeast		Nat	Nation	
do hands-on activities or	Teacher	Student	Teacher	Student	Teacher	Student	
investigations in science?	Percentage and Average Scale Score						
Never or hardly ever	9 (2.5) 131 (3.6)!	30 (1.6) 129 (2.0)	2 (1.4)	22 (3.0) 128 (2.2)	1 (0.6) 119 (4.0)!	18 (1.1) 134 (1 <i>.</i> 2)	
Once or twice a month	29 (3.5) 136 (3.3)	34 (1.2) 140 (1.6)	39 (7.7) 140 (4.0)!	32 (3.1) 146 (2.7)	16 (2.4) 140 (3.4)	32 (1.5) 152 (1.5)	
Once or twice a week	51 (4.0) 134 (2.0)	23 (1.3) 132 (2.0)	43 (6.2) 146 (2.0)!	28 (2.3) 147 (1.8)	64 (3.5) 153 (1.5)	33 (1.3) 155 (1.2)	
Almost every day	10 (2.0) 136 (3.1)!	12 (1.1) 127 (2.2)	15 (5.1) 146 (4.0)!	17 (3.1) 146 (2.9)!	19 (3.2) 152 (2.2)	18 (1.1) 151 (1.5)	

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



- About three quarters of the eighth-grade students in Mississippi (79 percent) had teachers who reported assigning science projects in school which take a week or more to complete. About half of the students (52 percent) reported receiving such assignments and their average scale score was 131.
- The average scale score of students who reported doing science projects or investigations that take a week or more (131) was not significantly different from that of students who did not (135).



TABLE 5.4

Public School Teachers' and Students' Reports on Long-Term Science Projects

Do you ever assign (do) individual or group science projects or investigations in school that take a week or more?	Missi	ssippi	Sout	heast	Nat	ion
	Teacher	Student	Teacher	Student	Teacher	Student
		Percent	age and Av	erage Scal	e Score	

Yes	79 (3.5)	52 (2.2)	70 (7.2)	62 (3.6)	82 (2.6)	63 (2.8)
	134 (1.8)	131 (1.6)	143 (1.9)	145 (2.1)	151 (1.3)	151 (1.3)
No	21 (3.5)	48 (2.2)	30 (7.2)	38 (3.6)	18 (2.6)	37 (2.8)
	135 (2.2)	135 (1.8)	143 (5.0)!	136 (2.6)	147 (3.4)	148 (1.7)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



- In Mississippi, the eighth-grade students who reported designing and carrying out their own scientific investigations once a week or more frequently (13 percent) received an average scale score of 118.
- The average scale score for Mississippi students who reported designing and carrying out their own science investigations once a week or more (118) was lower than that for students who reported doing this once or twice a month (135).



TABLE 5.5

Public School Students' Reports on Independent Science Investigations

When you study science in school, how often	Mississippi	Southeast	Nation		
do you design and carry out your own science investigations?	Percentage and Average Scale Score				
Never or hardly ever	67 (1.2)	61 (3.2)	63 (1.1)		
	136 (1.5)	141 (2.1)	151 (1.0)		
Once or twice a month	19 (1.0) 135 (2.0)	23 (2.4)	23 (0.8)		
Once or twice a week	8 (0.6)	11 (1.0)	10 (0.6)		
	119 (3.5)	140 (2.9)	142 (2.3)		
Almost every day	5 (0.4)	5 (0.5)	5 (0.4)		
	116 (3.7)	132 (6.0)	137 (2.5)		

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



CHAPTER 6

Influences Beyond School that Facilitate Learning Science

The home environment can be an important support for the school environment. To examine the relationship between science scale scores and home factors, data regarding students' responses to questions about home factors and principals' responses to questions about parental involvement in the school were examined. The student questionnaires also asked students how often they had changed schools because of household moves to examine the impact of student mobility on academic achievement.

Students' attitudes toward science can influence their performance in the assessment. For example, in a recent large scale science assessment, students who agreed that science learning is useful for the future and that science should be required in school performed better than those who disagreed with these statements.⁶⁰ These attitudes toward science may be attributed to factors within the school and external influences. The beliefs and general impressions that secondary school students form about science can affect not only their performance in assessments but also their decisions about pursuing scientific careers in the future.⁶¹



Campbell, J.R., C.M. Reese, C. O'Sullivan, and J.A. Dossey. NAEP 1994 Trends in Academic Progress. (Washington, DC: National Center for Education Statistics, 1996).

⁶¹ Gallagher, S.A. "Middle School Classroom Predictors of Science Persistence." Journal of Research in Science Teaching, 1994, 33. pp. 721-734.

Discussing Studies at Home

The importance of schoolwork for students and their families can by measured by how often it is discussed at home. When students discuss academic work at home, they create an important link between home and school. Recent NAEP assessments in various subject areas have found a positive relationship between discussing studies at home and student performance.⁶²

The NAEP 1996 assessment asked students to report on how frequently they discuss schoolwork at home. As shown in Table 6.1, the results for eighth graders attending public schools in Mississippi indicate that:

- Less than half of the eighth graders (42 percent) said they discussed their schoolwork at home almost every day. This percentage was greater than the percentage who said they never or hardly ever had such discussions (21 percent).
- The average scale score for students who discussed their schoolwork almost every day (135) was not significantly different from that for students who never or hardly ever did so (132).



TABLE 6.1

Public School Students' Reports on Discussing Studies at Home

How often do you discuss things	Mississippi	Southeast	Nation
you have studied in school with someone at home?	Percent	e Score	
Never or hardly ever	21 (0.8)	23 (1.3)	21 (0.8)
	132 (1.9)	132 (2.6)	141 (1.5)
Once or twice a month	9 (0.5)	10 (0.6)	9 (0.4)
	132 (2.8)	143 (3.4)	149 (1.6)
Once or twice a week	27 (0.9)	29 (1.5)	28 (1.0)
	134 (1.7)	146 (2.2)	151 (1.3)
Almost every day	42 (1.1)	38 (1.0)	41 (1.1)
	135 (1.7)	145 (1.8)	153 (1 <i>.</i> 2)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

⁶² Campbell, J.R., P.L. Donahue, C.M. Reese, and G.W. Phillips. NAEP 1994 Reading Report Card for the Nation and the States. (Washington, DC: National Center for Education Statistics, 1996); Beatty, A.S., C.M. Reese, H.R. Persky, and P. Carr. NAEP 1994 U.S. History Report Card. (Washington, DC: National Center for Education Statistics, 1996); Persky, H.R., C.M. Reese, C.Y. O'Sullivan, S. Lazer, J. Moore, and S. Shakrani. NAEP 1994 Geography Report Card. (Washington, DC: National Center for Education Statistics, 1996).



Literacy Materials in the Home

Students can learn much about science by reading materials outside the classroom. For example, scientific information can often be found in mainstream newspaper and magazine articles. Also, the availability of reading and reference materials at home may be an indicator of the value placed on learning by the parents.¹ In recent NAEP assessments, a positive relationship has been reported between print materials in the home and average scale scores.²

The NAEP science assessment asked students whether their families had more than 25 books, an encyclopedia, a newspaper, or any magazines in their home. Table 6.2 shows the percentages of eighth-grade public school students reporting that their families have all four types, only three types, or two or fewer types of these literacy materials. The table also presents students' corresponding average scale scores. Based on their responses:

- Less than half of the students in Mississippi (41 percent) reported having all four types of literacy materials in their homes. This percentage was smaller than the percentage for the nation (47 percent).
- The percentage of students in Mississippi reporting having two or fewer types of these materials (27 percent) was smaller than the percentage having all four types (41 percent). The percentage having two or fewer types was somewhat greater than the percentage for the nation (24 percent).
- The average science scale score for students in Mississippi with all four types of literacy materials (142) was higher than that for students with two or fewer types (121).

² Campbell, J.R., P.L. Donahue, C.M. Reese, and G.W. Phillips. NAEP 1994 Reading Report Card for the Nation and the States. (Washington, DC: National Center for Education Statistics, 1996); Beatty, A.S., C.M. Reese, H.R. Persky, and P. Carr. NAEP 1994 U.S. History Report Card. (Washington, DC: National Center for Education Statistics, 1996); Persky, H.R., C.M. Reese, C.Y. O'Sullivan, S. Lazer, J. Moore, and S. Shakrani. NAEP 1994 Geography Report Card. (Washington, DC: National Center for Education Statistics, 1996).



¹ Rogoff, B. Apprenticeship in Thinking: Cognitive Development in Social Context. (New York: Oxford University Press, 1990).



TABLE 6.2

Public School Students' Reports on Literacy Materials in the Home

How many of the following types of reading materials are in your home	Mississippi Southeast Nation			
(more than 25 books, an encyclopedia, a newspaper, magazines)?	Percent	age and Average Scal	le Score	

Zero to two	27 (1.0)	29 (1.6)	24 (0.7)
	121 (2.1)	126 (2.6)	132 (1.2)
Three	32 (1.1)	30 (0.7)	29 (0.8)
	133 (1.9)	143 (2.1)	149 (1.0)
Four	41 (1.4)	42 (1.8)	47 (1.1)
	142 (1.3)	151 (1.7)	158 (1.2)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

Television Viewing Habits

Past NAEP assessments have shown that more than 40 percent of eighth-grade students reported watching four or more hours of television each day. A major concern is that watching television reduces the time spent on homework and related academic activities. Although the effects of such extensive television exposure are difficult to document, a generally negative relationship exists between NAEP score results and number of television hours watched.3 The recent TIMSS assessment shows a similar pattern for most countries. In general, beyond one to two hours of daily television viewing, the more that eighth graders reported watching, the lower their science achievement.⁴

Students were asked how much television (including videotapes) they usually watched each school day. The results for eighth-grade public school students in Mississippi are shown in Table 6.3 and indicate the following:



Campbell, J.R., P.L. Donahue, C.M. Reese, and G.W. Phillips. NAEP 1994 Reading Report Card for the Nation and the States. (Washington, DC: National Center for Education Statistics, 1996); Beatty, A.S., C.M. Reese, H.R. Persky, and P. Carr. NAEP 1994 U.S. History Report Card. (Washington, DC: National Center for Education Statistics, 1996); Persky, H.R., C.M. Reese, C.Y. O'Sullivan, S. Lazer, J. Moore, and S. Shakrani. NAEP 1994 Geography Report Card. (Washington, DC: National Center for Education Statistics, 1996); Campbell, J.R., C.M. Reese, C.Y. O'Sullivan, and J.A. Dossey. NAEP 1994 Trends in Academic Progress. (Washington, DC: National Center for Education Statistics, 1996).

Beaton, A.E., M.O. Martin, I.V.S. Mullis, E.J. Gonzalez, T.A. Smith, and D.L. Kelly. Science Achievement in the Middle School Years: IEA's Third International Mathematics and Science Study (TIMSS). (Chestnut Hill, MA: TIMSS International Study Center at Boston College, 1996).

- Among eighth graders, 29 percent reported watching six or more hours of television on a typical day. This percentage was greater than the percentage who reported watching one hour or less (12 percent).
- The percentage of eighth graders in Mississippi who reported watching six or more hours of television a day (29 percent) was greater than the percentage for the nation (17 percent).
- The average science scale score for eighth-grade students who reported watching two to three hours of television a day (140) was higher than that for students who reported watching one hour or less (133).
- The average science scale score for eighth graders who reported watching two to three hours of television a day (140) was higher than that for students who reported watching six hours or more (123).



TABLE 6.3

Public School Students' Reports on Television Viewing Habits

On a school day, about how many	Mississippi	. Southeast	Nation	
hours do you usually watch TV or videotapes outside of school hours?	Percentage and Average Scale Score			
1 hour or less	12 (0.8)	14 (0.8)	19 (1.0)	
	133 (2.6)	143 (3.1)	156 (2.0)	
2 to 3 hours	32 (1.0)	36 (1.8)	40 (1.3)	
·	140 (2.0)	148 (1.9)	154 (1.2)	
4 to 5 hours	27 (1.0)	28 (1.2)	24 (0.6)	
	137 (1.6)	147 (2.0)	148 (1.0)	
6 hours or more	29 (1.1)	22 (1.6)	17 (0.7)	
	123 (1.9)	123 (2.1)	130 (1.1)	

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

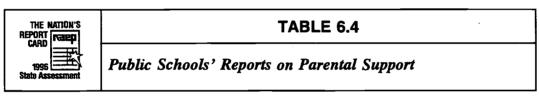


Parental Support

When parents are involved in their children's education, both children and parents are likely to benefit. Research on students at risk has shown that parents' participation in their child's education has more effect on the child's performance than parents' income or education.⁶⁷ Parental involvement is naturally part of the home environment, but it is also increasingly sought in the school.

As part of the NAEP assessment, the principals of participating students were asked about parental involvement in their schools. Table 6.4 presents the results for eighth graders in public schools in Mississippi. According to these results:

- Overall, a large majority of the eighth-grade students attended schools where principals characterized parental support as very positive (18 percent) or somewhat positive (67 percent).
- The average scale score for eighth graders attending school where parental support was characterized as very positive (140) was higher than that for the 15 percent of students whose principals reported somewhat to very negative parental support (127).



How would you characterize	Mississippi	Southeast	Nation	
parental support for student achievement within your school?	Percentage and Average Scale Score			
Somewhat to very negative	15 (3.6)	0 (****)	7 (2.6)	
	127 (2.7)!	*** (**,*)	154 (2.1)!	
Somewhat positive	67 (4.7)	67 (6.3)	61 (5.6)	
	133 (1.9)	137 (2.3)	148 (1.4)	
Very positive	18 (3.4)	33 (6.3)	31 (4.7)	
	140 (3.8)	151 (3.4)!	151 (3.3)	

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). ! Interpret with caution — the nature of the sample does not allow accurate determination of the variability of this statistic. *** Sample size is insufficient to permit a reliable estimate. **** Standard error estimates cannot be accurately determined.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

⁶⁷ U.S. Department of Education. Mapping out the National Assessment of Title I: The Interim Report — 1996. (Washington, DC: Office of Educational Research and Improvement, 1996).



Student Mobility

The United States has long been a nation "on the move." Research indicates that moving more than once or twice during a school career lowers student performance. Students who attend the same school throughout their careers are most likely to graduate, while the most mobile of the school populations have the highest rates of failure and dropping out. The effects of high mobility are far-reaching; schools with high mobility rates depress performance even for students who do not move. 68

To examine the relationship between mobility and science performance, the NAEP assessment asked students how many times since starting first grade they had changed schools due to changes in where they lived. Table 6.5 shows results for eighth-grade public school students in Mississippi.

- In terms of student mobility, 50 percent of eighth graders reported not moving since starting first grade while 5 percent of students reported moving six or more times. The students with the highest reported mobility had an average scale score (130) that did not differ significantly from that of students who reported not moving (133).
- The percentage of students in Mississippi who reported moving six or more times (5 percent) was not significantly different from the percentage for the nation (6 percent).

THE NATION'S REPORT REPORT	TABLE 6.5
1996 State Assessment	Public School Students' Reports on Mobility

Since you started first grade, how many times have you changed	Mississippi	Southeast	Nation	
schools, not counting when you were promoted to the next grade?	Percentage and Average Scale Score			
	50 (4 8)			

50 (1.2)	44 (2.7)	44 (1.2)
133 (1.6)	145 (2.7)	153 (1.3)
18 (0.7)	18 (0.9)	19 (0.8)
133 (2.5)	148 (2.6)	154 (1.4)
8 (0.7)	11 (0.8)	10 (0.4)
133 (3.3)	140 (3.1)	145 (1.4)
9 (0.7)	11 (0.9)	11 (0.6)
133 (2.8)	131 (2.1)	141 (2.3)
9 (0.6)	11 (0.7)	10 (0.5)
136 (2.6)	135 (2.7)	142 (1.7)
5 (0.5)	6 (0.6)	6 (0.3)
130 (4.0)	138 (3.4)	141 (2.0)
	133 (1.6) 18 (0.7) 133 (2.5) 8 (0.7) 133 (3.3) 9 (0.7) 133 (2.8) 9 (0.6) 136 (2.6) 5 (0.5)	133 (1.6) 145 (2.7) 18 (0.7) 18 (0.9) 133 (2.5) 148 (2.6) 8 (0.7) 11 (0.8) 133 (3.3) 140 (3.1) 9 (0.7) 11 (0.9) 133 (2.8) 131 (2.1) 9 (0.6) 11 (0.7) 136 (2.6) 135 (2.7) 5 (0.5) 6 (0.6)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



81

ERIC Clearinghouse on Urban Education. Highly Mobile Students: Educational Problems and Possible Solutions. (New York: ERIC Clearinghouse on Urban Education, ERIC/CUE Digest, Number 73, 1991).

URL: http://www.ed.gov/databases/ERIC_Digests/ed338745.html. See also The Condition of Education 1995/indicator46 at URL: http://www.ed.gov/NCES/pubs/ce/c9546a01.html.

Students' Views About Science

Science educators have been interested in the relationship between students' attitudes and student performance for several decades. A considerable body of research has shown a correlation between students' attitudes and their performance in science, with positive attitudes typically being associated with higher performance. Therefore, the 1996 NAEP science assessment asked several questions to gauge students' attitudes toward science. Table 6.6 shows the responses for eighth graders in Mississippi.

- In Mississippi, 41 percent of eighth graders agreed that science is useful for solving everyday problems. The average scale score for these students (135) was not significantly different from that for students who were unsure about this statement or who did not agree with it (132).
- In Mississippi, 42 percent of the students agreed that learning science is mostly memorizing facts. The average scale score for eighth graders who felt that learning science is mostly memorizing (131) was not significantly different from the average scale score of students who were unsure or disagreed with this statement (134).



TABLE 6.6

Public School Students' Views About Science

How much do you agree with the following statements?	Mississippi Southeast Nation						
	Percentage and Average Scale Score						

Science is useful for solving everyday problems.	·		
Disagree	25 (1.0)	25 (1.4)	25 (1.0)
	127 (1.7)	132 (2.7)	139 (1.5)
Not sure	33 (1.0)	37 (1.0)	-35 (0.7)
	136 (2.0)	142 (2.0)	150 (0.9)
Agree	41 (1.2)	38 (1.9)	40 (1.1)
-	135 (1.5)	147 (1.9)	155 (1.1)
Learning science is mostly memorizing.			
Disagree	27 (1.1)	29 (1.8)	30 (0.8)
	137 (1.7)	143 (2.5)	150 (1.3)
Not sure	32 (1.2)	38 (1.3)	37 (0.5)
	132 (2.1)	141 (2.0)	148 (1.1)
Agree	42 (1.3)	33 (1.4)	33 (0.9)
_	131 (1.7)	141 (2.3)	149 (1.1)

The NAEP science scale ranges from 0 to 300. The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details).

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.

⁶⁹ Weinburg, M. "Gender Differences in Student Attitudes Toward Science: A Meta Analysis of the Literature from 1970 to 1991." Journal of Research in Science Teaching 1985, 32. pp. 387-398.



APPENDIX A

Reporting NAEP 1996 Science Results

A.1 Participation Guidelines

As was discussed in the Introduction, unless the overall participation rate for a jurisdiction is sufficiently high, the assessment results for that jurisdiction may be subject to appreciable nonresponse bias. Moreover, even if the overall participation rate is high, significant nonresponse bias may exist if the nonparticipation that does occur is heavily concentrated among certain types of schools or students. The following guidelines concerning school and student participation rates in the state assessment program were established to address four significant ways in which nonresponse bias could be introduced into the jurisdiction sample estimates.

The first three guidelines describe the determination of whether a jurisdiction is eligible to have its results published. Guidelines 4-11 describe conditions under which a jurisdiction's published results will include a notation. Such a notation would indicate the possibility of bias in particular results, due to nonresponse from segments of the sample. Note that in order for a jurisdiction's results to be published without notations, that jurisdiction must comply with all guidelines. (A thorough discussion of the NAEP participation guidelines can be found in the Technical Report of the NAEP 1996 State Assessment Program in Science.)

Guidelines on the Publication of NAEP Results

Guideline 1 — Publication of Public School Results

A jurisdiction will have its public school results published in the NAEP 1996 Science Report Card (or in other reports that include all state-level results) if and only if its weighted participation rate for the initial sample of public schools is greater than or equal to 70 percent. Similarly, a jurisdiction will receive a separate NAEP 1996 Science State Report if and only if its weighted participation rate for the initial sample of public schools is greater than or equal to 70 percent.



Guideline 2 — Publication of Nonpublic School Results

A jurisdiction will have its nonpublic school results published in the NAEP 1996 Science Report Card (or in other reports that include all state-level results) if and only if its weighted participation rate for the initial sample of nonpublic schools is greater than or equal to 70 percent AND meets minimum sample size requirements.\(^1\) A jurisdiction eligible to receive a separate NAEP 1996 Science State Report under guideline 1 will have its nonpublic school results included in that report if and only if that jurisdiction's weighted participation rate for the initial sample of nonpublic schools is greater than or equal to 70 percent AND meets minimum sample size requirements. If a jurisdiction meets guideline 2 but fails to meet guideline 1, a separate NAEP 1996 Science State Report will be produced containing only nonpublic school results.

Guideline 3 — Publication of Combined Public and Nonpublic School Results

A jurisdiction will have its combined results published in the NAEP 1996 Science Report Card (or in other reports that include all state-level results) if and only if both guidelines 1 and 2 are satisfied. Similarly, a jurisdiction eligible to receive a separate NAEP 1996 Science State Report under guideline 1 will have its combined results included in that report if and only if guideline 2 is also met.

Guidelines for Notations of NAEP Results

Guideline 4 — Notation for Overall Public School Participation Rate

A jurisdiction that meets guideline 1 will receive a notation if its weighted participation rate for the initial sample of public schools was below 85 percent **AND** the weighted public school participation rate after substitution was below 90 percent.

Guideline 5 — Notation for Overall Nonpublic School Participation Rate

A jurisdiction that meets guideline 2 will receive a notation if its weighted participation rate for the initial sample of nonpublic schools was below 85 percent AND the weighted nonpublic school participation rate after substitution was below 90 percent.

Minimum participation size requirements for reporting nonpublic school data consist of two components: (1) a school sample size of six or more participating schools and (2) an assessed student sample size of at least 62.



Guideline 6 — Notation for Strata-Specific Public School Participation Rate

A jurisdiction that is not already receiving a notation under guideline 4 will receive a notation if the sample of public schools included a class of schools with similar characteristics that had a weighted participation rate (after substitution) of below 80 percent, and from which the nonparticipating schools together accounted for more than five percent of the jurisdiction's total weighted sample of public schools. The classes of schools from each of which a jurisdiction needed minimum school participation levels were determined by degree of urbanization, minority enrollment, and median household income of the area in which the school is located.

Guideline 7 — Notation for Strata-Specific Nonpublic School Participation Rate

A jurisdiction that is not already receiving a notation under guideline 5 will receive a notation if the sample of nonpublic schools included a class of schools with similar characteristics that had a weighted participation rate (after substitution) of below 80 percent, and from which the nonparticipating schools together accounted for more than five percent of the jurisdiction's total weighted sample of nonpublic schools. The classes of schools from each of which a jurisdiction needed minimum school participation levels were determined by type of nonpublic school (Catholic versus non-Catholic) and location (metropolitan versus nonmetropolitan).

Guideline 8 — Notation for Overall Student Participation Rate in Public Schools

A jurisdiction that meets guideline 1 will receive a notation if the weighted student response rate within participating public schools was below 85 percent.

Guideline 9 — Notation for Overall Student Participation Rate in Nonpublic Schools

A jurisdiction that meets guideline 2 will receive a notation if the weighted student response rate within participating nonpublic schools was below 85 percent.



86

Guideline 10 — Notation for Strata-Specific Student Participation Rates in Public Schools

A jurisdiction that is not already receiving a notation under guideline 8 will receive a notation if the sampled students within participating public schools included a class of students with similar characteristics that had a weighted student response rate of below 80 percent, and from which the nonresponding students together accounted for more than five percent of the jurisdiction's weighted assessable public school student sample. Student groups from which a jurisdiction needed minimum levels of participation were determined by the age of the student, whether or not the student was classified as a student with a disability (SD) or of limited English proficiency (LEP), and the type of assessment session (monitored or unmonitored), as well as school level of urbanization, minority enrollment, and median household income of the area in which the school is located.

Guideline 11 — Notation for Strata-Specific Student Participation Rates in Nonpublic Schools

A jurisdiction that is not already receiving a notation under guideline 9 will receive a notation if the sampled students within participating nonpublic schools included a class of students with similar characteristics that had a weighted student response rate of below 80 percent, and from which the nonresponding students together accounted for more than five percent of the jurisdiction's weighted assessable nonpublic school student sample. Student groups from which a jurisdiction needed minimum levels of participation were determined by the age of the student, whether or not the student was classified as a student with a disability (SD) or of limited English proficiency (LEP), and the type of assessment session (monitored or unmonitored), as well as type and location of school.



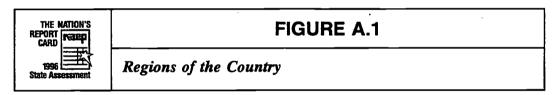
A.2 NAEP Reporting Groups

The NAEP state assessment program provides results for groups of students defined by shared characteristics — region of the country, gender, race/ethnicity, parental education, type of school, and participation in federally funded Title I programs and the free/reduced-price lunch component of the National School Lunch Program. Based on criteria described later in this appendix, results are reported for subpopulations only when sufficient numbers of students and adequate school representation are present. For public school students, there must be at least 62 students in a particular subgroup from at least 5 primary sampling units (PSUs). For nonpublic school students, the minimum requirement is 62 students in a particular subgroup from at least 6 different schools. However, the data for all students, regardless of whether their subgroup was reported separately, were included in computing overall results for Mississippi. Definitions of the subpopulations referred to in this report are presented on the following pages.

Region

Results are reported for four regions of the nation: Northeast, Southeast, Central, and West. The states included in each region are shown in Figure A.1. All 50 states and the District of Columbia are listed. Territories and the two Department of Defense Education Activity jurisdictions were not assigned to a region.

Regional results are based on national assessment samples, not on aggregated state assessment program samples. Thus, the regional results are based on a *different* and *separate* sample from that used to report the state results.



NORTHEAST	NORTHEAST SOUTHEAST		WEST
Connecticut Delaware District of Columbia Maine Maryland Massachusetts New Hampshire New Jersey New York Pennsylvania Rhode Island Vermont Virginia*	Alabama Arkansas Florida Georgia Kentucky Louisiana Mississippi North Carolina South Carolina Tennessee Virginia* West Virginia	Illinois Indiana Iowa Kansas Michigan Minnesota Missouri Nebraska North Dakota Ohio South Dakota Wisconsin	Alaska Arizona Califomia Colorado Hawaii Idaho Montana Nevada New Mexico Oklahoma Oregon Texas Utah Washington Wyoming

Note: The part of Virginia that is included in the Washington, DC, metropolitan area is included in the Northeast region; the remainder of the state is in the Southeast region.



² For the State Assessment Program, a PSU is most often a single school; for the national assessment, a PSU is a selected geographic region (a county, group of counties, or metropolitan statistical area).

Gender

Results are reported separately for males and females.

Race/Ethnicity

The racial/ethnic results presented in this report attempt to provide a clear picture based on several sources. The race/ethnicity variable is an imputed definition of race/ethnicity derived from up to three sources of information. This variable is used for race/ethnicity subgroup comparisons. Two questions from the student demographics questionnaire were used in the determination of derived race/ethnicity:

If you are Hispanic, what is your Hispanic background?

- ° I am not Hispanic.
- Mexican, Mexican American, or Chicano
- Puerto Rican
- ° Cuban
- Other Spanish or Hispanic background

Students who responded to this question by filling in the second, third, fourth, or fifth oval were considered Hispanic. For students who filled in the first oval, did not respond to the question, or provided information that was illegible or could not be classified, responses to the question below were examined in an effort to determine race/ethnicity.

Which best describes you?

- White (not Hispanic)
- Black (not Hispanic)
- Hispanic ("Hispanic" means someone who is from a Mexican, Mexican American, Chicano, Puerto Rican, Cuban, or other Spanish or Hispanic background.)
- Asian or Pacific Islander ("Asian or Pacific Islander" means someone who is from a Chinese, Japanese, Korean, Filipino, Vietnamese, or other Asian or Pacific Island background.)
- American Indian or Alaskan Native ("American Indian or Alaskan Native" means someone who is from one of the American Indian tribes, or one of the original people of Alaska.)
- ° Other (specify)



88

Students' race/ethnicity was then assigned on the basis of their response. For students who filled in the sixth oval ("Other") or provided illegible information or information that could not be classified, or did not respond at all, race/ethnicity was assigned as determined by school records.³

Derived race/ethnicity could not be determined for students who did not respond to either of the demographic questions and for whom a race/ethnicity designation was not provided by the school.

The details of how race/ethnicity classifications are derived is presented so that the readers can determine the usefulness of the results for their particular uses. It should be noted that a nonnegligible number of students indicated a Hispanic background (e.g., Puerto Rican or Cuban) and indicated that a racial/ethnic category other than Hispanic best described them. These students were classified as Hispanic according to the rules described above. Also, information from the schools did not always correspond to students' descriptions of themselves.

Parents' Highest Level of Education

The variable representing level of parental education is derived from responses to two questions from the set of general background questions. Students were asked to indicate the extent of their mothers' education:

How far in school did your mother go?

- ° She did not finish high school.
- She graduated from high school.
- She had some education after high school.
- She graduated from college.
- ° I don't know.

Students were asked a similar question about their fathers' education level:

How far in school did your father go?

- ° He did not finish high school.
- He graduated from high school.
- He had some education after high school.
- He graduated from college.
- ° I don't know.



)

³ The procedure for assigning race/ethnicity was modified for Hawaii. See the Technical Report for the NAEP 1996 State Assessment Program in Science for details.

This information was combined into one parental education reporting variable through the following procedure. If a student indicated the extent of education for only one parent, that level was included in the data. If a student indicated the extent of education for both parents, the higher of the two levels was included in the data. For students who did not know the level of education for both parents or did not know the level for one parent and did not respond for the other, the parental education level was classified as "I don't know." If the student did not respond for either parent, the student was recorded as having provided no response.

It should be noted that, nationally, approximately one-tenth of eighth graders reported not knowing the education level of either of their parents.

Type of School

Samples for the 1996 state assessment program were expanded to include students attending nonpublic schools (Catholic schools and other religious and private schools) in addition to students attending public schools. The expanded coverage was instituted for the first time in 1994. Samples for the 1990 and 1992 Trial State Assessment programs had been restricted to public school students only. For those jurisdictions meeting pre-established participation rate standards (see earlier section of this appendix), separate results are reported for public schools, for nonpublic schools, and for the combined public and nonpublic school samples. The combined sample for each jurisdiction also contains students attending Bureau of Indian Affairs (BIA) schools and Department of Defense Domestic Dependent Elementary and Secondary Schools (DDESS) in that jurisdiction. These two categories of schools are not included in either the public or nonpublic school samples.

. Note that eighth graders in the DDESS and Department of Defense Dependents Schools (DoDDS)⁴ were assessed in 1996 as separate jurisdictions and reported as jurisdictions with public school samples only.

⁴ The Department of Defense Dependents Schools (DoDDS) refers to overseas schools (i.e., schools outside the United States). Department of Defense Domestic Dependent Elementary and Secondary Schools (DDESS) refers to domestic schools (i.e., schools in the United States). DoDDS and DDESS fourth grades were also assessed in science, for a special report.



Title I Participation

On the basis of available school records, students were classified either as currently participating in a Title I program or receiving Title I services, or as not receiving such services. The classification only refers to the school year when the assessment was administered (i.e., the 1995—96 school year) and is not based on participation in previous years. If the school did not offer any Title I programs or services, all students in that school were classified as not participating.

Free/Reduced-Price School Lunch Program Eligibility

On the basis of available school records, students were classified either as currently eligible for the Department of Agriculture's free/reduced-price lunch program or not. The classification refers only to the school year when the assessment was administered (i.e., the 1995—96 school year) and is not based on eligibility in previous years. If the school did not participate in the program or if school records were not available, all students in that school were classified as "Information not available."

A.3 Guidelines for Analysis and Reporting

This report describes science performance for eighth graders and compares the results for various groups of students within this population — for example, those who have certain demographic characteristics or who responded to a specific background question in a particular way. The report examines the results for individual demographic groups and individual background questions. It does not include an analysis of the relationships among combinations of these subpopulations or background questions.

Drawing Inferences from the Results

Because the percentages of students in these subpopulations and their average scale scores are based on samples — rather than on the entire population of eighth graders in a jurisdiction — the numbers reported are necessarily estimates. As such, they are subject to a measure of uncertainty, reflected in the standard error of the estimate. When the percentages or average scale scores of certain groups are compared, it is essential to take the standard error into account, rather than to rely solely on observed similarities or differences. Therefore, the comparisons discussed in this report are based on statistical tests that consider both the magnitude of the difference between the averages or percentages and the standard errors of those statistics.



One of the goals of the science state assessment program is to estimate scale score distributions and percentages of students in the categories described in A.2 for the overall populations of eighth-grade students in each participating jurisdiction based on the particular samples of students assessed. The use of confidence intervals, based on the standard errors, provides a way to make inferences about the population average scale scores and percentages in a manner that reflects the uncertainty associated with the sample estimates. An estimated sample average scale score ± 2 standard errors approximates a 95 percent confidence interval for the corresponding population average or percentage. This means that one can conclude with approximately 95 percent confidence that the average scale score of the entire population of interest (e.g., all eighth-grade students in public schools in a jurisdiction) is within ± 2 standard errors of the sample average.

As an example, suppose that the average science scale score of the students in a particular jurisdiction's eighth-grade sample were 156 with a standard error of 1.2. A 95 percent confidence interval for the population average would be as follows:

Average
$$\pm$$
 2 standard errors = $156 \pm 2 \times (1.2) = 156 \pm 2.4 =$
 $156 - 2.4$ and $156 + 2.4 = (153.6, 158.4)$

Thus, one can conclude with 95 percent confidence that the average scale score for the entire population of eighth-grade students in public schools in that jurisdiction is between 153.6 and 158.4.

Similar confidence intervals can be constructed for percentages, if the percentages are not extremely large or extremely small. For extreme percentages, confidence intervals constructed in the above manner may not be appropriate, and accurate confidence intervals can be constructed only by using procedures that are quite complicated.

Extreme percentages, defined by both the magnitude of the percentage and the size of the sample from which it was derived, should be interpreted with caution. (The forthcoming Technical Report of the NAEP 1996 State Assessment Program in Science contains a more complete discussion of extreme percentages.)



Analyzing Subgroup Differences in Averages and Percentages

The statistical tests determine whether the evidence, based on the data from the groups in the sample, is strong enough to conclude that the averages or percentages are actually different for those groups in the population. If the evidence is strong (i.e., the difference is statistically significant), the report describes the group averages or percentages as being different (e.g., one group performed higher than or lower than another group), regardless of whether the sample averages or sample percentages appear to be about the same or not. If the evidence is not sufficiently strong (i.e., the difference is not statistically significant), the averages or percentages are described as being not significantly different — again, regardless of whether the sample averages or sample percentages appear to be about the same or widely discrepant. The reader is cautioned to rely on the results of the statistical tests rather than on the apparent magnitude of the difference between sample averages or percentages when determining whether those sample differences are likely to represent actual differences between the groups in the population.

In addition to the overall results, this report presents outcomes separately for a variety of important subgroups. Many of these subgroups are defined by shared characteristics of students, such as their gender or race/ethnicity. Other subgroups are defined by the responses of the assessed students' science teachers to questions in the science teacher questionnaire.

In Chapter 1 of this report, differences between the jurisdiction and the nation were tested for overall science scale score and for each of the fields of science. In Chapter 2, significance tests were conducted for the overall scale score for each of the subpopulations. In Chapters 3 through 6, comparisons were made across subgroups for responses to various background questions.

As an example of comparisons across subgroups, consider the question: Do students who reported discussing studies at home almost every day exhibit higher average science scale scores than students who report never or hardly ever doing so?

To answer the question posed above, begin by comparing the average science scale score for the two groups being analyzed. If the average for the group that reported discussing their studies at home almost every day is higher, it may be tempting to conclude that that group does have a higher science scale score than the group that reported never or hardly ever discussing their studies at home. However, even though the averages differ, there may be no real difference in performance between the two groups in the population because of the uncertainty associated with the estimated average scale scores of the groups in the sample. Remember that the intent is to make a statement about the entire population, not about the particular sample that was assessed. The data from the sample are used to make inferences about the population as a whole.



As discussed in the previous section, each estimated sample average scale score (or percentage) has a degree of uncertainty associated with it. It is therefore possible that if all students in the population (rather than a sample of students) had been assessed or if the assessment had been repeated with a different sample of students or a different, but equivalent, set of questions, the performances of various groups would have been different. Thus, to determine whether there is a real difference between the average scale score (or percentage of students with a certain attribute) for two groups in the population, an estimate of the degree of uncertainty associated with the difference between the scale score averages or percentages of those groups must be obtained for the sample. This estimate of the degree of uncertainty — called the standard error of the difference between the groups — is obtained by taking the square of each group's standard error, summing these squared standard errors, and then taking the square root of this sum.

In a manner similar to that in which the standard error for an individual group average or percentage is used, the standard error of the difference can be used to help determine whether differences between groups in the population are real. The difference between the mean scale score or percentage of the two groups — 2 standard errors of the difference — represents an approximate 95 percent confidence interval. If the resulting interval includes zero, there is insufficient evidence to claim a real difference between groups in the population. If the interval does not contain zero, the difference between groups is statistically significant (different) at the 0.05 level.

As another example, to determine whether the average science scale score of eighth-grade males is higher than that of eighth-grade females in a particular jurisdiction's public schools, suppose that the sample estimates of the average scale scores and standard errors for males and females were as follows:

Group	Average Scale Score Standard Erro				
Males	148	0.9			
Females	146	1.1			

The difference between the estimates of the average scale scores of males and females is two points (148 - 146). The standard error of this difference is

$$\sqrt{0.9^2 + 1.1^2} = 1.4$$

Thus, an approximate 95 percent confidence interval for this difference is

Mean difference ± 2 standard errors of the difference =

$$2 \pm 2 \times (1.4) = 2 \pm 2.8 = 2 - 2.8$$
 and $2 + 2.8 = (-0.8, 4.8)$



Ò

95

The value zero is within this confidence interval, which extends from -0.8 to 4.8 (i.e., zero is between -0.8 and 4.8). Thus, there is insufficient evidence to claim a difference in average science scale score between the populations of eighth-grade males and females in public schools in the hypothetical jurisdiction.

Throughout this report, when the average scale scores or percentages for two groups were compared, procedures like the one described above were used to draw the conclusions that are presented in the text.⁵ If a statement appears in the report indicating that a particular group had a higher (or lower) average scale score than a second group, the 95 percent confidence interval for the difference between groups did not contain zero. An attempt was made to distinguish between group differences that were statistically significant but rather small in a practical sense and differences that were both statistically and practically significant. A procedure based on effect sizes was used. Statistically significant differences that are rather small are described in the text as somewhat higher or somewhat lower. When a statement indicates that the average scale score or percentage of some attribute was not significantly different for two groups, the confidence interval included zero, and thus no difference could be assumed between the groups. The reader is cautioned to avoid drawing conclusions solely on the basis of the magnitude of the difference. A difference between two groups in the sample that appears to be slight may represent a statistically significant difference in the population because of the magnitude of the standard errors. Conversely, a difference that appears to be large may not be statistically significant.

The procedures described in this section, and the certainty ascribed to intervals (e.g., a 95 percent confidence interval), are based on statistical theory that assumes that only one confidence interval or test of statistical significance is being performed. However, in each chapter of this report, many different groups are being compared (i.e., multiple sets of confidence intervals are being calculated). In sets of confidence intervals, statistical theory indicates that the certainty associated with the entire set of intervals is less than that attributable to each individual comparison from the set if considered individually. To hold the certainty level for the set of comparisons at a particular level (e.g., 0.95), modifications (called multiple comparison procedures) must be made to the methods described in the previous section. One such procedure — the Bonferroni method — was used in the analyses described in this report to form confidence intervals for the differences between groups whenever sets of comparisons were considered.⁶ Using this method, the confidence intervals in the text that are based on sets of comparisons are more conservative than those described on the previous pages. In other words, some comparisons that were individually statistically significant using the methods previously described may not be statistically significant when the Bonferroni method was used to take the number of related comparisons into account.



The procedure described above (especially the estimation of the standard error of the difference) is, in a strict sense, appropriate only when the statistics being compared come from independent samples. For certain comparisons in the report, the groups were not independent. In those cases, a different (and more appropriate) estimate of the standard error of the difference was used.

⁶ Miller, R.G. Simultaneous Statistical Inference. (New York: McGraw-Hill, 1966).

Most of the multiple comparisons in this report pertain to relatively small sets or "families" of comparisons. For example, when comparisons were discussed concerning students' reports of parental education, six comparisons were conducted — all pairs of the four parental education levels. In these situations, Bonferroni procedures were appropriate. However, the maps in Chapter 1 of this report display comparisons between Mississippi and all other participating jurisdictions. The "family" of comparisons in this case was as many as 46. To control the certainty level for a large family of comparisons, the False Discovery rate (FDR) criterion was used. Unlike the Bonferroni procedures which control the familywise error rate (i.e., the probability of making even one false rejection in the set of comparisons), the Benjamini and Hochberg (BH) approach using the FDR criterion controls the expected proportion of falsely rejected hypotheses as a proportion of all rejected hypotheses. Bonferroni procedures may be considered conservative for large families of comparisons.8 In other words, using the Bonferroni method would produce more statistically nonsignificant comparisons than using the BH approach. Therefore, the BH approach is potentially more powerful for comparing Mississippi to all other participating jurisdictions. A more detailed description of the Bonferroni and BH procedures appears in the Technical Report of the NAEP 1996 State Assessment Program in Science.

Statistics with Poorly Estimated Standard Errors

Not only are the averages and percentages reported in NAEP subject to uncertainty, but their standard errors are as well. In certain cases, typically when the standard error is based on a small number of students or when the group of students is enrolled in a small number of schools, the amount of uncertainty associated with the standard errors may be quite large. Throughout this report, estimates of standard errors subject to a large degree of uncertainty are followed by the symbol "!". In such cases, the standard errors — and any confidence intervals or significance tests involving these standard errors — should be interpreted cautiously. Further details concerning procedures for identifying such standard errors are discussed in the Technical Report of the NAEP 1996 State Assessment Program in Science.

Williams, V.S.L., L.V. Jones, and J.W. Tukey. Controlling Error in Multiple Comparisons, with Special Attention to the National Assessment of Educational Progress. (Research Triangle Park, NC: National Institute of Statistical Sciences, December 1994).



Benjamini, Y. and Y. Hochberg. "Controlling the false discovery rate: A practical and powerful approach to multiple testing." Journal of the Royal Statistical Society, Series B, 57(1). (pp. 289-300, 1994).

Minimum Subgroup Sample Sizes

Results for science performance and background variables were tabulated and reported for groups defined by gender, race/ethnicity, parental education, type of school, and participation in federally funded Title I programs and the free/reduced-price school lunch component of the National School Lunch Program. NAEP collects data for five racial/ethnic subgroups (White, Black, Hispanic, Asian/Pacific Islander, and American Indian/Alaskan Native) and four levels of parents' education (Graduated From College, Some Education After High School, Graduated From High School, and Did Not Finish High School) plus the category "I Don't Know."

In many jurisdictions, and for some regions of the country, the number of students in some of these groups was not sufficiently high to permit accurate estimation of performance and/or background variable results. As a result, data are not provided for the subgroups with students from very few schools or for the subgroups with very small sample sizes. For results to be reported for any state assessment subgroup, public school results must represent at least 5 primary sampling units (PSUs) and nonpublic school results must represent at least 6 schools. For results to be reported for any national assessment subgroup, at least 5 PSUs must be represented in the subgroup. In addition, a minimum sample of 62 students per subgroup is required. For statistical tests pertaining to subgroups, the sample size for both groups has to meet the minimum sample size requirements.

The minimum sample size of 62 was determined by computing the sample size required to detect an effect size of 0.5 total-group standard deviation units with a probability of 0.8 or greater. The effect size of 0.5 pertains to the *true* difference between the average scale score of the subgroup in question and the average scale score for the total eighth-grade public school population in the jurisdiction, divided by the standard deviation of the scale score in the total population. If the *true* difference between subgroup and total group mean is 0.5 total-group standard deviation units, then a sample size of at least 62 is required to detect such a difference with a probability of 0.8. Further details about the procedure for determining minimum sample size appear in the *Technical Report of the NAEP 1996 State Assessment Program in Science*.



Describing the Size of Percentages

Some of the percentages reported in the text of the report are given qualitative descriptions. For example, the number of students currently taking a biology class might be described as "relatively few" or "almost all," depending on the size of the percentage in question. Any convention for choosing descriptive terms for the magnitude of percentages is to some degree arbitrary. The descriptive phrases used in the report and the rules used to select them are shown below.

Percentage	Descriptive Term Used in Report
p = 0 $0 8 13 18 22 27 30 36$	None A small percentage Relatively few Less than one fifth About one fifth About one quarter Less than one third About one third Less than half
47 $53 64 71 79 89 p = 100$	About half More than half About two thirds About three quarters A large majority Almost all All



98

APPENDIX B

The NAEP 1996 Science Assessment

The science framework for the 1996 National Assessment of Educational Progress was produced under the auspices of the National Assessment Governing Board through a consensus process. The consensus process, managed by the Council of Chief State School Officers, with the National Center for Improving Science Education and the American Institutes for Research, developed the framework over a ten-month period between October 1990 and August 1991. The following factors guided the process for developing consensus on the science framework:

- The active participation of individuals such as curriculum specialists, science teachers, science supervisors, state supervisors, administrators, individuals from business and industry, government officials, and parents;
- The representation of what is considered essential learning in science, and the recommendation of innovative assessment techniques to probe the critical abilities and content areas;
- The recognition of the lack of agreement on such things as common scope of instruction and sequence, components of scientific literacy, important outcomes of learning, and the nature of overarching themes in science.

While maintaining some conceptual continuity with the 1990 NAEP Science Assessment, the 1996 framework takes into account the current reforms in science education, as well as documents such as the science framework used for the 1991 International Assessment of Educational Progress. In addition, the Framework Steering Committee recommended that a variety of strategies, including the following, be used for assessing students' performance.²



Science Framework for the 1996 National Assessment of Educational Progress. (Washington, DC: National Assessment Governing Board, 1993).

² Ibid.

- Performance tasks that allow students to manipulate physical objects and draw scientific understanding from the materials before them
- Constructed-response questions that provide insights into students' levels
 of understanding and ability to communicate in the sciences as well as
 their ability to generate, rather than simply recognize, information
 related to scientific concepts and their interconnections
- Multiple-choice items that probe students' conceptual understanding and ability to connect ideas in a scientifically sound way

B.1 Percentage of Assessment Time by Domain

The framework for the 1996 science assessment can be described as a two-dimensional matrix. The three fields of science (earth, physical, and life) make up the first dimension and ways of knowing and doing science (conceptual understanding, scientific investigation, and practical reasoning) make up the second dimension. Every question or task in the assessment is classified according to the two major dimensions. There are also two overarching domains — nature of science (that includes nature of technology) and themes (systems, models, and patterns of change).

In addition to describing the content of the assessment, the framework also recommends what percentage of time should be devoted to each field of science, each way of knowing and doing science, the nature of science, and themes.

In this section, each figure describes an element of the framework, and is followed by a table showing the *actual* distribution of assessment time as well as the distribution recommended by the framework. Care was taken to ensure congruence between the proportions actually used in the assessment and those recommended in the assessment specifications. Note that the tables represent all three grades assessed nationally; only grade 8 was assessed at the state level.

Figure B.1 describes the fields of science and Table B.1 shows the actual and recommended distribution of assessment time across each field. The ways of knowing and doing science are outlined in Figure B.2. The distribution of assessment time for this dimension, both actual and recommended, is depicted in Table B.2.





FIGURE B.1

Description of the Three Fields of Science

Earth Science

The earth science content assessed centers on objects and events that are relatively accessible or visible. The concepts and topics covered are solid Earth (lithosphere), water (hydrosphere), air (atmosphere), and the Earth in space. The solid Earth consists of composition; forces that alter its surface; the formation, characteristics and uses of rocks; the changes and uses of soil; natural resources used by humankind; and natural forces within the Earth. Concepts and topics related to water consist of the water cycle; the nature of oceans and their effects on water and climate; and the location of water, its distribution, characteristics, and effect of and influence on human activity. The air is broken down into composition and structure of the atmosphere (including energy transfer); the nature of weather; common weather hazards; and air quality and climate. The Earth in space consists of setting of the Earth in the solar system; the setting and evolution of the solar system in the universe; tools and technology that are used to gather information about space; apparent daily motions of the Sun, the Moon, the planets and the stars; rotation of the Earth about its axis, and the Earth's revolution around the Sun; and tilt of the Earth's axis that produces seasonal variations in the climate.

Physical Science

The physical science component relates to basic knowledge and understanding concerning the structure of the universe as well as the physical principles that operate within it. The major sub-topics probed are matter and its transformations, energy and its transformations, and the motion of things. Matter and its transformations are described by diversity of materials (classification and types and the particulate nature of matter); temperature and states of matter; properties and uses of material (modifying properties, synthesis of materials with new properties); and resource management. Energy and its transformations involve different forms of energy; energy transformations in living systems, natural physical systems, and artificial systems constructed by humans; and energy sources and use, including distribution, energy conversion, and energy costs and depletion. Motion is broken down into an understanding of frames of reference; force and changes in position and motion; action and reaction; vibrations and waves as motion; general wave behavior; electromagnetic radiation; and the interactions of electromagnetic radiation with matter.

Life Science

The fundamental goal of life science is to attempt to understand and explain the nature and function of living things. The major concepts assessed in life science are change and evolution, cells and their functions (not at grade 4), organisms, and ecology. Change and evolution includes diversity of life on Earth; genetic variation within a species; theories of adaptation and natural selection; and changes in diversity over time. Cells and their functions consists of information transfer; energy transfer for the construction of proteins; and communication among cells. Organisms are described by reproduction, growth and development; life cycles; and functions and interactions of systems within organisms. The topic of ecology centers on the interdependence of life — populations, communities, and ecosystems.

SOURCE: Science Framework for the 1996 National Assessment of Educational Progress. (Washington, DC: National Assessment Governing Board, 1993).





TABLE B.1

Distribution of Assessment Time by Field of Science

		Earth	P	nysical		Life
	Actual	Recommended	Actual	Recommended	Actual	Recommended
Grade 4	33%	33%	34%	33%	33%	33%
Grade 8	30%	30%	30%	30%	40%	40%
Grade 12	33%	33%	33%	33%	34%	33%



FIGURE B.2

Description of Knowing and Doing Science

Conceptual Understanding

Conceptual understanding includes the body of scientific knowledge that students draw upon when conducting a scientific investigation or engaging in practical reasoning. Essential scientific concepts involve a variety of information including facts and events the student learns from science instruction and experiences with the natural environment and scientific concepts, principles, laws, and theories that scientists use to explain and predict observations of the natural world.

Scientific Investigation

Scientific investigation probes students' abilities to use the tools of science, including both cognitive and laboratory tools. Students should be able to acquire new information, plan appropriate investigations, use a variety of scientific tools, and communicate the results of their investigations.

Practical Reasoning

Practical reasoning probes students' ability to use and apply science understanding in new, real-world applications.

SOURCE: Science Framework for the 1996 National Assessment of Educational Progress. (Washington, DC: National Assessment Governing Board, 1993).



TABLE B.2

Distribution of Assessment Time by Knowing and Doing Science

Conceptua	onceptual Understanding		Scientific Investigation		l Reasoning	
Actual	Recommended	Actual	Recommended	Actual	Recommended	
45%	45%	38%	45%	17%	10%	
45%	45%	29%	30%	26%	25%	
44%	45%	28%	30%	28%	25%	
	45% 45%	45% 45% 45% 45%	Actual Recommended Actual 45% 45% 38% 45% 45% 29%	Actual Recommended Actual Recommended 45% 45% 38% 45% 45% 45% 29% 30%	Actual Recommended Actual Recommended Actual 45% 45% 38% 45% 17% 45% 45% 29% 30% 26%	



The two overarching dimensions are described and accounted for by Figure B.3 and Table B.3, which describe the nature of science and the themes that transcend the scientific disciplines.



FIGURE B.3

Description of Overarching Domains

The Nature of Science

The nature of science incorporates the historical development of science and technology, the habits of mind that characterize these fields, and methods of inquiry and problem-solving. It also encompasses the nature of technology that includes issues of design, application of science to real-world problems, and trade-offs or compromises that need to be made.

Themes

Themes are the "big ideas" of science that transcend the various scientific disciplines and enable students to consider problems with global implications. The NAEP science assessment focuses on three themes: systems, models, and patterns of change.

- Systems are complete, predictable cycles, structures or processes occurring in natural
 phenomena. Students should understand that a system is an artificial construction
 created to represent, or explain a natural occurrence. Students should be able to identify
 and define the system boundaries, identify the components and their interrelationships
 and note the inputs and outputs to the system.
- Models of objects and events in nature are ways to understand complex or abstract phenomena. As such they have limits and involve simplifying assumptions but also possess generalizability and often predictive power. Students need to be able to distinguish the idealized model from the phenomenon itself and to understand the limitations and simplified assumptions that underlie scientific models.
- Patterns of change involve students' recognition of patterns of similarity and differences, and recognize how these patterns change over time. In addition, students should have a store of common types of patterns and transfer their understanding of a familiar pattern of change to a new and unfamiliar one.

SOURCE: Science Framework for the 1996 National Assessment of Educational Progress. (Washington, DC: National Assessment Governing Board, 1993).





TABLE B.3

Distribution of Assessment Time by Overarching Domains

	Nature o	f Science	Themes			
	Actual	Recommended	Actual*	Recommended		
Grade 4	19%	≥15%	53%	33%		
Grade 8	21%	≥15%	49%	50%		
Grade 12	31%	≥15%	55%	50%		

^{*} Several of the hands-on tasks were classified as themes.

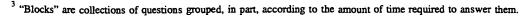
SOURCE: Science Framework for the 1996 National Assessment of Educational Progress. (Washington, DC: National Assessment Governing Board, 1993).

B.2 The Assessment Design

The state science assessment used booklets that were identical to those used at grade 8 for the national assessment. Each student in the state assessment program in science received a booklet containing six sections. Three of these sections were blocks³ of cognitive questions that assessed the knowledge and skills outlined in the framework, and the other three sections were sets of background questions. Two of the three cognitive sections were paper-and-pencil, and the third section consisted of a hands-on task with related questions. In the state assessment at grade 8, students were allowed 30 minutes to complete each cognitive block. (For the national assessment, students at grades 8 and 12 were allowed 30 minutes, while students at grade 4 were given cognitive blocks that each required 20 minutes to complete.)

At each grade level there were 15 different sections or blocks of cognitive questions, but each student's booklet contained only three of these blocks of items. Every block consisted of both multiple-choice and constructed-response questions. Short constructed-response questions required a few words or a sentence or two for an answer (e.g., briefly stating how nutrients move from the digestive system to the tissues) while the extended constructed-response questions generally required a paragraph or more (e.g., outlining an experiment to test the effect of increasing the amount of available food on the rate of increase of the hydra population). Some constructed-response questions also required diagrams, graphs, or calculations. It was expected that students could adequately answer the short constructed-response questions in about 2 to 3 minutes and the extended constructed-response questions in about 5 minutes.

105





Other features were built into the blocks of cognitive questions. Four of the blocks were hands-on tasks in which students were given a set of equipment and asked to conduct an investigation and answer questions relating to the investigation. Every student was assessed on one of these four blocks. A second feature was the inclusion of three theme blocks — one assessing systems, one assessing models, and one assessing patterns of change. For example, students were shown a simplified model of part of the Solar System with a brief description, and then asked a number of questions based on this scenario. Theme blocks were randomly placed in booklets, but not in all booklets. No student received more than one theme block.

Each booklet in the assessment also included three sets of student background questions. The first, consisting of general background questions, asked students about such things as mother's and father's level of education, reading materials in the home, homework, and school attendance. The second, consisting of science background questions, asked students questions about their classroom learning activities such as hands-on exercises, courses taken, use of specialized resources such as computers, and views on the utility and value of science. Students were given five minutes to complete each of these questionnaires. The third set contained five questions about students' motivation to do well on the assessment, their perception of the difficulty of the assessment, and their familiarity with the types of cognitive questions asked. This section took three minutes or less to complete.

Using information gathered from the field test, the booklets were carefully constructed to balance time requirements for the question types in each block. For more information on the design of the assessment, the reader is referred to Appendix C.



B.3 Usage of Question Types

The data in Table B.4 reflect the number of questions by type and by grade level for the 1996 assessment. One hundred and sixty-five multiple-choice (MC), 219 short constructed-response (SCR), and 59 extended constructed-response (ECR) questions make up the assessment, giving a total of 443 unique questions in the pool. Some of these questions were used at more than one grade level; thus, the sum at each grade level is greater than the total number of unique questions. For the state assessment program at grade 8, students responded to subsets (determined by booklet) of 74 multiple-choice questions, 100 short constructed-response questions, and 20 extended constructed-response tasks.

THE NATION'S REPORT CARD	TABLE B.4
1996 State Assessment	Distribution of Items by Question Type

•	Grade 4			Grade 8			Grade 12		
	мс	SRC	ERC	MC	SRC	ERC	MC	SRC	ERC
Grade 4 only	42	57	12						
Grades 4 & 8 overlap	9	16	4	9	16	4			
Grade 8 only				44	58	13			
Grades 8 & 12 overlap				21	26	3	21	26	3
Grade 12 only		`					49	62	27
TOTAL by grade	51	73	16	74	100	20	70	88	30

MC — multiple-choice questions; SRC — short constructed-response questions; ERC — extended constructed-response questions



107

APPENDIX C

Technical Appendix: The Design, Implementation, and Analysis of the 1996 State Assessment Program in Science

C.1 Overview

The purpose of this appendix is to provide technical information about the 1996 state assessment program in science. It describes the design of the assessment and gives an overview of the steps used to implement the program, from the planning stages through the analysis of the data.

This appendix is one of several documents that provide technical information about the 1996 state assessment program. Readers interested in more details are referred to the *Technical Report of the NAEP 1996 State Assessment Program in Science*. Theoretical information about the models and procedures used in NAEP can be found in the special NAEP-related issue of the *Journal of Educational Statistics* (Summer 1992/Volume 17, Number 2) as well as previous national technical reports.

Educational Testing Service (ETS) was awarded the cooperative agreement for the 1996 NAEP programs, including the state assessment program. ETS was responsible for overall management of the programs as well as for development of the overall design, the cognitive questions and questionnaires, data analysis, and reporting. National Computer Systems (NCS) was a subcontractor to ETS on both the national and state NAEP programs. NCS was responsible for printing, distributing, and receiving all assessment materials, and for scanning and scoring the assessments. The National Center for Education Statistics (NCES) awarded a separate cooperative agreement to Westat, Inc., for handling all aspects of sampling and field operations for the national and state assessments for 1996.



Organization of the Technical Appendix

This appendix has the following organization:

- Section C.2 provides an overview of the design of the 1996 state assessment program in science.
- Section C.3 discusses the partially-balanced incomplete block (PBIB) spiral design used to assign cognitive questions to assessment booklets and assessment booklets to students.
- Section C.4 outlines the sampling design used for the 1996 state assessment program.
- Section C.5 summarizes Westat's field administration procedures.
- Section C.6 describes the flow of the data from receipt at NCS through data entry and professional scoring.
- Section C.7 summarizes the procedures used to weight the assessment data and to obtain estimates of the sampling variability of subpopulation estimates.
- Section C.8 describes the initial analyses performed to verify the quality of the data.
- Section C.9 describes the item response theory scales and the overall science composite scale created for the final analyses of the state assessment program data.
- Section C.10 provides an overview of the linking of the scaled results from the state assessment program in science to those from the national assessment.

C.2 Design of the NAEP 1996 State Assessment Program in Science

The design for the state assessment program in science included the following major aspects:

- Participation at the jurisdiction level was voluntary, except for a few jurisdictions for which NAEP has been mandated by the state legislature.
- Students from public and nonpublic schools were assessed. Nonpublic schools included Catholic schools, other religious schools, and private schools. Separate representative samples of public and nonpublic schools were selected in each participating jurisdiction and students were randomly sampled within schools. The size of a jurisdiction's nonpublic school samples was proportional to the percentage of students in that jurisdiction attending such schools.



- The eighth-grade science assessment instruments used for the state assessment program and the national assessment consisted of 15 blocks of questions, of which 4 were hands-on tasks. Each block could contain a mixture of question types — constructed-response or multiple-choice - that was determined by the nature of the task. In addition, the constructed-response questions were of two types: constructed-response questions required students to respond to a question with a few words or a few sentences, while extended constructed-response questions required students to respond to a question with a paragraph or more, sometimes including graphs or calculations. The hands-on tasks were similar to laboratory exercises. Each student was given 2 of the 11 cognitive blocks of questions, and one of the four hands-on blocks.
- A complex form of matrix sampling called a partially balanced incomplete block (PBIB) spiraling design was used. With PBIB spiraling, students in an assessment session received different booklets containing 3 of the 15 blocks. This provided for greater science content coverage without imposing an undue testing burden by administering an identical set of questions to each student.
- Sets of background questions given to the students, the students' science teachers, and the principals or other school administrators provided a variety of contextual information. The background questionnaires for the state assessment program were identical to those used in the national eighth-grade assessment.
- The total assessment time for each student was approximately two hours, including cleanup and collection of materials from hands-on tasks. Each assessed student was assigned a science booklet that contained 3 of the 15 blocks of science questions requiring 30 minutes each (including a hands-on task block in the last position), followed by a 5-minute general background questionnaire, a 5-minute science background questionnaire, and a 3-minute motivation questionnaire. Thirty-seven different booklets were assembled.
- The assessments were administered in the five-week period between January 29 and March 4, 1996. One-fourth of the schools in each jurisdiction were assessed each week throughout the first four weeks. Because of the severe weather throughout much of the country, the fifth week was used for regular testing as well as for makeup sessions.
- Data collection was, by law, the responsibility of each participating jurisdiction. Security and uniform assessment administration were high priorities. Extensive training of state assessment personnel was conducted to assure that the assessment would be administered under standard, uniform procedures. For jurisdictions that had participated in previous NAEP state assessments, 25 percent of both public and nonpublic school assessment sessions were monitored by Westat staff. For the jurisdictions new to NAEP, 50 percent of both public and nonpublic school sessions were monitored.



C.3 Assessment Instruments

The student assessment booklets contained six sections and included both cognitive and noncognitive questions. The assembly of cognitive questions into booklets and their subsequent assignment to assessed students were determined by a matrix sampling design using a variant of a balanced incomplete block design (BIB), with spiraled administration. Each assessed student received a booklet containing 3 of the 15 cognitive blocks according to a design that ensured that each block was administered to a representative sample of students within each jurisdiction. The third cognitive block was always one of the four hands-on blocks; this requirement meant that the BIB was partially balanced (PBIB).

In addition to two 30-minute sections of cognitive questions and the 30-minute performance task section, each booklet included two 5-minute sets of general and science background questions designed to gather contextual information about students, their experiences in science, and their attitudes toward the subject, and one 3-minute section of motivation questions designed to gather information about the student's level of motivation while taking the assessment.

In addition to the student assessment booklets, three other instruments provided data relating to the assessment: a science teacher questionnaire, a school characteristics and policies questionnaire, and an SD/LEP student questionnaire (for students categorized as students with disabilities or with limited English proficiency).

The teacher questionnaire was administered to the science teachers of the eighth-grade students participating in the assessment. The questionnaire consisted of three sections and took approximately 20 minutes to complete. The first section focused on the teacher's general background and experience; the second, on the teacher's background related to science; and the third, on classroom information about science instruction.

The school characteristics and policies questionnaire was given to the principal or other administrator in each participating school and took about 20 minutes to complete. The questions asked about the principal's background and experience, school policies, programs, and facilities, and the demographic composition and background of the students and teachers.

The SD/LEP student questionnaire was completed by the staff member most familiar with any student selected for the assessment who was classified in either of two ways: students with disabilities (SD) had an Individualized Education Plan (IEP) or equivalent special education plan (for reasons other than being gifted and talented); students with limited English proficiency were classified as LEP students. The questionnaire took approximately three minutes to complete and asked about the student and the special programs in which the student participated. It was completed for all selected SD or LEP students regardless of whether or not they participated in the assessment. Selected SD or LEP students participated in the assessment if they were determined by the school to be able to participate, considering the terms of their IEP and accommodations provided by the school or by NAEP.



C.4 The Sampling Design

The sampling design for NAEP is complex, in order to minimize burden on schools and students while maximizing the utility of the data. For further details see the Technical Report for the NAEP 1996 State Assessment Program in Science. The target populations for the state assessment program in science consisted of eighth-grade students enrolled in either public or nonpublic schools. The representative samples of public school eighth graders assessed in the state assessment program came from about 100 schools (per grade) in each jurisdiction. If a jurisdiction had fewer than 100 public schools with a particular grade, all or almost all schools were asked to participate. If a jurisdiction had smaller numbers of students in each school than expected, more than 100 schools were selected for participation. The nonpublic school samples differed in size across the jurisdictions, with the number of schools selected proportional to the nonpublic school enrollment within each jurisdiction. Typically, about 25 nonpublic schools were included for each jurisdiction. The school samples in each state were designed to produce aggregate estimates for the jurisdiction and for selected subpopulations (depending upon the size and distribution of the various subpopulations within the jurisdiction) and also to enable comparisons to be made, at the jurisdiction level, between administration of assessment tasks with monitoring and without monitoring. The public schools were stratified by urbanization, percentage of Black and Hispanic students enrolled, and median household income within the ZIP code area of the school. The nonpublic schools were stratified by type of control (Catholic, private/other religious, other nonpublic), metropolitan status, and enrollment size per grade.

The national and regional results are based on nationally representative samples of eighth-grade students. The samples were selected using a complex multistage sampling design involving the sampling of students from selected schools within selected geographic areas across the country. The sample design had the following stages:

- (1) selection of geographic areas (a county, group of counties, or a metropolitan statistical area);
- (2) selection of schools (public and nonpublic) within the selected areas; and
- (3) selection of students within selected schools.



Each selected school that participated in the assessment, and each student assessed, represent a portion of the population of interest. To make valid inferences from student samples to the respective populations from which they were drawn, sampling weights are needed. Discussions of sampling weights and how they are used in analyses are presented in sections C.7 and C.8.

The state results provided in this report are based on state-level samples of eighth-grade students. The samples of both public and nonpublic school students were selected based on a two-stage sample design that entailed selecting students within schools. The first-stage samples of schools were selected with a probability proportional to the eighth-grade enrollment in the schools. Special procedures were used for jurisdictions with many small schools and for jurisdictions with a small number of schools. As with the national samples, the state samples were weighted to allow for valid inferences about the populations of interest.

The results presented for a particular jurisdiction are based on the representative sample of students who participated in the 1996 state assessment program. The results for the nation and regions of the country are based on the nationally and regionally representative samples of students who were assessed as part of the national NAEP program. Using the national and regional results from the 1996 national assessment was necessary because of the voluntary nature of the state assessment program. Because not every state participated in the program, the aggregated data across states did not necessarily provide representative national or regional results.

In most jurisdictions, up to 30 students were selected from each school, with the aim of providing an initial sample size of approximately 3,000 public school students per jurisdiction for the eighth grade. The student sample size of 30 for each school was chosen to ensure that at least 2,000 public school students participated from each jurisdiction, allowing for school nonresponse, exclusion of students, inaccuracies in the measures of enrollment, and student absenteeism from the assessment. In jurisdictions with fewer schools, larger numbers of students per school were often required to ensure initial samples of roughly 3,000 students. In certain jurisdictions, all eligible eighth graders were targeted for assessment. Jurisdictions were given the option to reduce the expected student sample size in order to reduce testing burden and the number of multiple-testing sessions for participating schools. At grade 8, four jurisdictions (Alaska, Delaware, Hawaii, and Rhode Island) elected to exercise this option. Using this option can involve compromises such as higher standard errors and accompanying loss of precision.



In order to provide for wider inclusion of students with disabilities and limited English proficiency, the 1996 state assessments both in mathematics and science involved dividing the sample of students at each grade level into two subsamples, referred to as S1 and S2. S1 provided continuity with the 1992 mathematics assessment and thus allowed for the reporting of performance over time by using the same exclusion criteria for students with disabilities and limited English proficiency as was used in that assessment. S2 provided for wider inclusion of students with disabilities and limited English proficiency by incorporating new exclusion rules.

The NAEP 1996 science assessment was developed using a new framework, and therefore does not include reporting of performance over time. However, in order to make the sample design identical for both subjects at the state level, both S1 and S2 were included. For further discussion, see the NAEP 1996 Science Report Card.

The 1996 national assessment in science used only the more inclusive S2 guidelines for student participation. The national assessments in mathematics and science both involved an additional subsample, S3, in which accommodations were provided for certain students with disabilities or limited English proficiency, again in order to make NAEP more inclusive.

For the national science assessment, scaling and analysis procedures (discussed in sections C.8 through C.10) were applied to all assessed students from S2. For the state science assessment, scaling and analysis procedures were applied to a combination of all assessed students from S2 and students who were not identified as SD or LEP from S1. This combination of segments of the S1 and S2 subsamples maximized the usefulness of available data while allowing for comparisons to the student population in the national sample. This combination, referred to as the "reporting sample," was the sample used to link the state science assessment to the national assessment (see Section C.10), as well as for scaling and reporting.

Additional analyses will be conducted on the national samples to study the effects of changing the exclusion rules and allowing the use of accommodations. Preliminary discussion can be found in the NAEP 1996 Science Report Card and the NAEP 1996 Mathematics Report Card; more detailed discussion will follow in future NAEP publications.



C.5 Field Administration

Administering the 1996 program required collaboration among staff in the participating jurisdictions and schools and the NAEP contractors, especially Westat, the field administration contractor.

Each jurisdiction volunteering to participate in the 1996 state assessment program appointed a state coordinator to serve as liaison between NAEP staff and the participating schools. In addition, Westat hired and trained a supervisor for each jurisdiction and six field managers who worked with groups of jurisdictions. The state supervisors worked with the state coordinators, overseeing assessment activities, training school district personnel to administer the assessment, and coordinating quality control monitoring efforts. Each field manager worked with the state coordinators from seven to eight jurisdictions and the state supervisors assigned to those jurisdictions. An assessment administrator prepared and conducted the assessment session in one or more schools. These individuals were usually school or district staff and were trained by Westat. Westat also hired and trained three to five quality control monitors in each jurisdiction. For jurisdictions that had previously participated in the state assessment program, 25 percent of the public and nonpublic school sessions were monitored. For jurisdictions new to the program, 50 percent of all sessions were monitored. The assessment sessions were conducted during a five-week period beginning in late January 1996.

C.6 Materials Processing, Professional Scoring, and Database Creation

Upon completion of each assessment session, school personnel shipped the assessment booklets and forms to NCS for professional scoring, entry into computer files, and checking. The files were then sent to ETS for creation of the database.

After NCS received all appropriate materials from a school, they were forwarded to the professional scoring area where the responses to the constructed-response question were evaluated by trained staff using guidelines prepared by ETS. Each constructed-response question had a unique scoring guide that defined the criteria to be used in evaluating students' responses. The extended constructed-response questions were evaluated with four- or five-level rubrics. Some of the short constructed-response questions were rated according to three-level rubrics that permit partial credit to be given; other short constructed-response questions were scored as either acceptable or unacceptable.

For the national science assessment and the state assessment program in science, over 4.1 million constructed responses were scored. This figure includes rescoring to monitor interrater reliability. The overall percentage of agreement between scorers for the reliability sample was 93 percent for the tasks in the cognitive blocks and 95 percent for the hands-on tasks.



Data transcription and editing procedures were used to generate the disk and tape files containing various assessment information, including the sampling weights required to make valid statistical inferences about the population from which the state assessment program sample was drawn. Prior to analysis, the data from these files underwent a quality control check at ETS. The files were then merged into a comprehensive, integrated database.

C.7 Weighting and Variance Estimation

A complex sample design was used to select the students who were assessed in each of the participating jurisdictions. The properties of a sample selected through a complex design are very different from those of a simple random sample in which every student in the target population has an equal chance of selection and in which the observations from different sampled students can be considered to be statistically independent of one another. Therefore, the properties of the sample for the complex state assessment program design were taken into account during the analysis of the assessment data.

One way that the properties of the sample design were addressed was by using sampling weights to account for the fact that the probabilities of selection were not identical for all students. All population and subpopulation characteristics based on the state assessment program data used sampling weights in their estimation. These weights included adjustments for school and student nonresponse.

Not only must appropriate estimates of population characteristics be derived, but appropriate measures of the degree of uncertainty must be obtained for those statistics. One component of uncertainty results from sampling variability, which is a measure of the dependence of the results on the particular sample of students actually assessed. Because of the effects of cluster selection (schools are selected first, then students are selected within those schools), observations made on different students cannot be assumed to be independent of each other (and, in fact, are generally positively correlated). As a result, classical variance estimation formulas will produce incorrect results. Thus, a jackknife variance estimation procedure that accounts for the characteristics of the sample was used for all analyses.

Jackknife variance estimation provides a reasonable measure of uncertainty for any statistic based on values observed without error. Statistics such as the percentage of students correctly answering a given question meet this requirement, but other statistics based on estimates of student science performance, such as the average science scale score of a subpopulation, do not. Because each student typically responds to relatively few questions from a particular field of science (e.g., physical or life science), a nontrivial amount of imprecision exists in the measurement of the scale score of a given student. This imprecision adds another component of variability to statistics based on estimates of individual performance.



C.8 Preliminary Data Analysis

After the computer files of student responses were received and merged into an integrated database, all cognitive and noncognitive questions were subjected to an extensive item analysis. For each cognitive question, this analysis yielded the number of respondents, the percentage of responses in each category, the percentage who omitted the question, the percentage who did not reach the question, and the correlation between the question score and the block score. In addition, the item analysis program provided summary statistics for each block of cognitive questions, including a reliability (internal consistency) coefficient. These analyses were used to check the scoring of the questions, to verify that the difficulty level of the questions was appropriate, and to ensure that students had received adequate time to complete the assessment. The results were reviewed by knowledgeable project staff in search of aberrations that might signal unusual results or errors in the database.

The question and block-level analyses were conducted using rescaled versions of the final sampling weights provided by Westat (see Section C.7). The rescaling was implemented for each jurisdiction. The sum of the sampling weights for the public school students within each jurisdiction was constrained to be equal. The same transformation was applied to the weights of the nonpublic school students in that jurisdiction. The sum of the weights for each of the Department of Defense (DoDEA) samples (i.e., DDESS and DoDDS) was constrained to equal the same value as the public school students in other jurisdictions. Using rescaled weights does not alter the value of statistics calculated separately within each jurisdiction. However, for statistics obtained from samples that combine students from different jurisdictions, using rescaled weights results in a roughly equal contribution of each jurisdiction's data to the final value of the estimate. Equal contribution of each jurisdiction's data to the results of the item response theory (IRT) scaling was viewed as a desirable outcome. The original final sampling weights provided by Westat were used in reporting.

Additional analyses that compared the data from the monitored sessions with those from the unmonitored sessions were conducted to determine the comparability of the assessment data from the two types of administrations. Differential item functioning (DIF) analyses were carried out using the national assessment data. DIF analyses identified questions that were differentially difficult for various subgroups, so that these questions could be re-examined for their fairness and their appropriateness for inclusion in the scaling process.



C.9 Scaling the Assessment Questions

The primary analysis and reporting of the results from the state assessment program used item response theory (IRT) scale-score models. Scaling models quantify a respondent's tendency to provide correct answers to the domain of questions that contribute to a scale as a function of a parameter called performance, estimated by a scale score. The scale scores can be viewed as a summary measure of performance across the domain of questions that make up the scale. Three distinct IRT models were used for scaling: three-parameter logistic models for multiple-choice questions; two-parameter logistic models for short constructed-response questions that were scored correct or incorrect; and generalized partial credit models for short and extended constructed-response questions that were scored on a multipoint scale (i.e., greater than two levels).

Three distinct scales were created for the state assessment program in science to summarize eighth-grade students' abilities according to the three defined fields of science (earth, physical, and life). These scales were defined identically to, but separately from, those used for the scaling of the national NAEP eighth-grade science data. Although the questions composing each scale were identical to those used in the national assessment program, the item parameters for the state assessment program scales were estimated from combined public school data from the jurisdictions participating in the state assessment program. Item parameter estimation was carried out on an item calibration subsample. The calibration subsample consisted of a sample drawn from approximately 25 percent sample of all available public school data. To ensure equal representation in the scaling process, each jurisdiction contributed the same number of students to the item calibration sample. Within each jurisdiction, 25 percent of the calibration sample was taken from monitored administrations while the remaining 75 percent came from unmonitored administrations.

Within each scale, the estimates of the empirical item characteristic functions were compared with the theoretical curves to determine how well the IRT model fit the observed data. For correct-incorrect questions, nonmodel-based estimates of the expected proportions of correct responses to each question for students with various levels of scale proficiency were compared with the fitted item response curve. For the short and extended partial-credit constructed-response questions, the comparisons were based on the expected proportions of students with various levels of scale proficiency who achieved each score level. In general, the scaling models fit the question-level results well.



¹ For the creation of scales, schools from the DoDEA jurisdictions are considered nonpublic, so the responses from these students were not included in the item calibration sample.

Using the item parameter estimates, estimates of various population statistics were obtained for each jurisdiction. The NAEP methods use random draws ("plausible values") from estimated proficiency distributions for each student to compute population statistics. Plausible values are not optimal estimates of individual student proficiencies; instead, they serve as intermediate values to be used in estimating population characteristics. Under the assumptions of the scaling models, these population estimates will be consistent, in the sense that the estimates approach the model-based population values as the sample size increases, which would not be the case for population estimates obtained by aggregating optimal estimates of individual performance.

The 1996 science assessment was developed using a new framework. Because it was not appropriate to compare results from the 1996 assessment to those of previous NAEP science assessments, no attempt was made to link or align scores on the new assessment to those of previous assessments. Therefore, it was necessary to establish a new scale for reporting. Earlier NAEP assessments (such as the current mathematics assessment and the 1994 reading assessment) were developed with a cross-grade framework, in which the trait being measured is conceptualized as cumulative across the grades of the assessment. This concept was reflected in the scaling. The score scales developed for these assessments were cross-grade scales on a single 0-500 scale for all three grades in the assessment.

In 1993, the National Assessment Governing Board (NAGB) determined that future NAEP assessments should be developed using within-grade frameworks. This removes the constraint that the trait being measured is cumulative, and there is no need for overlap of questions across grades. Consistent with this view, NAGB also declared that scaling be performed within-grade. Any items which happened to be the same across grades in the assessment were scaled separately for each grade, thus allowing common items, potentially, to function differently in the separate grades. The 1994 NAEP history and geography assessments were developed and scaled within-grade. After scaling, the scales were aligned so that grade 8 had a higher mean than did grade 4, and grade 12 had a higher mean than grade 8. The results were reported on a final 0-500 scale that looked similar to those used in mathematics and reading, in spite of the differences in development and scaling. This definition of the reporting scale was a source of potential confusion and misinterpretation.

The 1996 science assessment was also developed and scaled using within-grade procedures. A new reporting metric was adopted to differ from the 0-to-500 reporting scales used in other NAEP subject areas in order to minimize confusion with other common test scales and to discourage cross-grade comparisons. For each grade in the national assessment, the mean for each field of science was set at 150 and the standard deviation was set at 35. First, the reporting metric was developed using data from the national assessment program; the results for the state assessment program were then linked to that scale using procedures described in Section C.10.



In addition to the plausible values for each scale, a composite of the three fields of science scales was created as a measure of overall science performance; as for the individual fields of science scales, the mean of the composite scale was set to 150 with a standard deviation of 35.² This composite was a weighted average of the plausible values for the three fields of science scales. The scales were weighted proportionally to the relative importance assigned to each field of science in the science framework (see Table B.1). The definition of the composite for the state assessment program was identical to that used for the national eighth-grade science assessments.

C.10 Linking the State Results to the National Results

A major purpose of the state assessment program was to allow each participating jurisdiction to compare its 1996 results with those for the nation as a whole and with those for the region of the country where it is located. For meaningful comparisons to be made between each jurisdiction and the relevant national sample, results from these two assessments had to be expressed in terms of a similar system of scale units.

The results from the state assessment program were linked to those from the national assessment through linking functions determined by comparing the results for the aggregate of all students assessed in the state assessment program with the results for eighth-grade students within the National Linking Sample of the national NAEP. The National Linking Sample of the national NAEP is a representative sample of the population of all grade-eligible public school students within the aggregate of 43 participating states and the District of Columbia. (Guam and the two DoDEA jurisdictions were not included in the National Linking Sample.) Specifically, the National Linking Sample for science consisted of all eighth-grade students in public schools in the states and the District of Columbia who were assessed in the national cross-sectional science assessment.

A linear equating within each field of science scale was used to link the results of the state assessment program to the national assessment. For each scale, the adequacy of the linear equating was evaluated by comparing the distribution of science scale scores based on the aggregation of all assessed students at each grade from the participating states and the District of Columbia with the equivalent distribution based on the students in the National Linking Sample. In the estimation of these distributions, the students were weighted to represent the target population of public school students in the specified grade in the aggregation of the states and the District of Columbia. If a linear equating were adequate, the distribution for the aggregate of states and the District of Columbia and that for the National Linking Sample would have, to a close approximation, the same shape in terms of the skewness, kurtosis, and higher moments of the distributions. The only differences in the distributions allowed by linear equating would be in the means and variances. Generally, this has been found to be the case.

Thus, each field of science scale was linked by matching the scale mean and standard deviation of the scale scores across all students in the state assessment (excluding Guam and the two DoDEA jurisdictions) to the corresponding mean and standard deviation across all students in the National Linking Sample.



² The national average of students in public and nonpublic schools combined is 150. The national average seen in the tables in this report is based on the average for public schools only (148).

APPENDIX D

Teacher Preparation

Because teachers are key to improving science education, their background and professional development should be examined. Eighth-grade science teachers completed questionnaires about their background and training, including their experience, certification, undergraduate and graduate course work in science, and involvement in pre-service education.

Consistent with procedures used throughout this report, the student was the unit of analysis. That is, the science teachers' responses were linked to their students, and the data reported are the percentages of students taught by these teachers rather than the percentages of teachers.

The tables in Appendix D represent only a few of the questions in the teacher questionnaire, and this small selection can give only a sketchy profile of the teachers. A report scheduled to appear in early 1998 will explore more of the questions related to school and classroom policy and practices and should give a better picture of the nation's teachers.



The interested reader can obtain additional information on teachers' characteristics and qualifications and the conditions under which they teach in SASS by State (NCES 96-312) from the 1993-94 Schools and Staffing Survey. URL: http://www.ed.gov/NCES/pubs/96312.html.



Public School Teachers' Reports on Their Highest Level of Education

What is the highest academic	Mississippi	Southeast	Nation
degree you hold? 		Percentage	
Bachelor's degree	65 (4.2)	45 (7.0)	55 (4.2)
Master's degree	30 (4.3)	37 (5.8)	34 (4.0)
Education specialist's or professional diploma	5 (1.6)	15 (9.2)	9 (3.4)
Doctorate or professional degree	0 (****)	3 (1.6)	1 (0.5)

The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). **** Standard error estimates cannot be accurately determined.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



TABLE D.2

Public School Teachers' Reports on Their Major Fields of Study

What were your major fields of	Mississippi	Southeast	Nation	
etudy? (multiple responses possible)		Percentage		
Undergraduate				
Education or elementary education	52 (3.7)	54 (6.7)	38 (3.7)	
Secondary education	32 (3.8)	34 (6.4)	41 (4.5)	
Science education	37 (3.1)	31 (7.7)	36 (4.2)	
Life science	26 (3.7)	35 (4.3)	43 (5.1)	
Physical science	17 (3.4)	22 (5.3)	19 (5.0)	
Earth science	13 (2.5)	20 (5.6)	22 (4.1)	
Other ·	28 (3.4)	33 (5.9)	35 (4.7)	
Graduate				
Education or elementary education	30 (4.0)	22 (4.5)	27 (3.8)	
Secondary education	16 (2.8)	17 (5.8)	26 (3.4)	
Science education	24 (3.4)	19 (5.5)	28 (5.0)	
Life science	12 (2.1)	9 (3.1)	10 (1.8)	
Physical science	9 (2.5)	6 (3.1)	5 (1.5)	
Earth science	4 (1.6)	5 (2.3)	9 (2.4)	
Other	22 (3.1)	35 (9.2)	42 (4.5)	
No graduate study	28 (3.6)	27 (6.0)	13 (2.4)	

The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.





Public School Teachers' Reports on Their Teaching Certification

Mississippi	Southeast	Nation
	Percentage	

What type of teaching certification do you have in this state in your main assignment field?			
l don't have a certificate in my main assignment field.	1 (0.8)	2 (1.3)	1 (0.5)
Certification by an accreditation body other than the state	0 (****)	0 (****)	0 (****)
Temporary, provisional, or emergency state certificate	6 (2.1)	2 (****)	4 (1.3)
Probationary state certificate (Initial certificate)	2 (1.3)	3 (1.5)	3 (1.3)
Regular or standard state certificate	74 (3.7)	62 (7.4)	79 (3.5)
Advanced professional certificate	17 (2.8)	32 (6.7)	13 (3.0)
Do you have teaching certification in any of the following areas that is recognized by the state in which you teach? (multiple responses possible) Elementary or middle/junior high school education	78 (4.0)	68 (7.1)	66 (5.9)
Elementary science	16 (3.3)	11 (4.4)	25 (4.3)
Middle/junior high school or secondary science	80 (4.1)	89 (5.0)	95 (1.6)
Other	33 (6.4)	62 (7.0)	51 (6.3)

The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). **** Standard error estimates cannot be accurately determined.

error estimates cannot be accurately determined.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



123

, t



Public School Teachers' Reports on Years of Teaching Experience

Counting this year, how many years	Mississippi	Southeast	Nation
have you		Percentage	

aught at either the elementary r secondary level? ¹			
2 years or less	13 (3.1)	14 (3.5)	9 (2.2)
3-5 years	· 17 (3.1)	11 (3.8)	9 (1.7)
6-10 years	11 (2.4)	20 (7.0)	22 (3.2)
11-24 years	40 (3.9)	38 (5.9)	36 (4.1)
25 years or more	19 (3.1)	17 (3.8)	24 (3.2)
aught science? ²			
2 years or less	16 (3.2)	18 (3.5)	13 (2.4)
3-5 years	20 (3.5)	13 (3.4)	11 (2.2)
6-10 years	16 (3.1)	28 (7.7)	30 (3.2)
11-24 years	39 (4.1)	33 (8.0)	26 (3.4)
25 years or more	10 (2.2)	9 (2.0)	20 (3.0)

The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). Teachers were instructed to include part-time teaching experience. Teachers were instructed to include full-time and part-time assignments, but not substitute assignments.

SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment



TABLE D.5

Public School Teachers' Reports on Recent Course Taking

During the last two years, how	Mississippi	Southeast	Nation
many college or university courses have you taken in science or science education?		Percentage	,

None	58 (4.6)	54 (7.0)	59 (3.4)
One	15 (3.7)	16 (3.1)	14 (2.8)
Two	10 (2.5)	19 (5.2)	11 (2.4)
Three or more	17 (3.3)	11 (2.8)	16 (2.8)

The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.





Public School Teachers' Reports on Professional Development Activities

Mississippi	Southeast	Nation
	Percentage	

During the past two years, have you taken college or university courses in any of the following?		·	
Methods of teaching science	17 (3.1)	16 (5.1)	12 (2.2)
Biology/life science	17 (3.4)	15 (5.2)	14 (2.7)
Chemistry	11 (3.1)	13 (5.0)	6 (1.7)
Physics	7 (2.0)	18 (5.1)	8 (1.8)
Earth science	9 (2.6)	16 (5.2)	9 (2.0)
During the past five years, have you taken courses or participated in professional development activities in any of the following? Use of computers for data acquisition	34 (4.0)	57 (6. 2)	50 (4.6)
Use of computers for data analysis	32 (4.4)	48 (6.2)	54 (4.4)
Use of multimedia for science education	31 (3.2)	58 (7.9)	54 (4.5)
Laboratory management or safety	23 (3.6)	30 (8.2)	28 (3.8)
Integrated science instruction	60 (4.3)	49 (6.6)	46 (4.2)

The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



125



Public School Teachers' Reports on Professional Development

Ouring the last year, how much time	Mississippi	Southeast	Nation
rofessional development rorkshops or seminars in science r science education?		Percentage	
None	9 (2.0)	4 (1.5)	8 (2.5)
None Less than six hours	9 (2.0) 18 (3.6)	4 (1.5) 19 (6.9)	8 (2.5) 16 (4.2)
*****	'''	, , ,	
Less than six hours	18 (3.6)	19 (6.9)	16 (4.2)

The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within \pm 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



TABLE D.8

Public School Teachers' Reports on Membership in Professional Societies

Mississippi	Southeast	Nati <i>o</i> n
	Percentage	
41 (4.3)	52 (9.0) 48 (9.0)	57 (4.5) 43 (4.5)
		Percentage 41 (4.3) 52 (9.0)

The standard errors of the statistics appear in parentheses. It can be said with about 95 percent confidence that, for each population of interest, the value for the entire population is within ± 2 standard errors of the estimate for the sample. In comparing two estimates, one must use the standard error of the difference (see Appendix A for details). SOURCE: National Center for Education Statistics, National Assessment of Educational Progress (NAEP), 1996 Science Assessment.



ACKNOWLEDGMENTS

This report is the culmination of the effort of many individuals who contributed their considerable knowledge, experience, and creativity to the NAEP 1996 science assessment. The NAEP 1996 science assessment was a collaborative effort among staff from the National Center for Education Statistics (NCES), the National Assessment Governing Board (NAGB), Educational Testing Service (ETS), Westat, Inc., and National Computer Systems (NCS). In addition, the program benefited from the contributions of hundreds of individuals at the state and local levels — governors, chief state school officers, state and district test directors, state coordinators, and district administrators — who tirelessly provided their wisdom, experience, and hard work. Most importantly, NAEP is grateful to the over 109,000 students and the teachers and administrators from 4,400 schools in 47 jurisdictions who made the science assessment possible.

The NAEP 1996 science assessment was funded through NCES, in the Office of Educational Research and Improvement of the U.S. Department of Education. The Commissioner of Education Statistics, Pascal D. Forgione, and the NCES staff — Sue Ahmed, Peggy Carr, Arnold Goldstein, Steven Gorman, Larry Ogle, Gary W. Phillips, Sharif Shakrani, Maureen Treacy — worked closely and collegially with the authors to produce this report. The authors were also provided guidance by the members of the National Assessment Governing Board and NAGB staff. In particular, the authors are indebted to Arnold Goldstein of NCES for his daily efforts to coordinate the activities of the many people who contributed to this report.

The NAEP project at ETS is housed in the Center for the Assessment of Educational Progress under the direction of Paul Williams. The NAEP 1996 assessments were directed by Stephen Lazer and John Mazzeo. Tom Corley, Lee Jones, Tim Ligget, Beth Nichols, Christine O'Sullivan, Amy Pearlmutter, Will Pfeiffenberger, Mario Yepes-Baraya, and Ann Marie Zolandz worked with the Science Instrument Development committee to develop the assessment instrument. Sampling and data collection activities were conducted by Westat under the direction of Rene Slobasky, Nancy Caldwell, Keith Rust, and Dianne Walsh. Printing, distribution, scoring, and processing activities were conducted by NCS under the direction of Brad Thayer, Patrick Bourgeacq, Charles Brungardt, Jay Happel, Mathilde Kennel, Linda Reynolds, and Brent Studer.



The complex statistical and psychometric activities necessary to report results for the NAEP 1996 Science Assessment were directed by Nancy Allen, John Barone, James Carlson, and Juliet Shaffer. The analyses presented in this report were led by John Donoghue and Steven Isham with assistance from Spencer Swinton, Lois Worthington, Inge Novatkoski, Kate Pashley, David Freund, and Norma Norris.

Laura Jerry was responsible for the development and creation of the computer-generated reports, with assistance from Xiaohui Wang, Laura Jenkins, Phillip Leung, Inge Novatkoski, Bruce Kaplan, and Alfred Rogers. Two of the reports were skillfully produced by Karen Damiano. A large group of NAEP staff at ETS checked the data, text, and tables. Debbie Kline coordinated the technical appendices. The overall production efforts were completed by Carol Errickson, Barbette Tardugno, Loretta Casalaina, Kelly Gibson (cover design), Sharon Davis-Johnson, and Alice Kass. Editorial assistance was provided by Walt Brower and by John Calderone of Aspen Systems. The World Wide Web version of the state reports was produced by Philip Leung and Pat O'Reilly with assistance from Debbie Kline, Craig Pizzuti, and Christine Zelenak.

Many thanks are due to the numerous reviewers, both internal and external to NCES and ETS. The comments and critical feedback of the following reviewers are reflected in this report: Sue Ahmed, Peggy Carr, Mary Frase, Arnold Goldstein, Andrew Kolstad, Michael Ross, and Shi-Chang Wu of NCES; Rolf Blank of CCSSO; Audrey Champagne of the State University of New York in Albany; Michelle Leon of the Connecticut Department of Education; Will Pfeiffenberger of ETS; Senta Raizen of the National Center for Improving Science Education; and Mistilina Sato of Stanford University.



NAEP 1996 Science Instrument Development Committee

An Instrument Development Committee was convened to oversee the development of items and scoring rubrics. Committee members wrote assessment exercises and ensured that the instrument adhered to the assessment framework and specifications. In addition, the committee made certain that the instrument was developmentally appropriate for each grade and that it was relevant to curricular and instructional goals. The members are to be commended for their diligence and dedication to the lengthy process of producing the instrument:

Gail Baxter, University of Michigan
Ron Bonnstetter, University of Nebraska
Audrey Champagne, State University of New York at Albany
Richard Clark, Minnetonka, Minnesota
Sally Crissman, Shady Hill School, Cambridge, Massachusetts
Pat Dung, Los Angeles Educational Partnership
Michael Johnson, Science Skills Center
Michael Jojola, Isleta, New Mexico
Clifton Poodry, University of California at Santa Cruz
Senta Raizen, National Center for Improving Science Education
Douglas Reynolds, Rensselaer, New York
Realista Rodriguez, Stuart High School, Falls Church, Virginia
Mistilina Sato, Stanford University
Gerald Weaver, University City High School, Philadelphia, Pennsylvania
Mary Louise Bellamy, National Association of Biology Teachers



DONE

ERRATA NOTICE

Date: December 29, 1997

To: Participants in the NAEP 1996 Science State Assessment

From: Nada Ballator

Center for the Assessment of Educational Progress at Educational Testing Service

1-800-223-0267

Re: Replacement pages attached for NAEP 1996 Science State Reports, correcting

error in national and regional data in Table 6.2 and associated text

An error was recently discovered in the *national and regional* data presented in Table 6.2 of the 1996 science state reports. *For all states and jurisdictions, the data are correct*; however, incorrect national data made it necessary to recompute comparisons between state and national results. The error involved the student background item, "About how many books are in your home?" which is reported in the *NAEP 1996 Science State Report* in Table 6.2, as well as in the bullets comparing your jurisdiction with the nation.

Attached to this memo are the two corrected pages to insert into your printed reports. If you received camera-ready copy of the NAEP 1996 science state report, we have also enclosed pages for insertion there. The pages are for Chapter 6 in the section on "Literacy Materials in the Home" which includes Table 6.2; they contain revised comparisons to national data, and revised national and regional data in the table. We apologize for the publication of inaccurate data, and for the extra effort its correction will cause you.

The state science reports appear on the NCES web site (http://nces.ed.gov/naep). All affected reports on the web were corrected on December 17. There is now a Revised logo beside the reports on the Index of Results and Summary Data web page (http://nces.ed.gov/naep/rsdindex.shtml) and on the Current Assessment Results web page (http://nces.ed.gov/naep/naep1996.html), and an Errata Notice containing a brief description of the repair on the NAEP 1996 Science State Reports web page (http://nces.ed.gov/naep/96state/97499.shtml).

Also on the web site, the student data tables for national science results for public schools have been revised. On the web page for NAEP 1996 Summary Data Tables, Student Data (http://nces.ed.gov/naep/tables96/index.shtml), you will see an Errata Notice describing the repair. Please alert anyone who may be using national 1996 science student data to this revision concerning the raw variable, "How many books are in your home," and the derived variable HOMEEN3, "Home environment - Articles (of 4) in home."

We very much regret the extra work that this error may have necessitated in your jurisdiction; we will redouble our efforts to prevent such things happening again.



NCES 97-499 MS







U.S. DEPARTMENT OF EDUCATION

Office of Educational Research and Improvement (OERI) Educational Resources Information Center (ERIC)



NOTICE

REPRODUCTION BASIS

	This document is covered by a signed "Reproduction Release (Blanket)" form (on file within the ERIC system), encompassing all or classes of documents from its source organization and, therefore, does not require a "Specific Document" Release form.
g	This document is Federally-funded, or carries its own permission to reproduce, or is otherwise in the public domain and, therefore, may be reproduced by ERIC without a signed Reproduction Release form (either "Specific Document" or "Blanket").

