

The image shows the front cover of the book. The title is prominently displayed at the top. Below the title, there is a subtitle and the name of the committee. The cover has a clean, professional design with a light background and dark text.

## **Cognitive Aspects of Survey Methodology: Building a Bridge Between Disciplines**

Committee on National Statistics, National Research Council

ISBN: 0-309-59179-1, 192 pages, 8.5 x 11, (1984)

**This PDF is available from the National Academies Press at:**  
<http://www.nap.edu/catalog/930.html>

Visit the [National Academies Press](http://www.nap.edu) online, the authoritative source for all books from the [National Academy of Sciences](http://www.nap.edu), the [National Academy of Engineering](http://www.nap.edu), the [Institute of Medicine](http://www.nap.edu), and the [National Research Council](http://www.nap.edu):

- Download hundreds of free books in PDF
- Read thousands of books online for free
- Explore our innovative research tools – try the “[Research Dashboard](#)” now!
- [Sign up](#) to be notified when new books are published
- Purchase printed books and selected PDF files

**Thank you for downloading this PDF. If you have comments, questions or just want more information about the books published by the National Academies Press, you may contact our customer service department toll-free at 888-624-8373, [visit us online](#), or send an email to [feedback@nap.edu](mailto:feedback@nap.edu).**

**This book plus thousands more are available at <http://www.nap.edu>.**

Copyright © National Academy of Sciences. All rights reserved.  
Unless otherwise indicated, all materials in this PDF File are copyrighted by the National Academy of Sciences. Distribution, posting, or copying is strictly prohibited without written permission of the National Academies Press. [Request reprint permission for this book](#).

---

---

# **Cognitive Aspects of Survey Methodology: Building a Bridge Between Disciplines**

**Report of the Advanced Research Seminar on Cognitive Aspects of Survey  
Methodology**

Thomas B. Jabine, Miron L. Straf, Judith M. Tanur, and Roger Tourangeau Editors

Committee on National Statistics  
Commission on Behavioral and Social Sciences and Education  
National Research Council

NATIONAL ACADEMY PRESS  
Washington, D.C. 1984

---

NOTICE: The project that is the subject of this report was approved by the Governing Board of the National Research Council, whose members are drawn from the councils of the National Academy of Sciences, the National Academy of Engineering, and the Institute of Medicine. The members of the committee responsible for the report were chosen for their special competences and with regard for appropriate balance.

This report has been reviewed by a group other than the authors according to procedures approved by a Report Review Committee consisting of members of the National Academy of Sciences, the National Academy of Engineering, and the Institute of Medicine.

The National Research Council was established by the National Academy of Sciences in 1916 to associate the broad community of science and technology with the Academy's purposes of furthering knowledge and of advising the federal government. The Council operates in accordance with general policies determined by the Academy under the authority of its congressional charter of 1863, which established the Academy as a private, nonprofit, self-governing membership corporation. The Council has become the principal operating agency of both the National Academy of Sciences and the National Academy of Engineering in the conduct of their services to the government, the public, and the scientific and engineering communities. It is administered jointly by both Academies and the Institute of Medicine. The National Academy of Engineering and the Institute of Medicine were established in 1964 and 1970, respectively, under the charter of the National Academy of Sciences.

Available from:

Committee on National Statistics

National Research Council

2101 Constitution Avenue, N.W.

Washington, D.C. 20418

Printed in the United States of America

## CONTENTS

	PARTICIPANTS AND GUESTS, ADVANCED RESEARCH SEMINAR ON COGNITIVE ASPECTS OF SURVEY METHODOLOGY	v
	COMMITTEE ON NATIONAL STATISTICS	vii
	PREFACE	ix
CHAPTER 1	THE SEMINAR	1
	Background	1
	Introduction	3
	Surveys as a Vehicle for Cognitive Research	6
	Improving Survey Methods	10
	Issues for the National Health Interview Survey	21
CHAPTER 2	AFTER THE SEMINAR	25
	Laboratory-Based Research on the Cognitive Aspects of Survey Methodology	26
	Center for Health Statistics (Monroe Sirken and Robert Fuchsberg)	
	Cognitive Processes in Survey Responding: Project Summaries	35
	Roger Tourangeau, William Salter, Roy D'Andrade, Normal Bradburn, and associates	
	The Intersection of Personal and National History	38
	Howard Schuman and Philip Converse	
	A Proposal for the Development of a National Memory Inventory	44
	Endel Tulving and S. James Press	
	Protocol Analysis of Responses to Survey Recall Questions	61
	Elizabeth Loftus	
	Thoughts and Research on Estimates About Past and Future Behavior	65
	Lee Ross	
	Outreach Activities	69

---

CONTENTS	iv
APPENDIX A BACKGROUND PAPERS	71
Cognitive Sciences and Survey Methods	73
Roger Tourangeau	
Potential Contributions of Cognitive Research to Survey Questionnaire Design	101
Norman Bradburn and Catalina Danis	
Record Checks for Sample Surveys	130
Kent Marquis	
APPENDIX B DESIGNING AND BUILDING THE BRIDGE	149
APPENDIX C BACKGROUND MATERIALS FOR THE SEMINAR	157
APPENDIX D BIOGRAPHICAL SKETCHES OF PARTICIPANTS	165
INDEX	171

---

## ADVANCED RESEARCH SEMINAR ON COGNITIVE ASPECTS OF SURVEY METHODOLOGY

### PARTICIPANTS

JUDITH M. TANUR (Chair), Department of Sociology, State University of New York, Stony Brook  
NORMAN M. BRADBURN, National Opinion Research Center, Chicago  
PHILIP E. CONVERSE, Institute for Social Research, University of Michigan  
ROY G. D'ANDRADE, Department of Anthropology, University of California, San Diego  
STEPHEN E. FIENBERG, Department of Statistics, Carnegie-Mellon University  
ROBERT FUCHSBERG, National Center for Health Statistics, U.S. Department of Health and Human Services  
THOMAS B. JABINE, Consultant, Committee on National Statistics  
WILLETT KEMPTON, Institute of Psychology, University of Washington  
ALBERT MADANSKY, Graduate School of Business, University of Chicago  
ROBERT MANGOLD, Bureau of the Census, U.S. Department of Commerce  
KENT MARQUIS, Bureau of the Census, U.S. Department of Commerce  
ANDREW ORTONY, Center for the Study of Reading, University of Illinois  
S. JAMES PRESS, Department of Statistics, University of California, Riverside  
LEE ROSS, Department of Psychology, Stanford University  
WILLIAM J. SALTER, Bolt, Beranek, and Newman, Inc.  
HOWARD SCHUMAN, Institute for Social Research, University of Michigan  
MONROE G. SIRKEN, National Center for Health Statistics, U.S. Department of Health and Human Services  
ANNE M. SPRAGUE, Administrative Secretary, Committee on National Statistics  
MIRON L. STRAF, Research Director, Committee on National Statistics  
ROGER TOURANGEAU, National Opinion Research Center, New York  
ENDEL TULVING, Department of Psychology, University of Toronto  
Biographical sketches of seminar participants appear in [Appendix D](#).

### GUESTS

MURRAY ABORN, National Science Foundation  
EARL F. BRYANT, National Center for Health Statistics, U.S. Department of Health and Human Services

JACOB J. FELDMAN, National Center for Health Statistics, U.S. Department of Health and Human Services  
GARY G. KOCH, Department of Biostatistics, University of North Carolina  
DAWN W. NELSON, Bureau of the Census, U.S. Department of Commerce  
SARA B. NERLOVE, National Science Foundation  
ROBERT PEARSON, Social Science Research Council  
PHILLIP STONE, Department of Psychology and Social Relations, Harvard University  
KATHERINE K. WALLMAN, Council of Professional Associations on Federal Statistics  
THOMAS C. WALSH, Bureau of the Census, U.S. Department of Commerce

## COMMITTEE ON NATIONAL STATISTICS

STEPHEN E. FIENBERG (Chair), Department of Statistics, Carnegie-Mellon University

LEO BREIMAN, Department of Statistics, University of California, Berkeley

JOEL E. COHEN, Laboratory of Populations, The Rockefeller University

WAYNE A. FULLER, Department of Statistics, Iowa State University

F. THOMAS JUSTER, Institute of Social Research, University of Michigan

GARY G. KOCH, Department of Biostatistics, University of North Carolina

PAUL MEIER, Department of Statistics, University of Chicago

JANE A. MENKEN, Office of Population Research, Princeton University

LINCOLN E. MOSES, Department of Statistics, Stanford University

JOHN W. PRATT, Graduate School of Business, Harvard University

CHRISTOPHER A. SIMS, Department of Economics, University of Minnesota

BURTON H. SINGER, Department of Statistics, Columbia University

COURTENAY M. SLATER, CEC Associates, Washington, D.C.

JUDITH M. TANUR, Department of Sociology, State University of New York at Stony Brook

EDWIN D. GOLDFIELD, Executive Director

MIRON L. STRAF, Research Director



About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.

## PREFACE

Cross-disciplinary collaboration is difficult. Practitioners from different disciplines could be said to live in different cultures. They see different things as important (or trivial); they use different research techniques; they have different backgrounds of understood and taken-for-granted knowledge; and they have knowledge of different specific languages in which the specialized terms of one discipline may be meaningless to another--or worse, have well-defined but differing meanings in the other discipline. Thus the comingling of disciplines can give rise to culture shock, resulting in either bewilderment and apathy or in ethnocentrism, but in either case yielding communication failure and resulting frustration. Yet most of us agree that collaboration between practitioners trained in different disciplines can engender research projects that have exceptional promise, both for enriching the cultures of the parent disciplines and for creating a hybrid culture that attains its own viability and establishes its own research tradition. But how are the cultural barriers between disciplines to be overcome? The following pages report on what we consider an experiment in encouraging cross-disciplinary collaboration. While the Advanced Research Seminar on Cognitive Aspects of Survey Methodology--CASM--was surely not an experiment in the statistical sense of the term, consonant with the cross-disciplinary aims of the project, we shall continue to call it an experiment.

The results of the experiment are not yet in. Indeed we shall have to wait some years to see whether a whole new field arises; only a few projects are carried out; or the seminar becomes merely a fond memory for the participants, little influencing their own work or that of others. But the early prospects for our results seem bright. The body of this report sketches many ideas for collaborative research that arose during the seminar and outlines some more fully formed projects that CASM participants plan to undertake. We hope that these are only the first artifacts of the comingled culture created by the seminar. Thus, while we are not yet able to state definitively whether our experiment was a success or something less, we can see already that some progress has been made in disciplinary cross-fertilization. We present these results in the hope that readers will see the new field as one in which they might wish to carry out research and perhaps even take up some of our ideas. We have been reporting on the seminar at various professional meetings

and in various journals, as outlined at the end of [Chapter 2](#). We are proselytizing.

Further, we hope that readers planning other cross-disciplinary endeavors may find in our experiment some guidance for procedures. To this end, [Appendix B](#) of our report details the steps we took and the ingredients we sought in order to maximize the possibility of a favorable result of our experiment. Let me briefly sketch those steps and ingredients here.

Clearly, the most important ingredient is the people who participated in the experiment. Our effort had the benefit of participation by top-flight researchers in several of the cognitive sciences, survey methodology, and applied statistics, as well as dedicated participants from the government agencies involved--the Bureau of the Census and the National Center for Health Statistics (NCHS)--and an enormously competent staff. (Note that these categories are by no means mutually exclusive; see [Appendix D](#) for biographical sketches of participants.)

Second, what is needed is willingness to work hard and time in which to get the work done. Staff devoted much time before the meeting of the seminar to inviting participants, to assembling (and in some cases, preparing) background material, and to detailed and careful planning and coordination of an agenda and of social events. During the seminar, their efforts as rapporteurs kept all participants informed of everyone's progress. Participants were asked to give time beforehand to familiarize themselves with the background materials, to spend almost a week together at the seminar, to spend time after the seminar reviewing draft reports and writing up research plans, and to reassemble seven months later to review their progress. (See [Appendix A](#) for the background papers prepared for the seminar and [Appendix C](#) for a list of other background materials.)

A third ingredient is what might be called the ambiance of the seminar meeting. We were together first in an attractive setting, St. Michaels, Maryland, many miles both from our offices and urban distractions or other intellectual intrusions. Our discussions--formal and informal--often continued late into the night. Feelings of trust, intellectual respect, friendship, and tolerance for idiosyncrasies grew, and with them the excitement of shared research ideas and prospects. We had been warned that people from different disciplines would not listen to each other, but we developed a vocabulary that made cross-disciplinary communication not only possible but inviting; the image of a newly formed common culture is only slightly exaggerated. One example of the acculturation is supplied by a participant who joined the enterprise with a research agenda that was advertised as unalterable by any outcome of the seminar. By the time the seminar reconvened in Baltimore six months later, this participant had actually carried out research to test some ideas generated at St. Michaels.

Thus we feel we can write a prescription to maximize the chances of success in cross-disciplinary collaboration. Through the support and vision of the National Science Foundation (especially from Murray Aborn) and the efforts of members and staff of the Committee on National Statistics (especially chair Stephen Fienberg and Miron Straf), together with an informal group of advisers (who came to be known as friends of the seminar and included, especially, Robert Abelson and Phillip Stone),

we were able to assemble an excellent group of participants. Detailed planning (by all of those already listed and especially by Tom Jabine and the authors of the formal background papers--Norman Bradburn, Catalina Danis, and Roger Tourangeau) provided an appropriate and stimulating agenda for the participants. Preparatory work by the NCHS in supplying background materials and the survey instruments for the National Health Interview Survey and by the Census Bureau in interviewing participants gave us further common experiences to draw upon in our discussions. The participants themselves worked very hard in an atmosphere ideally conducive to the encouragement of cross-disciplinary collaboration. The volunteer respondents who consented to the videotaping of an interview using the NHIS questionnaire must remain anonymous and thus cannot be singled out for thanks. They should know, nevertheless, that their contribution was invaluable in providing the participants with a shared resource that was used repeatedly in our discussions. All of these contributors--and others too numerous to mention but who were enormously helpful--receive heartfelt thanks from me as chair of the seminar. My fond hope is that their efforts will be so successful that they will also earn the gratitude of practitioners of the new cross-disciplinary research tradition they will have helped to start.

Judith M. Tanur  
Montauk, New York  
June 14, 1984

About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.

# CHAPTER 1

## THE SEMINAR

### BACKGROUND

Desirable linkages between disciplines do not always develop spontaneously; deliberate efforts are often needed to encourage them. This report describes one such effort. The two primary disciplines involved were cognitive psychology and survey research, but other cognitive scientists and statistical methodologists also played important roles.

The Advanced Research Seminar on Cognitive Aspects of Survey Methodology (CASM) was convened by the Committee on National Statistics (CNSTAT) with funding from the National Science Foundation. The seminar, held in St. Michaels, Maryland, on June 15-21, 1983, and a follow-up meeting held in Baltimore on January 12-14, 1984, were the main elements of the CASM project, whose goal was to foster a dialogue between cognitive scientists and survey researchers and to develop ideas and plans for collaborative research.

This is the report of the CASM project. The primary audience for this report consists of cognitive scientists, survey researchers, and others with substantive interests in these fields. A second audience consists of persons interested in the broad question of how to foster interdisciplinary communication and collaboration. For the benefit of the latter group, [Appendix B](#) gives procedural details: it explains how the seminar and follow-up meeting were organized and conducted.

The cognitive sciences are concerned with the study of such processes as understanding language, remembering and forgetting, perception, judgment, and inferring causes. Because all of these and other cognitive processes are important in survey research interviews, it would not be surprising to find a fairly long history of collaboration between cognitive scientists and survey researchers in problems of mutual interest. Strangely enough, however, until a few years ago members of the two disciplines appear to have had little contact.

In recent years one begins to find some convergence between the two groups, at least in the sense of discussing problems of joint concern. In 1978 the Social Science Research Council (United Kingdom) and the Royal Statistical Society sponsored a one-day seminar to discuss the problems associated with the collection and interpretation of retrospective and recall data in social surveys, with participation by psychologists and social scientists as well as survey researchers (see Louis Moss and Harvey Goldstein, editors, [The Recall Method in Social](#)

Surveys, NFER Publishing Co., Ltd., 1979). In 1980 the Committee on National Statistics convened a panel on survey measurement of subjective phenomena to explore issues related to the validity and reliability of such measures. The panel included, along with survey researchers and statisticians, one cognitive psychologist and several members from the social sciences. One of the panel's recommendations called for intellectual input from social and cognitive scientists in an extensive interdisciplinary investigation of the subjective aspects of survey questions (Recommendation 12 in the panel's summary report, Surveys of Subjective Phenomena, National Academy Press, 1981; the complete report and papers of the panel, Surveying Subjective Phenomena [2 vols.], are scheduled for publication by Russell Sage Foundation in 1984).

In connection with ongoing work on the redesign of the National Crime Survey (NCS), the Bureau of Social Science Research, with support from the Bureau of the Census and the Bureau of Justice Statistics, convened a two-day workshop in September 1980 that brought together a number of cognitive scientists and survey statisticians (see [Appendix C](#), Item 1). The discussions at the workshop focused on factors affecting the success or failure of respondents in NCS interviews in trying to recall incidents of victimization and remember details of those incidents. A number of suggestions were made that illustrated how knowledge of cognitive processes might be applied in a survey context. There was strong agreement that survey questionnaire design could benefit from the application of ideas from cognitive sciences, and, conversely, that cognitive researchers could benefit by thinking of surveys as experiments for testing their theories.

The CASM project can be regarded as a logical outgrowth of the 1980 workshop. Those who first proposed the project believed that an effective effort to construct an interdisciplinary bridge between cognitive sciences and survey research should have the following four characteristics:

- (1) It should attempt to develop ideas and plans for collaborative research involving cognitive scientists and survey researchers.
- (2) In addition to recall, which was the primary topic of the 1980 workshop, it should consider other cognitive processes that take place in survey interviews, such as comprehension and judgment.
- (3) A small group of experts from the two disciplines, accompanied by a few applied statisticians and representatives of other relevant fields, should meet for an extended period to further their understanding of the areas of intersection between the cognitive sciences and survey research and to stimulate ideas for relevant research.
- (4) Above all, participation in the project should offer potential benefits to members of both disciplines: for survey researchers, through the application of cognitive research to data collection problems; for cognitive scientists, through exploration of the potential uses of surveys as vehicles for cognitive research.

These four criteria provided the framework for the planning and conduct of the seminar.

There were 22 participants in the seminar, including cognitive scientists, survey researchers, applied statisticians, staff of the Bureau of the Census and the National Center for Health Statistics, and staff of CNSTAT (see [Appendix D](#) for biographical sketches of the participants).

Many of the discussions were focused on a specific survey, the National Health Interview Survey (NHIS) of the National Center for Health Statistics. Prior to the seminar, participants had been interviewed in that survey, and two interviews with volunteer respondents were videotaped for viewing at the seminar.

The remainder of this chapter can be thought of as the proceedings of the seminar. It is a synthesis of the discussions at the 1983 St. Michaels seminar and the 1984 follow-up meeting in Baltimore. The views of the participants and their ideas for research are organized by topic and are not attributed to specific individuals. No citations from the literature are included, although the discussions obviously drew extensively on the findings of researchers in relevant disciplines (see the papers in [Appendix A](#) for many pertinent references). The immediate purpose of the seminar was the development of proposals for cross-disciplinary research: the format of the proceedings reflects the informal seminar format that was considered best suited for that purpose.

[Chapter 2](#) is about outcomes. It describes research activities undertaken and research plans developed by CASM participants, individually or in small groups, after the St. Michaels meeting. Other outcomes, such as papers presented at scientific meetings or published, are also described.

Two background papers that examine relationships between the cognitive sciences and survey research were prepared and distributed ahead of time to the St. Michaels seminar participants. A third paper, which looks at the role of record checks in measuring the validity of self-reported data in surveys and experiments, was developed from a presentation on this topic at the seminar. All three of these papers (in somewhat revised form) are included in [Appendix A](#). As noted above, [Appendix B](#) describes the preparations for and the conduct of the seminar, [Appendix C](#) is a selected bibliography prepared for seminar participants, and [Appendix D](#) contains biographical sketches of the participants.

## INTRODUCTION

During the six-day seminar at St. Michaels and the two-day follow-up meeting in Baltimore, a vast number of topics were broached, arguments aired, and suggestions made. It is difficult to do justice to the full range of the enthusiastic and wide-ranging discussions. This account attempts to strike a balance between comprehensiveness and depth. A few common threads ran through much of the discussion, and these are considered in some detail. Other ideas, though no less valuable, cannot be so easily stitched into a general pattern, and these are given more cursory treatment.

This section is organized under three headings: surveys as a vehicle for cognitive research; methods for improving surveys; and issues



particularly relevant to the National Health Interview Survey (NHIS). Several of the recurring themes of the conference fall under each rubric.

Under the first rubric, for example, we note the issue of the barriers to the use of survey methods within the cognitive sciences. These barriers include differences in methods used for research on survey methodology and those used for research in the cognitive sciences, especially cognitive psychology. For the most part, reputable survey research is carried out with probability samples that represent broad population groups; cognitive research, by contrast, is typically conducted with samples of convenience drawn from the student population at a single college or university. Survey research is usually carried out in field settings, such as the respondent's home, and involves everyday tasks, such as the recall of recent events; cognitive research often uses the controlled conditions and artificial tasks of the laboratory.

These obvious differences in research practice reflect more basic differences in what might be termed "research culture." Many surveys are carried out under government auspices and are concerned with pressing policy issues. Timeliness is often crucial. Most cognitive research, on the other hand, is conducted in academic settings and has theoretical rather than applied aims. Although researchers in both groups are concerned about error, they differ in their views of the sources of error and in their assessments of the relative importance of the various sources. Survey researchers typically use their sample data to make estimates for a target population and to calculate the sampling errors of those estimates. While they are well aware of and concerned about the potential effects of nonsampling (measurement) error on their estimates, these effects are infrequently quantified. Cognitive scientists, in contrast, make little use of the formal apparatus of sampling theory and are only secondarily concerned about the generalizability of their findings to broad populations. However, they often include quantitative measures of nonsampling error, such as reliability or validity coefficients, in presenting their results. Even when the two groups appear to share a concept--such as the concept of attitudes--it turns out on closer examination to have different meanings within the two disciplines. And when one reflects on the important subcultures within each research community, one can easily understand the difficulties in bridging the gap between survey researchers and cognitive scientists.

Differences between the practices and concepts of the survey and cognitive research communities have a number of concrete consequences for the study of topics that are of concern to both fields. Within cognitive psychology, for example, memory researchers know the items that the subjects are trying to recall, and much of their work on retrieval cues presupposes this knowledge. Survey researchers, on the other hand, do not usually know the facts that the respondents are trying to recall, and methods for providing the respondent with useful cues to retrieval under these circumstances are not well known. The question of how to help someone remember something without knowing what the something is has not yet received much attention within cognitive psychology, but it is a provocative one that arises from a consideration of the problems faced by survey research.

It would be as wrong to overstate the differences between these two research communities as to overlook them. They share many concerns, including one that encompassed most of the discussion in *St. Michaels and Baltimore: What are the uses and limitations of self-report data?* This central methodological question transcends the boundaries of these two disciplines and is no less relevant to clinical psychologists, historians, and sociologists than to survey and cognitive researchers.

A few other major themes are worth mentioning before we present a more detailed review. The seminar group identified a number of key issues concerning the limits and organization of memory that are relevant in the survey context. One general issue concerns ignorance about the distribution of memory ability across types of events and types of people. There is little detailed knowledge of what kinds of real-world events are likely to be remembered accurately and which are likely to be forgotten. There is similar ignorance about the distribution of memory ability in the general population, which is, of course, the population of interest in most surveys. A related issue concerns the organization of memory. In informal settings, people talk about events using narrative or script-like structures to organize their accounts; there is evidence that memory is organized along similar lines. Clearly, this “natural” organization is quite different from the usual organization of questions in a survey questionnaire. A number of participants raised the question of how recall would be affected if interviews more closely paralleled the structure of events in memory.

The clash between the organization of memory and the organization of an interview leads to another general theme--the need to understand the social psychology of the interview and the role of the interviewer. Interviews can be a frustrating experience for respondents, who may have fairly simple stories to tell but are forced to recast them to fit the terms of the questionnaire. Questions designed to prompt complete and unambiguous answers may be seen by respondents as repetitive and tedious; they may be felt as interruptions that actually impede the smooth flow of information from respondent to interviewer. The model of the interviewer's role that underlies much of survey research is that the interviewer should be a neutral recording device. Although this view has much to recommend it, particularly in opinion research where it is important for the interviewer not to bias the respondent, it may be less appropriate in other types of survey research. For example, some questions place a premium on accurate recall, and interviewers might help respondents by suggesting strategies for retrieval; other questions require estimates or judgments, and there interviewers might help respondents by giving them anchors for the judgments. The model of the interviewer as a kind of collaborator of the respondent also has implications for how an interview should be conducted. One possibility is a two-stage interview process. In the first stage, respondents would be invited to tell their stories in their own terms. In the second stage, an interviewer and a respondent would fill out the questionnaire together. The aims of such a procedure would be to humanize the interview situation and to reduce the frustration it can engender--leading, it is hoped, to fuller recall and more accurate reporting.

Another recurring theme in the work of the seminar concerned the content of the NHIS and might be labelled psychological and cultural factors in the definitions of health and illness. Respondents who see themselves as generally healthy may be more prone than others to underreport specific conditions; at the other extreme, it is plausible that respondents with serious, chronic ailments may underreport minor complaints because their threshold for counting a condition as an illness may be raised. Cultural groups may differ in the terms they use to refer to particular illnesses, or they may label as illnesses certain conditions that do not correspond to diseases in standard medical classifications. The rules for deciding whether a condition warrants medical treatment may depend on nonmedical considerations, including family roles (parents may decide for their children, wives for their husbands) and the mechanism for paying medical bills. These substantive issues were explored not only as interesting topics in their own right, but also as a means for improving the NHIS--the more that is known about the subjective side of health, the better information can be elicited about the objective side.

### **SURVEYS AS A VEHICLE FOR COGNITIVE RESEARCH**

Cognitive researchers have neglected the survey as a research tool, and a substantial portion of the seminar discussion concerned ways to change this situation. From the point of view of cognitive psychology, survey methodology offers a number of advantages over classical laboratory methods. Well-run surveys use probability samples selected from well-defined populations, and these samples are often much larger than those used in most laboratory studies. In addition, surveys require the use of processes (such as retrieval, over very long periods of time, of information concerning naturally occurring events) that are difficult to simulate in the laboratory setting. Consequently, some questions of considerable interest to cognitive researchers can best be studied in the context of large-scale surveys.

#### **Collection of Survey Data on Cognitive Abilities**

Data from a large national probability sample would have an immediate payoff for cognitive psychologists in the study of cognitive, especially memory, abilities. A number of questions about cognition are virtually impossible to answer without benchmark data from large, representative samples. It is not known, for example, whether memory ability generally declines with age and, if so, how sharp the decline is for different types of memory. Researchers at the seminar called for a national inventory of memory and cognitive ability to remedy this gap (see Tulving and Press, Chapter 2). The national inventory would administer a small set of standardized cognitive and memory tests to a probability sample of respondents, perhaps in conjunction with an ongoing health survey, such as the Health and Nutrition Examination Survey or the NHIS. A national inventory of cognitive and memory abilities would provide national norms

for cognitive and memory skills that are needed to address theoretical questions concerning the relationship between age and ability. The norms might have an important practical application as well: it is thought that one of the first symptoms of Alzheimer's Syndrome is the deterioration of memory; age-specific norms on standardized memory tests would greatly facilitate the diagnosis of this disorder.

A related proposal called for a national survey on cognitive failures or slips. All of us are prey to such mnemonic mishaps as losing our car keys or forgetting appointments. The survey would ask respondents to report on a number of these everyday cognitive failures and would be used to develop norms for these failures. In contrast to the proposed national inventory of cognitive and memory ability, the data would be based on self-reports rather than standardized measures. There has already been some experience with questions on these topics in surveys of selected population groups. The data would be useful in addressing questions concerning beliefs about memory (a topic referred to as metamemory). In conjunction with objective measures, the data could, for example, be used to determine whether older respondents believe that their memories are failing and the relationship, if any, between perceived and actual memory loss.

### Other Proposed Surveys

Several other surveys on topics of interest to cognitive scientists were proposed at the seminar. One concerns the relationship between public history and private recollection (see Schuman and Converse, Chapter 2). Such a survey would examine individual interpretations of national events, like the Great Depression. It would compare the perceptions of people who lived through the events with those of people who have merely read or heard about them. It would explore the impact of prior events on the interpretation of recent events: for example, it would seek to determine whether respondents who experienced the Vietnam War interpret its lessons differently from those who have only read or heard about it.

Another proposed survey would explore "naive" economic theories, general beliefs about how the economy works. Aside from its intrinsic interest as a study of the organization of belief systems, such a survey would examine the effect of economic beliefs on individual economic behavior and, in aggregate, on the economy.

The final proposal called for a survey of emotional experience. Such a survey would collect data on the range of emotional experiences in the general population and examine a variety of questions, such as the impact of emotion on mental and physical health and the relationship between emotions and their expression.

### Surveys as Cognitive Experiments

The proposals described so far have in common the collection of survey data on issues relevant to cognitive researchers; this section describes how surveys can provide a context for experimental cognitive research.

Surveys can be viewed as large-scale experiments on cognitive processes--such as memory and judgment--in relatively uncontrolled settings.

Surveys have some serious drawbacks as a setting for memory experiments. One obvious problem is the difficulty in determining in a survey whether recall is accurate, but this problem can, in many cases, be overcome. In the context of a longitudinal study, reports from later waves of the survey can be checked against presumably more accurate reports from earlier waves. Even in cross-sectional surveys, a subsample of the respondents are often reinterviewed in an effort to control the quality of data collection; these "validation" interviews provide a basis for assessing the accuracy of recall. When respondents are interviewed more than once as part of a longitudinal design or validation procedure, it is also possible for interviewers to make observations that can then be used to assess the accuracy of the respondents' reports. Sometimes records are available to help distinguish accurate from inaccurate reporting, but the difficulties in using record checks to determine the accuracy of reporting should not be underestimated; it is easy to exaggerate the degree of "forgetting" when either the records themselves or the procedures for matching records to respondent reports are not perfectly reliable (see Marquis, Appendix A).

Surveys, particularly pilot studies for surveys, often include methodological experiments as integral components, and these experiments have been underutilized by cognitive scientists as vehicles for research. Aside from opportunities for classical experiments, surveys provide a setting for quasi-experimental studies on the impact of situational factors on cognitive processes. Information on factors affecting individual interviews could be recorded by interviewers and used to measure the effect of situational variables, such as the presence of other family members or the length of the interview, on cognitive performances, such as judgment and recall.

### **Surveys as a Paradigm for Research**

One point made repeatedly during the seminar is that the interview is a theoretically interesting situation that is fundamentally different from the situations that usually confront respondents in laboratory settings. Several proposals--ranging from the very concrete to the very general--shared the notion that laboratory methods should be used to illuminate processes found in the interview situation.

One concrete proposal was for laboratory investigation of memory phenomena associated with reference periods. Survey researchers have found that underreporting is more marked for events that occur at the beginning of a reference period than it is for later events. Respondents who are asked to recall their health problems during the last year recall more problems that occurred during most recent six months than during the earlier six months. They also recall more problems for the most recent six months than respondents who are asked only about problems that arose during those six months. The effect is not limited to long reference periods--it is observed even for two-week reference periods, for which

underreporting is greater for events that occurred in the first week--and so appears amenable to research in laboratory settings. Another memory phenomenon observed by survey researchers but neglected by cognitive scientists is telescoping, the reporting of an event that actually occurred outside the bounds set by the reference period.

Another concrete proposal for experimental research involves studying the relationship between the organization of memory and the optimal retrieval strategy. A common design in survey research requires gathering parallel information about all members of a household. Experiments could help to determine which sequence of questions produces the fullest recall--person by person, topic by topic, or some other organization. It might also be possible to give respondents some flexibility, allowing them to choose one order or the other or to switch back and forth. Such studies would address a number of intriguing theoretical questions--How are memories for everyday events organized? What strategies do respondents use to retrieve such memories? Are some strategies more effective than others? Can reordering the questions influence the choice of strategy?--and would have practical implications for survey research as well. By the time of the Baltimore meeting, a pilot study along these lines had been conducted (see Loftus, Chapter 2) with interesting results.

A more general suggestion called for laboratory research to investigate the effects of "interfering" variables on the cognitive processes that are important in the survey setting. Laboratory experiments could, for example, examine the effect of the presence of other people on retrieval or comprehension. A number of such situational factors commonly present in survey interview settings are thought to affect cognitive processes, but there has been little systematic investigation of their impact.

Survey research suggests not only new phenomena for cognitive researchers to investigate, but also new methods of investigation. Since Ebbinghaus initiated the scientific study of memory in the 1880s, memory researchers have relied on a single paradigm: the experimenter has control over the information to be remembered and, in consequence, knows exactly what the subjects are trying to recall. In contrast, survey researchers typically do not know what the respondents know. A new type of memory research was suggested that would parallel the survey situation; in the new paradigm, the experimenter who structures the memory task would know only in a general way what material the subjects have learned. The aim of such research would be to determine whether experimenters can provide retrieval cues or suggest retrieval strategies that enhance recall without having exact knowledge of the material to be remembered.

A related suggestion called for memory research concerning everyday events for which records (or some other means for validation) are available. Memories for financial transactions, for example, can often be checked against entries in a checkbook; on a college campus, memories concerning doctor visits can often be compared with campus clinic records.

Not all of the proposed cognitive research concerned memory. One proposal focused on the study of judgment processes. Surveys often ask



respondents to make estimates, or judgments, regarding how frequently they have engaged in a particular behavior. An NHIS supplement, for example, asks respondents how many alcoholic beverages they drink in a typical two-week period. Psychologists have long been interested in the processes by which such judgments are made. One line of investigation has demonstrated that judgments are affected by the use of “anchors,” which serve as starting points for the judgment. In answering a question about their typical drinking behavior, respondents may try to recall how much they have drunk recently. This recollection is the basis for a preliminary estimate, or anchor, which is then adjusted upwards or downwards to reflect behavior during the typical period. When the anchor is misleading or the adjustment insufficient, the final estimate will be thrown off. Can the estimates be improved by providing accurate anchoring information? Survey researchers have balked at the idea of providing information, such as anchors, that might bias the respondents (a carryover, perhaps, from opinion research settings); the research on judgment indicates that respondents may bias themselves by generating misleading anchors. A line of research was proposed in which the information provided to respondents would be varied systematically: some respondents would be given detailed information about the distribution of responses, some would be given information about the mean response, and some would be given no information at all. The aim of the research would be to determine whether the information provided increases the accuracy of judgments.

### Summary

The survey method is a research tool that cognitive psychologists have neglected. Survey data can be collected on topics of interest to cognitive psychology, such as the distribution of cognitive abilities in the general population and the intersection of public and personal history. Surveys can provide a vehicle for experimental and quasi-experimental studies on cognitive processes in relatively uncontrolled settings. Finally, survey findings suggest new phenomena and new research paradigms for cognitive researchers to explore in laboratory settings.

### IMPROVING SURVEY METHODS

If survey research has much to offer the cognitive sciences, then the proposals made at the seminar indicate that the cognitive sciences also can contribute to survey methodology. The proposals are grouped into four categories: general strategies for improving survey methods; cognitive research that has special relevance to survey methodology; issues calling for further methodological research; and research tools that might be applied profitably to questions concerning survey methods.

### General Strategies

One of the suggested general strategies for improving surveys was to include methodological research as a component of every large survey. Other researchers bemoaned the “morselization” of methodological research. It is not enough to catalogue an ever larger number of response effects in surveys; instead, research on response effects must be more systematic and quantitative--survey researchers need to know not only what the potential problems are, but also when they are likely to arise and how seriously they will bias the results.

In addition, the impact of response effects or other measurement errors must be incorporated into assessments of the reliability of survey results. Statistical models have been developed to measure the likely effects of sampling error; similar models are needed to assess the impact of measurement error. One approach considered by the seminar group was to treat survey items as a selection from a population of potential items; standard errors for survey estimates would then reflect both the sampling of respondents and the sampling of items.

Many of the participants were struck by how much a standardized interview differs from the acquisition of information in a normal conversation. It was proposed that survey instruments be organized to follow the same principles that work well in everyday conversation. A prerequisite, then, for the design of a survey instrument would be the study of ordinary conversations about the survey topic. This proposal was related to two broader concerns. The first is the apparent frustration of respondents at the artificiality of the typical survey interview. Interviews that were structured more like conversations would be “humanized”--less mechanical for both respondents and interviewers. The other broad concern is poor recall. A general hypothesis that emerged during the conference was that survey questionnaires might induce more accurate recall if their organization paralleled the organization of the experience in memory. The flow of ordinary conversation would provide a good indication of how memories for a class of events are organized.

### Relevant Research from the Cognitive Sciences

As is already apparent, a number of areas of investigation from the cognitive sciences were seen as particularly relevant to survey methodology. One such area is research on scripts and schemata. In the cognitive sciences, most researchers share the view that the interpretation and memory of experience is governed by higher-level knowledge structures. These higher-level structures, referred to variously as scripts, frames, or schemata, codify shared knowledge about classes of things or events. For example, there may be a script for visits to the doctor, which represents general assumptions about why people go to the doctor and the sequence of events that a doctor's visit usually comprises. These scripts may vary from person to person depending on such variables as type of health care system. For example,



the sequence of typical events for families using health maintenance organizations may differ from the sequence for families using another type of health care system.

Scripts and related concepts were seen as relevant to several points under discussion. One issue concerned respondents' reactions to the survey interview situation: Does the interview situation evoke a script and, if so, what are its properties? There are a number of situations in which people are required to answer questions, ranging from applying for a job to voting, and each one suggests different ways of responding. It is one thing to see an interview as just another situation in which the demands of an inflexible bureaucracy must be satisfied, quite another to see it as a situation in which people have their say. A number of proposals were made to explore basic interview scripts and the individual and cultural variations on them.

Scripts were seen as relevant to another methodological issue, the issue of the accuracy of memory. If questions in an interview were asked in an order that paralleled the sequence of events in the relevant script, then recall for the events might be improved. A common hypothesis was that recall would probably be facilitated for events that fit the stereotypical pattern of the script but might be hindered for events that departed from the script in some way. This would suggest an experiment in which two factors--the routineness of the events to be recalled and the structure of questionnaires (script versus usual order)--were simultaneously manipulated. Aside from any effect on memory, questionnaires that followed the ordering of the relevant script might speed up interviews and reduce the level of respondent burden.

Judgmental heuristics (the strategies used in rendering a judgment) is a second topic of cognitive research that has broad applicability to questions of survey methods. There is a large body of results indicating that people use a number of strategies, each with its characteristic shortcomings, to estimate frequencies and probabilities. We noted above that respondents may use a strategy of anchoring and adjustment to answer questions concerning their typical drinking habits, and their answers might be improved if they were provided with accurate anchoring information. The general point is that it is important to investigate what strategies, or heuristics, respondents use when they answer different types of questions and what errors result from those heuristics. One method that has proven useful in the study of problem solving is to have subjects think out loud during their attempts to reach a solution. Respondents in a pretest could be similarly encouraged to think out loud as they answer questions. An analysis of the protocols (transcripts) of these pretest interviews might reveal the strategies respondents commonly use to answer different types of questions. With a better understanding of the process by which survey questions are answered, it should be possible to improve response quality. Questions might be rephrased to encourage the use of strategies that yield more accurate estimates or such strategies might be directly suggested to respondents.

Several other topics in cognitive psychology were discussed. Since underreporting is a very serious problem for many surveys, it was natural that techniques to produce hypermnesia (increases over time in the amount

of recall) were considered. For example, repeated recall attempts can produce increases in the amount of material recalled. In addition, there was considerable interest in the part-set cuing effect: part-set cuing refers to the use of some members of a class as prompts for the retrieval of the other members. Giving too many class members as cues may hinder performance--subjects given a list of 25 of the states may recall fewer of the remaining 25 than subjects simply asked to name all 50. Questionnaires often provide examples in order to facilitate understanding and recall; research on the part-set cuing effect suggests that, under some circumstances, this strategy may reduce rather than increase recall. Research under survey conditions might give guidance about the optimum number of examples to provide.

### **Methodological Issues and Research Proposals**

Cognitive research can only go so far in answering questions about survey methods. In order to settle methodological issues, it is necessary to conduct new methodological research. A large number of topics needing further investigation were identified; in some cases, specific studies were proposed. Many of the topics concerned response error; those topics can be conveniently grouped into three categories--comprehension, recall, and judgment--that correspond to the tasks that respondents perform in answering questions. The remaining topics concerned the behavior of the interviewer or features of the interview situation itself.

#### **Comprehension**

Survey researchers are faced with the difficult problem of individual and cultural variations in the interpretation given to particular terms in survey questions. Not only do abstract terms, like "big government," evoke a wide range of meanings, but even relatively concrete terms, like "doctor," can have varied interpretations. The aim of survey researchers is to reduce these ambiguities so that all respondents are answering the same questions. Research is needed to see whether questions (or answers) need to be translated to reflect cultural differences in usage; in other cases, concepts may differ so radically among respondent groups that new question orientations may be required rather than mere rephrasing.

Survey instruments often provide examples to clarify the intent of a question, a practice that raised a number of issues. When examples are given, it is unclear whether they should consist of typical instances, which may help to define the category of interest, or atypical ones, which may help to sharpen the category boundaries. In some cases, it is possible to give an exhaustive list; in others, it may be useful to provide some "non-instances" to indicate what is excluded by the question. There is little research to indicate how people decide whether to include ambiguous, boundary cases in a category: Is a chiropractor a doctor? How does one decide? The number and kinds of examples that result in the clearest understanding of the question may depend in part on the category of interest. With fuzzy categories, it may be best to

give some atypical examples (“include chiropractors”) as well as some “non-examples” (“exclude physical therapists”). With familiar categories that have sharply defined boundaries, examples may be unnecessary.

Sometimes visual aids can be used to reduce confusion. A supplement to the NHIS asks respondents to estimate their drinking behavior in terms of ounces. It would help to show them glasses of different sizes with the capacities labelled, but even this procedure would still leave room for other ambiguities.

The abstract terms used in attitude questions raise even more difficulty for comprehension. Terms like “social programs” probably mean as many different things as there are different respondents. Even for a single respondent, the same term may evoke different meanings on different occasions. One reason that question order may affect results is that earlier questions can provide an interpretive context for later ones. “Social programs” may be interpreted one way in a series of questions about waste in government and another way in a series of questions about the problems of disadvantaged people. Further, in attitude questions it may not be possible or even desirable to separate the meaning of a term from its evaluation--part of what it means to have an attitude is to have a propensity to view the object of the attitude in a particular light. Clearly, further research is needed to determine how respondents interpret and answer attitude questions. (For some suggested research on these issues, see Tourangeau et al., Chapter 2.)

## Recall

One of the central problems of surveys is that survey results are often no more accurate than the memories of the respondents. The question of how to improve recall was perhaps the central question of the seminar. We have already noted that the sequence of questions in a questionnaire may affect the accuracy and completeness of recall. Several specific studies were proposed to compare different question orders. With life-event histories, a topical order could be compared with a chronological order; with household interviews, such as the NHIS, a person-by-person order could be compared with a topical order. Items regarding individual events might be ordered to reflect the script for that class of events. Even when questions follow a chronological organization, it may make a difference if recall proceeds from the most recent events to the least recent rather than in the opposite order (see Loftus, Chapter 2).

Question order relates to another concept from cognitive psychology, the concept of proactive inhibition. When subjects in memory experiments are asked to recall lists of related items, performance gets worse on the later lists, a phenomenon referred to as the build-up of proactive inhibition. When the items on later lists bear little similarity to those on the earlier lists, the effect disappears; the effect can also be reduced by increasing the time interval between trials. These findings suggest that research might be carried out to see whether periodic changes of topic or rest periods would promote fuller recall in interviews.

Many questions for further research concerned the reference period. Research indicates that events may be dated more accurately if they can be tied to some landmark event. Would it be helpful, therefore, to give respondents a warm-up period during which they would think about where they were and what they were doing during the reference period? Even if they did not recall personal landmarks during the warm-up, respondents might be encouraged to think about general spatial and temporal cues that could facilitate recall. Researchers also suggested several variations on the current definition of the reference period. A three-week period could be used (instead of the current two); all episodes would then be dated and only those in the two more recent weeks would be retained. The use of such an extended reference period might reduce underreporting, which is thought to be greater at the beginning of a reference period. In another variation, a rolling reference period would be tried; rather than reporting about a period defined by fixed dates, respondents would report about the two weeks prior to the interview.

Another tack for possibly improving respondent recall involved forewarning respondents about the content of the interview. For example, with computer-assisted telephone interviewing (CATI), it would be possible to contact respondents at the beginning of the reference period. This initial contact would provide respondents with a landmark for dating events; it also would provide an opportunity to suggest strategies for improving recall (e.g., noting doctor visits on a calendar, thinking about health problems every night). Even if CATI were not being used, an advance letter could include the forewarning and suggestions for memory aids.

Another area for research proposed above, regarding the use of examples, also has implications for respondent recall; examples and lists not only illustrate the meaning of a question, but also serve to prod recall. It is unclear whether atypical members of a category are especially hard to recall or are especially memorable; the overall efficacy of different types of examples may depend on their effect on memory. The work on part-set cuing indicates that, in some cases, less is more--too many examples can inhibit recall.

## Judgment

Many survey questions require some judgment or estimate from respondents. Human judgments are, of course, fallible, and it is natural to ask how they are made and how they can be improved. One suggestion was made repeatedly: experiment with giving respondents a chance to revise their initial answers, or asking them for a second estimate. Anecdotal evidence suggests that an adjusted answer or second estimate is often more accurate. A related line of work concerns the use of qualitative, controlled feedback, in which respondents are informed about the reasons cited by other respondents for making a quantitative judgment. Respondents become more confident when they hear other respondents' reasons, but not their numerical assessments. Second-chance methods might also be used with questions that rely more on memory than judgment;

respondents could be asked at the end of an interview whether they had recalled additional events they had not reported earlier.

Some types of judgments present special problems. Attitude questions usually require an evaluative judgment, and little is known about how judgments are made. Is memory first searched for information about the attitude object? What is the “attitude space” that is searched? What is retrieved? How is the information combined to produce the final judgment? The answers to these questions are simply not known. One familiar type of opinion item asks for respondents to list issues according to their importance (e.g., “What are the most important problems facing our nation today?”). Respondents commonly omit problems that are important but not particularly salient (e.g., nuclear war). Once again, little is known about what determines the relative salience of different issues. Another type of judgment that presents special difficulty is the estimation of probability. Research has shown that probability estimates are often at odds with the dictates of probability theory and that the probabilities of rare events are often greatly overestimated. In addition, probability estimates are known to be sensitive to both the framing of questions and the type of response scale that is presented to respondents.

Aside from research to improve understanding of how respondents make particular judgments, two general approaches to reducing respondent error were suggested for investigation. The first approach involves suggesting strategies to respondents for making the estimates. Some estimation strategies are better than others and strategies that are known to reduce error could be suggested to respondents. One proposal along these lines has already been mentioned--give respondents the mean as an anchor for their individual estimates. The second approach involves collecting and using what might be termed ancillary information about the estimate. Respondents might give their answers and then rate their confidence in them. The confidence ratings might then be incorporated into the survey results. The nature of the statistical procedure for incorporating ancillary data remains to be worked out. A further prerequisite to the use of such adjustments would be a study that assesses the correlation between confidence ratings and the accuracy of reports, perhaps using record checks to evaluate accuracy. Other ancillary measures relating to respondent error could be collected. Some of the suggestions included incorporating “lie” scales (to measure the propensity of respondents to give clearly invalid answers: one such scale is used in the Minnesota Multiphasic Personality Inventory), “denial” scales (to measure the propensity to deny or minimize symptoms), and questions assessing item sensitivity (“Would you be embarrassed to report . . . ?”). As with confidence ratings, it would be necessary to assess the validity of the ancillary data and to develop statistical procedures for incorporating them.

We simply do not know much about how respondents answer survey questions, and this ignorance was an undertone in much of the discussion about judgment. Questions that are intended to trigger memories for specific events may, in fact, elicit estimates based on more general knowledge. A study carried out by one of the participants between the

two meetings of the seminar group (see Ross, Chapter 2) indicates that respondents have little more confidence in their answers to questions about past behaviors than in their answers to questions about the future. This finding suggests that both sets of answers may be produced through a similar process that relies more on judgment than memory.

### **Interviewer Behavior**

Although much of the discussion focussed on the respondent as a source of nonsampling error, the interviewer and the interview situation were also seen as potential areas for improvement in survey methods. With respect to interviewer performance, a number of points merited further study. There is much to be learned from good interviewers; several proposals were aimed at learning more about their characteristics and techniques.

It is possible to observe interviewers at work by videotaping interviews or by making arrangements for interviews with “planted” respondents. (Both of these techniques were used in connection with the seminar.) Interviewer effectiveness can be rated on a number of dimensions, including objective performance measures--such as completion rates, item nonresponse rates, and editing error rates--as well as more subjective ones--such as warmth, voice quality, appropriateness of probes, and methods for coping with respondent fatigue. Perhaps the methods used by good interviewers could be taught to all interviewers. On the other hand, if good interviewers are born and not made, videotapes could be used to identify characteristics that might serve as criteria in selecting and hiring new interviewers.

The drive to standardize interviewer behavior has left interviewers little room for discretion. One proposal called for a comparison of different levels of interviewer discretion in dealing with respondent uncertainty. Some interviewers would be instructed not to give any clarification to respondents, some would be given standardized instructions for giving clarification, and some would be given the freedom to decide how much clarification to give. Interviewer discretion might reduce bias but increase interviewer variance. Both effects would have to be measured and the tradeoff weighed. Interviewers could also be given some discretion in determining question order (e.g., topical versus chronological); it would be interesting to see how interviewers would order the questions if they were free to choose. CATI systems may offer a good method for providing interviewers with some discretion; CATI questionnaires could have alternative branching structures for respondents who show a preference for one order over the other, and it would be up to the interviewer to decide which branch to follow.

Several of these proposals imply a conception of the interviewer that differs sharply from the prevailing view. Rather than seeing interviewers as a kind of neutral recording device, they might be viewed as collaborators with the respondents, helping out in various ways. One could conduct experiments in which interviewers would suggest strategies for retrieval and estimation, provide anchoring information for judgments, and exercise discretion in administering interviews. Accuracy rather than absolute standardization would be the aim of such approaches.



## The Interview Situation

A variety of research issues were identified that deal with aspects of the interview situation. They range from the effect of different interviewing modes to respondent attitudes toward interviews. The issue of the mode of data collection (telephone, face-to-face, or self-administered) is particularly urgent for the NHIS, which is committed to a mixed approach, with some interviews being conducted over the telephone and the rest in person. One question is how to gain the cooperation of respondents in the critical first few seconds of the telephone call. Another concerns the use of incentives: Would respondents feel committed to participating in the interview if they were sent a payment or reward in the advance letter?

These questions about gaining respondent cooperation and the use of incentives relate to broader concerns about respondent motivation. Not everyone views an interview in the same light or approaches it with the same motives. Different views may be systematically related to demographic or cultural variables. Poor respondents, for example, may see the survey interview in the same terms as an intake interview for welfare--a view that richer respondents are unlikely to share. These differences among subgroups in attitudes toward the interview could be assessed in a study in which people rated the similarity of the interview situation to other situations. Another approach to subgroup differences in respondent motivation and behavior assumes that, because interviews are a familiar part of contemporary society, people have probably developed rules for appropriate interview behavior. Possible rules might include not bringing up topics unless they are first mentioned by the interviewer and not asking for clarification. Different subgroups may follow different rules. Implicit in this discussion of subgroup differences is the notion that the ways in which respondents view the interviews will affect the level of their cooperation, the amount of deliberate withholding, and the accuracy of their answers.

The discussion about respondent motivation reflected concerns not only about response error but also about nonresponse error. One general strategy to reduce nonresponse is to explore why people answer survey questions rather than why they refuse. A number of motives were suggested--the desire to be helpful, a sense of duty, the wish to present oneself in a favorable light; it is probable that different respondents cooperate for different reasons. From the point of view of survey researchers, it is not the case that all motives are equally desirable: a respondent who wants to be maximally informative is preferable to one who wants to make the best impression.

Because the trend in surveys is toward longer interviews, there is considerable interest in finding ways to maintain the level of the respondent motivation over the course of an interview. It would be helpful to know how respondents' moods and attitudes change during interviews and how these changes affect data quality. It is common in laboratory experiments on memory to restrict testing to 45 minutes or an hour on the grounds that performance over longer periods may deteriorate from fatigue. Experiments to quantify the fatigue effect would provide useful data for both cognitive scientists and survey researchers.

One hypothesis about how motivation changes during an interview is that motivation starts out high but then wanes. At the beginning of an interview, most respondents probably start out with a broad criterion for reporting events; even if they are not sure the events are appropriate, they report them anyway. As an interview wears on and respondents learn the consequences of reporting, their criterion probably narrows. Supplements and other material toward the end of a questionnaire may, therefore, be particularly prone to underreporting. Interviewers may be susceptible to a form of the same problem and neglect to note down trivial conditions. One possible remedy for such criterion shifts is to identify all the relevant conditions or events at the start of the interview, before details are collected on any one condition or event. Respondent motivation and performance may be affected by the pace as well as the length of the interview. Various methods of changing pace could be compared (such as longer questions, rest periods, or multiple contacts), particularly for their effects on recall.

Several participants called for research on procedures that might increase respondent motivation by humanizing the interview situation. One humanizing method might be to reduce the standardization of the interview. Earlier we noted several proposals that would examine the effects of allowing interviewers greater latitude. There were, in addition, suggestions to try tailoring questionnaires to different subcultures or to individual respondents with different scripts for dealing with a topic. In the NHIS, one version of the questionnaire might be suitable for respondents with minor medical problems, another for respondents with serious chronic conditions. Organizing question orders according to conversational principles would reduce the inflexibility that can result from standardization. The most radical proposal along these lines was to try allowing respondents to tell their stories before any detailed questions are asked. The interview would begin as a conversation in which respondents were asked a few general questions (e.g., about their family's health and recent medical problems); respondent and interviewer would then work together to fill out the questionnaire.

A household survey like the NHIS affords some flexibility in the choice of respondent; several researchers offered hypotheses about who that respondent should be. For some purposes, the best respondent for the household may be the person who pays the bills; for others, it may be the "gatekeeper" (e.g., the person who makes the appointments with the doctor.) Some events may be better recalled by children--the first few experiences in a category are often easiest to recall.

It is not always necessary to select a single respondent. On the assumption that several heads are better than one, it may be useful to have several household members present during the interview--what one member forgets, another may recall. On the other hand, household members may distract each other, reducing recall, and respondents may be less willing to answer sensitive questions when other members of the household are present. It is possible to conduct independent interviews with several members of a household to assess the reliability of their reports. Clearly, more research is needed to determine how to take advantage of the fact that the NHIS is a household interview, in which



several persons--or combinations of persons--may serve as respondents. A start could be made by collecting ancillary data on who gave answers to which questions during an interview.

### Tools for Methodological Research

A review of the methods proposed for carrying out methodological studies will help to summarize the discussion of methodological issues. Many of the proposed studies were conceived as experiments that would compare different question orders, levels of interviewer discretion, or respondent rules. Such split-ballot studies, in which portions of the survey sample are randomly assigned to different treatments, have a long history in survey research. Other studies were seen as quasi-experiments. In these studies, natural variations in interview length or setting (e.g., the presence of other household members) would be measured and related to differences in the quality of the data.

A number of the proposals concerned the processes by which survey questions are understood and answered. Random probes inserted immediately after a question can be used to study how respondents interpret the terms in the question. Protocol analysis can be used to investigate the processes respondents use when they answer survey questions. Respondents would think out loud as they answered questions, transcripts would be made, and those transcripts analyzed for clues as to process. Protocol analysis was seen as particularly useful in identifying strategies for answering questions requiring an estimate or judgment. A related technique is debriefing the respondent after the interview has been completed. Such postinterview debriefings can be a very useful method for understanding how respondents interpreted survey questions and for clarifying the meaning of their responses.

Other proposals focused on the interview process; videotapes were seen as an invaluable tool for research on this process. Videotapes could be used to study the relationship between interviewer characteristics and techniques, on one hand, and measures of interview quality (e.g., item nonresponse), on the other. The effects of interactions between household members during the interview could also be explored. Participants at the seminar had themselves viewed a videotaped NHIS interview, and this experience may provide a model for future research endeavors. Videotaped interviews are clearly provocative tools that can stimulate active collaboration between cognitive scientists and survey researchers.

Most of the proposed research concerned nonsampling errors in surveys; several techniques were suggested for assessing the magnitude of nonsampling errors. Respondent reports can sometimes be checked against administrative records, although a number of pitfalls in record-check studies (e.g., errors in the records) can bias the results. Sometimes it is possible for interviewers to make direct observations that can provide a basis for assessing response errors, and sometimes reinterviews can be used to explore the reliability of the interview process. A final method that was proposed involved including measures of validity (such as lie scales) or measures of confidence into survey instruments. These

measures could be used to adjust survey estimates or they could be incorporated into estimates of total survey error.

One major source of error in estimates is underreporting: one proposal called for the development of mathematical models to estimate the amount of underreporting; the model might embody assumptions about the incidence of events of different types and the forgetting curves for each type.

A final set of proposals suggested combining several methods of research on surveys. Researchers might begin with laboratory research on judgment, for example, and then conduct split-ballot experiments to compare several methods for improving the judgments of respondents in surveys. Ethnographic studies could be used to explore variations in terminology or to determine what groups of people are excluded by survey samples. (The Census Bureau has employed similar ethnographic studies to assess undercoverage in the decennial census.) Finally, a cross-disciplinary team could study a few families intensively. These families would be interviewed and videotaped over long periods of time; family records would be checked and direct observations made. The cross-disciplinary method would establish an upper limit on the quality of information available and could be used as the standard for assessing the shortcomings of questionnaire data.

### ISSUES FOR THE NATIONAL HEALTH INTERVIEW SURVEY

A good part of the discussion centered on issues specific to the NHIS, especially issues of content. A general question concerned how the data, especially data on conditions, are used. Several researchers noted omissions from the NHIS and called for items on emotional stress and mental illness. If the NHIS is viewed in part as a survey of attitudes, then the most serious omission is the area of conceptions of health and illness. NHIS supplements could provide answers to many questions about the subjective side of health: What conditions do people include under the headings of health, illness, and injury? What health-related conditions are regarded as nonevents? How do the schemata or scripts for one kind of health event (e.g., an injury) differ from those for other kinds (e.g., an acute illness or chronic condition)? How do different subcultures differ in the conceptions of health and illness and how do their taxonomies for illness differ? How do emotional states affect physical health? We group the issues specific to the NHIS according to content areas of the questionnaire: utilization of medical services, health conditions, and restrictions in activity. Then we turn to a single item in the questionnaire that asks respondents for an overall rating of their health.

#### Utilization

Researchers identified two major issues regarding the items on use of medical services--the process by which people decide to seek help and underreporting of utilization. A number of factors determine when

someone decides to seek help: among those suggested are the nature of the condition or problem, the person's view of the medical system and how he or she relates to it, and the mechanism for paying for medical care. People may "schedule" their illnesses when psychologically convenient and they may seek help when it is convenient or just before it becomes especially inconvenient. The relationship between a person and the medical system may be mediated by a household gatekeeper, the person who usually calls the doctor and makes the appointments for the family.

These issues are interesting in their own right, and they have implications for questions of survey methodology, especially underreporting, as well. For example, questions that focus on the decision-making process might reduce underreporting of medical utilization. They would also aid in identifying the people in a household who are most knowledgeable about utilization--the decisionmakers, the gatekeepers, the people who pay the bills or fill out the insurance forms. The reference period is also relevant to the issue of underreporting. For hospitalizations, NHIS has used a 13-month reference period in some years, a 12-month reference period in others, and in a recent pretest, a 6-month period. It would be worthwhile to examine the estimated distributions of discharges by month under the different reference periods and to compare them to estimates based on hospital records.

### Conditions

The NHIS asks respondents a series of questions concerning medical conditions. One problem with these items is the terminology itself. Respondents may know they have a problem (for example, a bad back) without knowing the appropriate medical term for it. Self-report data might be more accurate if the items asked for symptoms rather than conditions. Or a general item might be added to ask for symptoms that bother the respondents but for which they do not know the cause. Another strategy that might reduce underreporting is to allow respondents to describe each problem in their own terms before proceeding to more structured items.

There are subcultural variations in health terminology and, in some cases, it may be possible for local physicians to translate folk terms (e.g., "high blood") into standard terminology. Even for the same individual, there may be several scripts or schemata for different types of health events (chronic conditions versus injuries), and different question orders or wordings may be needed to prompt the fullest recall of different types of conditions. If standard condition lists continue to be used, it might be easier to put them on individual cards and to group them according to conditions that tend to occur together. Respondents may find it easier to sort cards than to listen to lengthy lists, enabling them to deal with more items; grouping of conditions may facilitate retrieval.

One thread running through most of these suggestions is a concern about underreporting, which can occur because a condition has not been diagnosed, because the respondent does not recognize the term for it,

because the condition has been forgotten, or because the respondent is unwilling to report it. One way to estimate the amount of underreporting is to compare prevalence rates based on NHIS data with those of other surveys--such as the Health and Nutrition Examination Survey (HANES), which includes medical examinations, and the National Medical Care Utilization and Expenditure Survey (NMCUES), which incorporates checks of physician records--or with expert rankings of prevalence. It would, of course, facilitate the comparisons if a common set of conditions were used. It was proposed that the NHIS condition items be included in the HANES interview so that the relation between self-reports of conditions and medical diagnoses could be explored. Less serious chronic conditions (e.g., sinus trouble) may be especially prone to underreporting. An experiment was suggested to compare reporting under the current methods with reporting under methods designed to enhance recall (for example, by leaving respondents a chronic conditions checklist that could be mailed in).

Another thread running through this discussion was the neglect of psychological factors in health. The NHIS includes few items that assess mental health, and it does not include any of the standard scales that measure depression, stressful life events, or physical symptoms associated with stress (e.g., somaticization scales). Because of this omission, the NHIS data cannot be used to monitor trends in the prevalence of mental health problems or to assess the relationship between physical conditions and psychological states. Participants suggested the inclusion of more mental health items in the NHIS, subject to constraints of response burden and cost.

### **Restricted Activity**

As with the utilization and condition items, there was considerable interest in the subjective side of the items concerned with restrictions in activity brought on by illness or injury and considerable concern about underreporting. One proposal was to use random probes to find out how people interpret the term "restricted activity." The present approach may fail to measure the effects of mental illness; it would be useful, therefore, to know whether respondents include mental illness when they think about "illness or injury." Several new approaches to the restricted activity questions were suggested, partly with a view toward reducing underreporting. For the questions on the loss of days from work, respondents could be asked first to report all days lost from work for any reason and then to say why each day was lost. Another approach would be to begin the restricted activity section with questions about normal activities during the reference period. For each activity that they normally engage in, respondents would be asked whether it was curtailed or extended during the reference period and the reason for the change. Some activities, such as reading or watching television, may increase during periods of illness. It might also be useful to broaden the scope of the restricted activity questions by asking respondents whether they had carried out their major activities during the reference

period with less than their customary efficiency and, if so, why the change occurred.

### **Self-Perception of Health**

One item on the NHIS asks respondents to rate their overall health; no other single question provoked as much discussion at the seminar. How do respondents make this judgment? Part of the answer probably involves a comparison process: respondents may compare their current health with their health at other times, or they may compare themselves with other people of the same age. Judgments of overall health are no doubt influenced by objective conditions, but the influence may be limited (e.g., respondents who have successfully adjusted to long-term conditions may discount them in evaluating their health) and perceptions of objective conditions may be as much influenced by the overall judgment as the reverse. Research on underreporting of conditions demonstrates the impact of the overall evaluation on the reporting of conditions: underreporting is greater for respondents who see themselves as healthy. Global judgments typically integrate information from several dimensions. Little is known about the subjective dimensions of health; some multidimensional scaling studies might shed considerable light on the issue.

It would not be surprising to find that the self-perceived health status item is affected by question context. The correlation between the condition items and ratings of overall health might be increased if the condition items came first in the interview. Even if there were a correlation under both question orders, a positivity bias might be expected, with respondents seeing themselves as healthier than their answers to the condition items would warrant.

### **TRANSLATING IDEAS INTO ACTION**

The free-flowing discussions at the St. Michaels and Baltimore meetings led to many ideas about ways for cognitive scientists and survey researchers to work together to their mutual benefit. Surveys can be used to collect data of interest to cognitive scientists and can serve as a vehicle for cognitive research using larger and more heterogeneous samples than those normally used in laboratory experiments. The National Health Interview Survey and other surveys might be improved by applying what cognitive scientists have already learned about comprehension, memory, and judgment. New research studies of the cognitive processes involved in answering survey questions and conducting survey interviews should provide a basis for further improvements.

The real challenge and the main goal of the CASM project was to translate some of these ideas into specific collaborative research programs and activities. Now, slightly more than one year after the St. Michaels meeting, it is evident that this is being done by CASM participants and others. The details are given in [Chapter 2](#).

## CHAPTER 2

### AFTER THE SEMINAR

Chapter 1 summarized the views expressed by the CASM project participants at the St. Michaels and Baltimore meetings, with emphasis on their suggestions for cross-disciplinary research by cognitive scientists and survey researchers. This chapter describes some of the outcomes of the project, including research plans and activities developed by participants after the St. Michaels meeting and dissemination of project results through the publication or presentation of papers and other means.

Each of the first six sections of this chapter describes a research program or activity initiated by CASM participants working as individuals or in small groups. The first four sections describe plans for rather substantial research efforts. The first section describes a multiyear collaborative research program involving cognitive scientists and survey researchers. The plan for this program, which is already under way, was developed by CASM participants Sirken and Fuchsberg for the National Center for Health Statistics (NCHS). The program described in the second section was developed by CASM participants Tourangeau, Salter, D'Andrade, and Bradburn, along with other cognitive scientists. The program, whose objective is to study the cognitive underpinnings of the survey interview process, is also funded and under way.

The third and fourth sections contain prospectuses for survey collections of data that would be of considerable interest to cognitive scientists. Converse and Schuman propose to investigate personal interpretation of recent historical events for a sample of the U.S. population; funding for this project is expected soon. Tulving and Press present a proposal for a national memory inventory in which memory capabilities and other cognitive abilities would be tested for a large probability sample of the U.S. population; although the authors are not now in a position to pursue their proposal, they welcome and would cooperate with efforts by others to undertake the proposed research.

The fifth and sixth sections describe research done by students under the direction of Loftus and Ross, two of the cognitive scientists who participated in the CASM project.

The last section of this chapter describes outreach activities: steps taken by the CASM participants to share the ideas developed during and after the seminar with others and to recruit new members of the interdisciplinary network that has been established.

## LABORATORY-BASED RESEARCH ON THE COGNITIVE ASPECTS OF SURVEY METHODOLOGY

National Center for Health Statistics

(Monroe Sirken and Robert Fuchsberg)

The research project outlined in this plan uses the National Health Interview Survey (NHIS) as a test bed for research and experimentation of the sort discussed at the CASM seminar.

### Purpose and Objectives

Questionnaire design and data collection procedures are among the weakest links in the survey measurement process, and past efforts to improve the quality of survey instruments and procedures have posed serious and difficult methodological problems that are unlikely to be resolved by traditional survey research methods. Therefore, it is essential to test nontraditional modes for conducting research on survey methods. The objective of this project is to investigate the cognitive laboratory as the setting for conducting research on the cognitive aspects of survey methodology. It will tackle three of the most important questions to emerge from CASM. Namely, under what conditions are laboratory methods likely to:

- (1) produce results similar to or different from traditional field methods?
- (2) succeed where traditional methods have failed?
- (3) enhance the results obtained by traditional methods?

Although survey researchers and cognitive scientists are both concerned with the manner in which individuals handle information, their approaches to the problem and the methods used to study the problem are quite different, and there has been very little communication between them. Survey researchers are concerned about the survey measurement process and use field experiments to test response effects in terms of the wording, response categories, and orderings of questions. They make very little, if any, use of controlled laboratory experiments to investigate the manner in which the respondents and interviewers process the information presented by the survey instrument. The traditional method of developing, testing, and evaluating survey instruments involves sizeable field pretests and pilot studies of questionnaires that are developed by survey statisticians and tested under "normal" survey conditions by trained interviewers. Cognitive scientists, on the other hand, are concerned about the system individuals use in processing information. Cognitive psychologists conduct controlled experiments in a laboratory setting involving direct and intensive interaction with relatively small samples of subjects to investigate the mental procedures by which information is processed. A major objective of this project is



to contribute to the advancement of both disciplines, and to effect communication between them, including collaboration in research studies.

The demonstration will conduct laboratory-based research on the cognitive aspects of survey design using the combined methods of the cognitive and statistical sciences. Cognitive knowledge and techniques will be used to gain a better understanding of the effects of cognitive factors in the survey measurement process. From these laboratory findings statistical models will be developed for controlling survey measurement errors.

### **Benefits**

The laboratory is the ideal setting for conducting interdisciplinary research in which the combined technologies of the cognitive, social, biological, and computer sciences can be applied in researching the cognitive aspects of survey methodology. Participation of NCHS staff in the interdisciplinary laboratory, as described later in this plan when discussing collaborative arrangements, will help to bridge the gap that currently exists between government agency survey methodologists and university survey researchers and social scientists. There will be potential benefits for all disciplines. The project will provide additional methodologies for researching cognitive issues in surveys, new phenomena to examine in basic research in the cognitive and related sciences, and tested strategies for producing improvements in the methods and statistics of federal statistical surveys in general, and in NHIS, in particular.

A note of caution is in order about the potential benefits of this project. It is not expected that the project will produce definitive substantive findings with respect to any cognitive issues, although it may provide important leads for subsequent research. The major emphasis will be methodological rather than substantive. Even so, it is recognized that the methodological findings obtained in this or any single study will not be conclusive until verified by other researchers in subsequent trials.

### **Collaborative Arrangements**

This demonstration project will be conducted in a collaborative mode. Since NCHS has neither a cognitive research laboratory nor a staff of cognitive scientists, it will make contractual arrangements with universities to have the experiments conducted in their laboratories and with their scientists. Not only will this arrangement be cost-effective for this project, but it will, as noted earlier, have the major long-term benefit of establishing closer research ties between the federal statistical establishment and universities.

The NCHS and university laboratory staffs will collaborate in all research aspects of this project, and in the preparation of research reports, many of which will be suitable for publication in scholarly



journals. The NCHS will be primarily responsible for the survey and statistical methods and the contractor for the cognitive methods.

### **Work Plan**

Experiments will be conducted to test the application of laboratory-based methods for two broad types of questionnaire design problems:

- (1) development of survey instruments
- (2) investigation of specific cognitive issues

The NHIS questionnaire will be used by the laboratory as the survey instrument for both types of experiments. The supplement to the NHIS questionnaire will be used to test the development of survey instruments. This part of the questionnaire collects information on specific health topics (child care, health promotion, prescribed medicine, etc.) and changes annually. The specific cognitive issues will be generic to surveys and could arise in either the NHIS supplement or the core of the NHIS questionnaire. The latter collects basic information about the nation's health (health status, utilization of health services, etc.) and undergoes virtually no change from year to year.

The workplans for developing and pretesting a supplement to the NHIS questionnaire and for conducting laboratory research on specific cognitive issues relating to the NHIS questionnaire, respectively, are discussed in the next two sections. These plans were developed within the context of the collaborative mode in which the project will be conducted. On the one hand, these plans are intended not to overly restrict the Center's or the subcontractor's creativity as the research progresses. This is a concession to the nature of this research project, and also a major advantage of having the laboratory research conducted outside NCHS. On the other hand, the plans were developed with the view to pursuing certain objectives and producing particular products within a specified time frame. This is a requirement necessary to ensure project accountability, and also a major advantage of having the project administered by NCHS.

### **Survey Instrument Development**

The Center's schedule of activities to develop and test the 1987 NHIS supplement and a proposed schedule of the laboratory's activities are presented below.

<u>Date</u>	<u>Schedule of NHIS Activities</u>	<u>Schedule of Laboratory Activities</u>
1/85	Develop analysis plan	
6/85	Complete first draft of supplement	Complete first draft of supplement
10/85	Prepare pretest version of supplement	Complete testing of first draft of supplement
12/85	Start OMB pretest clearance	--
3/86	--	Complete testing of pretest version of the supplement
4/86	Conduct field pretest	
	Participate in field pretest	
6/86	Prepare pilot study version of supplement	(contingent on a three-month extension of the laboratory subcontract)
7/86	Start of OMB pilot study clearance	
10/86	Conduct pilot study	
10/86	Design NHIS supplement	
1/87	Start 1987 NHIS	

In accordance with the Center's established timetable for constructing its annual supplements to the NHIS questionnaire, the topic for the 1987 NHIS supplement will be selected during 1984 and the literature search will be completed by January 1985, exactly two years before the NHIS commences. During the two-year period, January 1985-1987, the Center staff will be engaged in the tightly scheduled set of activities as noted above. These activities exemplify the traditional method of constructing survey instruments. The sine qua non of this method is that the instruments are field tested under conditions that simulate the actual survey conditions as closely as possible. This approach is in sharp contrast to the proposed laboratory activities which would be conducted under controlled laboratory conditions. The schedule of laboratory activities is linked to the NHIS schedule of established activities so as to maximize the laboratory's potential contributions in developing and testing the 1987 NHIS supplement, subject to the condition that these activities should not interfere with nor jeopardize the basic integrity of the NHIS established testing practices.

This phase of the project will delineate the potential role of laboratory-based research in developing and testing survey instruments. The project could result in the development of improved NHIS pretesting protocols, including improved field pretesting methods for training and debriefing interviewers, and "innovative" laboratory-based methods for conducting unstructured interviews and group interviews.

### **Developing the First Draft of the NHIS Supplement**

The laboratory will devote the five-month period, January-May 1985, to developing the first draft of the NHIS supplement. During this same period the NCHS staff will be independently developing its own first draft of the NHIS supplement.

In early January 1985, NCHS will provide the laboratory with the items of information that will be collected in the 1987 NHIS supplement, and the contractor will transform these items into survey questions and procedures. Using cognitive techniques, such as protocol analysis, the laboratory will investigate the manner in which respondents process the required information, and on the basis of these findings, draft questions and design its own first draft of the NHIS supplement.

Comparing the first drafts of the NHIS supplements that are developed separately by the NCHS and the laboratory will indicate the extent to which the laboratory method is a surrogate for the traditional NHIS method and to what extent it produces different results. Merging what appear to be the best features of both versions of the questionnaire and comparing the combined result with the questionnaire that was developed entirely by traditional methods will provide a basis for assessing the enhancement value, if any, of developing questionnaires in the laboratory as an adjunct to, or in place of, traditional developmental methods.

### **Testing the First Draft of the NHIS Supplement**

During the six-month period from June 1985 to November 1985, the first drafts of the NHIS supplement will be pretested in the laboratory. Possibly three different versions will be laboratory tested: one that was developed by the NHIS staff, another developed by the laboratory staff, and possibly a third which incorporates what are judged to be the best features of the other two versions.

The laboratory will assess whether the drafts of the NHIS supplement are eliciting the kinds of information they are supposed to elicit. Laboratory testing will be performed on a variety of subjects who will be selected because they are expected to experience different types of cognitive problems with the questionnaire. The criteria for selecting subjects will depend somewhat on the topic covered by the NHIS supplement, but certainly they will reflect demographic, ethn and economic differences in the population.

The laboratory pretest findings will be discussed with NCHS staff during November 1985, so that they can be incorporated into an improved draft of the questionnaire that the NHIS staff would be preparing to accompany its request for OMB clearance in December 1985 to conduct a field pretest in 1986.

### **Testing the Field Pretest Draft of the NHIS Supplement**

During the four-month period, December 1985 until the subcontract ends in March 1986, the field pretest versions of the NHIS supplements will be pretested in the laboratory and then independently field pretested during April 1986. Pretesting identical drafts of the NHIS supplement by both laboratory and field methods will make it possible to compare and evaluate how well each method assessed whether the NHIS supplement was doing what it is supposed to be doing and, if not, what revisions were needed. Relative costs and turnaround times of conducting pretests by each method would also be compared.

### Specific Cognitive Issues

As noted earlier, the manner in which people handle information is of common interest to cognitive and survey scientists, but objectives and methods of the two sciences are quite different. Survey scientists conduct field experiments to evaluate the quality of responses elicited by survey instruments. Cognitive scientists, on the other hand, establish generalizations about the mental systems people use for processing information by conducting laboratory experiments. The mission of this project is to demonstrate the enhancement value to both scientific fields of conducting laboratory research on particular cognitive issues that have been implicated by survey scientists as adversely affecting the quality of survey responses.

Cross-fertilization of the two scientific fields is the keynote of this project. Cognitive issues that arise in surveys are representative of wider classes of cognitive phenomena that are being studied in cognitive science, but under restricted and unnatural laboratory conditions. Therefore, it is believed that bringing the survey cognitive issues and the survey experience with these issues into the cognitive laboratory will generate ideas for basic and applied research in cognitive science. And feeding the laboratory research findings on these cognitive issues back to survey scientists will, in turn, stimulate the development of improved statistical models of survey measurement errors and improved methods of constructing survey instruments. For cognitive science, the ultimate benefit will be a better understanding of the way people process information, and for survey science, it will be improved control over the cognitive component of survey measurement.

Three well-known, but largely unresolved, survey problems are presented as possible candidates for laboratory research. They are: telescoping, conditioning, and the respondents' perceptions of the confidentiality of their responses. Each problem involves cognitive issues that are poorly understood and, as will be explained later, seem to present interesting material for laboratory research. For example, for unknown reasons the effects of conditioning and telescoping are asymmetric. The conditioning effects of ordering questions or response categories are often more pronounced when ordered one, rather than another, way. Similarly, the telescoping effects of allocating events either to earlier or later periods than those in which they actually occurred usually results in inaccurately allocating fewer events to less recent than to more recent periods.

#### Telescoping

Failure of respondents to recall events and to recall correctly when the events occurred are major sources of error in the collection of survey data. The errors associated with these two cognitive sources are often confounded in surveys, which may help to explain why the classical negative accelerated forgetting curve predicted by cognitive science does not necessarily hold in surveys.

The tendency of survey respondents to allocate events either to earlier or later periods than those in which the events actually occurred is called telescoping. Typically, survey respondents report retrospectively about events that occurred during a reference period, which is a calendar period of specified length that precedes the interview date. The telescoping phenomenon has been observed both for reference periods that are bounded by prior interviews and for unbounded reference periods.

The findings of survey research indicate that unbounded recall has a net forward telescoping effect, that is, more events are shifted forward in time and erroneously reported in the reference period than are shifted backward and erroneously not reported in the reference period. In addition, events that are correctly placed within the reference period tend to be reported as having occurred more recently than they actually did. The importance of the event, the length of the reference period, and the characteristics of respondents all appear to have an effect on the telescoping phenomenon. With bounded recall there appears to be telescoping within the reference period itself, with the net forward effect being greatest for the most recent part of the reference period.

The telescoping phenomenon is representative of a wider class of cognitive phenomena involving temporal judgments. There is no doubt that cognitive scientists appreciate that the process of making temporal judgments is an important component of event memory; however, their understanding of the event-dating process is based primarily on laboratory experiments which involve neither naturally occurring events nor long-term memory. Consequently, it is unclear how well existing cognitive theory on temporal judgment applies to the real world of personal events such as those respondents are asked to recall in surveys. Apparently, the telescoping phenomenon, per se, has not been investigated in the cognitive laboratory, and it is proposed that the survey experience with this phenomenon may offer interesting leads for designing innovative laboratory experiments.

### Conditioning

All scientific investigations are subject to the risk that the measuring instruments will disturb the phenomenon under observation and thereby affect the accuracy of its measurement. In this broad sense, the conditioning concept in survey science is analogous to Heisenberg's uncertainty principle in physics, but without the latter's specificity.

Conditioning in survey research usually refers to the distorting effect of the total survey measurement process on survey responses, but in the narrower sense used here it refers more specifically to the response effects of collecting an item or set of items of information on another item or set of information items. For example, it refers to the response effects of adding one set of questions to another set of questions, such as the effects of NHIS core questions on the questions in the NHIS supplement or vice versa. It also refers to the effects of ordering a particular set of questions or response categories, or reinterviewing the same respondents, as in panel and quality check

surveys. The narrower definition is adopted here because it makes the survey conditioning phenomenon more amenable to laboratory experimentation.

Although there are many examples of conditioning effects in NHIS and other scientific surveys, problems often arise unexpectedly since the phenomenon is not well understood by survey scientists. For example, in his CASM paper, Bradburn refers to a mysterious asymmetric effect of question and response ordering. He notes that a different effect is observed when questions are ordered in one way than when they are ordered in another way. It is proposed that this curious effect of conditioning in surveys may offer leads for designing critical laboratory experiments on the conditioning effects of rotating the order in which the material is presented.

### Perceptions of Confidentiality

The response effects of asking for information about sensitive topics is a major survey concern because (1) policy makers and other users of health survey data often require this type of information, and (2) respondents are usually reluctant to provide this information and the quality of the information reported is often suspect. Examples of sensitive topics are: illicit behavior such as drug use, drunk driving, low-esteem behavior such as excessive drinking, overeating, and diseases with social stigma such as cancer, venereal diseases, tuberculosis, mental illness, etc.

Although scientific surveys subscribe to a strict policy of protecting the confidentiality of the reported information, assurances of this policy are often insufficient to overcome the suspicions of respondents that their responses may be disclosed in an identifiable form to third parties or their reluctance to report socially undesirable behavior to an interviewer. Survey scientists try to reassure respondents by using data collection techniques that seek to preserve the anonymity of the persons for whom sensitive information is reported in household surveys. Although it seems obvious that the success of these techniques would be greatly affected by the respondents' perceptions of the confidentiality protection afforded by these techniques, their perceptions have not been subjected to in-depth research and hence they are largely unknown.

There are three survey techniques often used for preserving respondent anonymity:

- (1) Randomized response--a respondent is simultaneously presented the sensitive question and another non-sensitive question, each of which can be answered yes or no. He/she answers only one question and he/she alone knows which one, because he/she selected the question to be answered by a random process such as flipping a coin.
- (2) Network sampling--the respondent serves as an informant for other persons to whom he/she is linked by virtue of kinship, friendship, or some other designated relationship, but who are otherwise unidentified.

- (3) Self-enumeration--the respondent writes the answers to the sensitive questions on a blank sheet of paper that he/she seals in a self-addressed envelope and mails to the survey organization.

The problem of respondent compliance in surveys on sensitive topics is representative of a wider class of cognitive phenomena in which people are faced with the task and the risk of making decisions on the basis of information that they may neither fully comprehend nor believe. Numerous applications of the anonymity techniques in surveys on sensitive topics have produced mixed results that traditional survey research methods have been unable to explain satisfactorily. It is expected that laboratory research on respondents' perceptions of these techniques under varying conditions may improve the design of surveys on sensitive topics and may lead to an improved understanding of the cognitive processes by which people assess risks on the basis of incomplete information.



## COGNITIVE PROCESSES IN SURVEY RESPONDING: PROJECT SUMMARIES

Roger Tourangeau, William Salter, Roy D'Andrade, Norman Bradburn, and associates

Researchers at the National Opinion Research Center (NORC), Yale University, and the University of Chicago have proposed three interrelated research programs to carry out a series of studies on the cognitive underpinnings of the survey interview process. All three projects share a common framework, which is described here briefly. The framework assumes that respondents in surveys proceed through three major stages in answering survey questions: they interpret the question, retrieve the relevant information, and formulate a response. It is our shared belief that response effects in surveys can best be understood by examining these processes in detail. Since different types of questions make different demands at each stage, we have also adopted a simple scheme for classifying survey questions. We distinguish three broad classes: questions that elicit attitudes or opinions; questions that ask about behaviors; and questions that concern the causes or reasons for behavior. The framework thus suggests nine areas of investigation defined by the three stages of survey responding and the three types of questions.

The NORC research program, developed by Roger Tourangeau, Roy D'Andrade, and Norman Bradburn, is entitled "Cognitive Processes in Survey Responding: Attitudes and Explanations." It concerns two types of survey questions--those concerning attitudes and reasons--and includes studies on all three stages of survey responding--interpretation, retrieval, and judgment. The Yale program, developed by Robert Abelson, is entitled "Cognitive Processes in Survey Responding: Multiple Schemas and the Role of Affect." It deals with the same two classes of survey questions as the NORC research program and offers a complementary perspective on some of the same issues explored there. It dovetails with the NORC work in other ways--it extends the analysis of attitudes in terms of cognitive schemata and incorporates studies on the role of affect in survey responses.

The University of Chicago program, developed by William Salter, Steven Shevell, Lance Rips, and Norman Bradburn, is entitled "Cognitive Processes in Survey Responding: Time and Frequency Estimation." It focuses on the remaining class of survey questions, those that concern behavior. It includes studies on the retrieval and judgment processes and on how these processes interact when respondents must judge the timing or frequency of events.

All three research programs share the interdisciplinary perspective of the Advanced Research Seminar on Cognitive Aspects of Survey Methodology. The research teams for each project include researchers who have done cognitive research or survey research or work in both fields. In addition, all three include a commitment to the replication and extension of laboratory findings to the field setting. Each project incorporates plans for split-ballot studies to be conducted within the context of a national survey.

Although the research has been divided into three separate projects, we expect to carry out the work as collaborators and to hold periodic meetings of the entire group. We believe that collaboration will enrich the quality of all our work and that the resulting whole will be greater than the sum of the parts.

### **NORC Research Program**

Roger Tourangeau, Roy D'Andrade, and Norman Bradburn

We propose a series of laboratory studies that explore cognitive processes in survey responding. The studies focus on questions that elicit attitudes or opinions and on those that ask for the reasons or causes of behavior.

In our view, knowledge and beliefs about attitude issues often form organized cognitive structures, or schemata. For the enduring issues that appear regularly in opinion surveys, several competing schemata are often available within a culture. Many respondents subscribe weakly to several of the available schemata for such issues; which schema guides their response to an attitude question is affected by variations in question wording and context. In a prestudy, we will identify the prevailing schemata for several typical attitude issues. The main studies will explore how question wording and context influence which schema is triggered by the question and, thus, affect how respondents interpret the question, think about the issue, and formulate their answer.

Questions that elicit reasons initiate a process in which potential explanations are generated and then tested for plausibility. Because the generation step is not exhaustive and the testing step not rigorous, respondents often accept the most readily available explanation. Three studies will explore the interpretation of questions that ask for explanations and the process of generating and evaluating potential explanations.

The findings of the laboratory studies will be replicated and extended in a split-ballot field experiment. The research program will help resolve a number of methodological puzzles in the survey research literature, including the effects of question order, wording, and context.

### **Yale Research Program**

Robert Abelson

Research on cognitive processes is planned that will have direct and indirect bearing on the improvement of questionnaires used in surveys of attitudes and beliefs. Early phases involve laboratory and small-sample questionnaire research. In the later phase, national samples of survey respondents will be used.

Present plans call for three lines of research. The first is concerned with "affective carryover effects," in which the feelings associated with the response to one attitude question may influence

responses to later questions. The proposed studies experimentally create differing emotional contexts prior to critical target items. Results could have important implications for context artifacts in surveys. The second research area will address the circumstances under which respondents entertain multiple possible explanations for their own and others' behavior, rather than single "sufficient" explanations. The third research line proposes methods, drawn from cognitive science, for distinguishing between symbolic and instrumental attitudes. In the former, beliefs are more rigidly attached to attitudes, and emotions may play a large role. The two types of attitudes may have different relations to demographic and other variables and may be important to distinguish in surveys.

### **University of Chicago Research Program**

Lance Rips, Norman Bradburn, Steven Shevell, and William Salter

This research will examine the cognitive processes and knowledge representations that are implicit in responses to many survey items.

Survey questions that appear to require simple recall often involve complex estimation (e.g., "In that last year, how often have you seen or spoken to a medical doctor or assistant about your illness?"). This research will investigate the mental processes used to make estimates and will explore differences in response accuracy achieved by using various estimation strategies. One set of studies will examine processes and knowledge structures used to estimate dates and durations. Other research will investigate estimates of quantities. Additional work will explore processes used by respondents when revising initial estimates. Laboratory results will be tested in field settings. This research will guide improvements in survey design that can substantially increase the precision of results from surveys.

## THE INTERSECTION OF PERSONAL AND NATIONAL HISTORY

Howard Schuman and Philip Converse

Each of us carries in our memory accounts of those events that make up the recent history of our society. Depending on our age, we have memories and reconstructed images of the Depression, the sixties, the Vietnam War, the Carter years, and so forth. Some of these memories are based on proximate contact, as in the case of someone who served as a combat soldier in Vietnam or was an active demonstrator against that war. Other memories are more distal, as in the case of those who watched both the war and the demonstrations on television.

History at the individual level is necessarily fragmented and idiosyncratic, but it is also vivid and personally important in a way that the historian's larger view of events can seldom be for readers. Therefore, it is likely to be significant as a person faces new events and problems that require interpretation and solution. This is widely recognized at important national decision points, for example, in the form of debates about whether, say, Lebanon is or is not similar to Vietnam, and what that means. Vietnam becomes the model not only because of whatever similarities it has to the new problematic event, but also because Vietnam is resonant both to commentators and to much of their audience.

We are proposing to investigate systematically the versions of recent history that are in the heads of the American population. Such an investigation is of interest in its own right and will also be valuable in helping us to understand how the American population responds to new events, those present and those soon to appear. From a practical point of view, this latter concern has to do with how personal experience of the past influences views of, and therefore actions toward, the present and the future. We do not assume that this process of drawing on the past is either automatic or simple, since many other factors--novel events, charismatic leaders, unanticipated external forces--play a part and may in some circumstances be more crucial. What we do assume is that models of the past, as personally experienced and learned, are one important factor in evaluating new experiences and are therefore worth studying systematically.

We expect perceptions, interpretations, and inferences from the past to be patterned to an important degree along social and cultural lines. Most fundamental of all should be age itself. Our investigation will attempt to capture the range from being in the midst of current news to learning history that is embodied in the dimension of age. Given the exploratory nature of the study, there will be an uncommon interest even in the straight distributions of many of the responses, such as the relative salience of diverse familiar historical events, the levels of personal and national importance ascribed to each, and the nature of the impacts felt to have resulted. Nonetheless, given the time-dependent focus of the study, it will be important to do much analytic work having partitioned the sample by age cohorts.

Although age is expected to be the most basic explanatory variable in this study, there is good reason to expect other societal divisions to enter

significantly. For example, the civil rights changes of the 1950s and early 1960s should be more salient to black than to white Americans, and to white southerners than to white non-southerners.

Other societal divisions may play a part in the perception of other events: high unemployment for those who have themselves been unemployed, the women's movement for committed feminists, etc. In sum, our study will need to obtain and use a number of standard and probably some not-so-standard ways of separating the total population along lines relevant to differential personal experience and involvement in the past. Some other individual difference measures, such as global versus more detailed schema construction, may also be attempted.

### **What Population Would be Studied and How**

We are proposing to translate these general aims into a concrete research project by interviewing a probability sample of the United States adult population. We would like a sample size of 1,600, which would provide an adequate number of respondents at each decade of life (20s, 30s, 40s, 50s, 60s, 70 and above). If the initial investigation proves fruitful enough, it could well be extended to other age groups and cultures.

The study would be carried out using experienced professional Survey Research Center telephone interviewers. This will considerably reduce the cost over face-to-face interviewing, without appreciably reducing coverage or response rates. Extra training in nondirective probing will be provided because of the open-ended nature of several of our central questions. Our design will allow us to obtain interviews in blocks of 400 at different points over a six- to twelve-month period, and will therefore provide data to separate the momentarily salient from the event or change that has more importance. (This same feature of the design permits a study of recency effects, as described below.)

### **The Questions**

The questions developed for this investigation would be of several types. Since we are interested in obtaining the frame of reference that respondents use in organizing their memories of the past, it is important to allow maximum scope initially for spontaneous expression. At the same time, some more standardized questions are necessary in order to be able to make systematic comparisons in terms of certain major events and changes.

### **Initial Frame of Reference**

Our first goal would be to obtain from each individual his or her own identification of the major events of the past 50 years. The following question, already pretested and piloted, would be asked at or very near the beginning of the interview.

“First, I’d like you to think back and tell me what you personally feel have been the one or two most important national or world events or changes in the past 50 years.”

(If only one given: “Is there any other national or world event or change over the past 50 years that you personally feel was important?”)

The question is stated in a general enough way so that it applies both to older respondents who have lived through the full period and to younger respondents who know parts of it only second- or third-hand. In addition, both event and change are stressed, since we are equally interested in discrete events and more diffuse changes.

The answers to this question will provide initial evidence of the organizing scheme in the minds of respondents. We expect some events to be fairly frequent across most types of respondents, especially when age is held constant, while mention of other events should be characteristic of particular groups or categories of people.

### **Significance of Events**

Central to our investigation is how people remember/reconstruct the meaning of events. Do they see them as successful or unsuccessful in outcome, symbolic of something good or bad, providing a positive or a negative lesson for the future? Questions about perceived importance will be asked for each of the events mentioned.

### **Personal Connections**

One of our main interests is in whether and how the events or changes mentioned by respondents are tied to their own lives. Our pilot work indicates that some people note this spontaneously, but for those who do not, we will pose questions about how these events changed their lives or their way of thinking about things.

### **Opinions of the Past**

Although the concern of this research is primarily with memories, it is clear that we shall seldom be in a position to make direct comparisons between these memories and what they are supposed to be memories of. Indeed, this is not really the purpose. If age produces relationships that are meaningful in terms of direct experience vs. secondary or tertiary learning, the reported memories can be attributed to past experience, but of course this is inferential. However, this limitation, which would be a major one for psychologists concerned with memory per se, is not especially serious for the proposed research. In addition to the now generally accepted fact that all complex memories involve a considerable degree of

reconstruction, rather than being essentially photographs of the past, our primary concern is with the memories as such, not with their accuracy.

Nonetheless, for some events--especially controversial ones like the Vietnam War--it would be of interest to learn something about the drift in collective memory as the event has receded in time. Toward this end, we would like to resurrect one or more opinion poll items that were frequently used at the time to measure the division of sentiment, and pose it again to our "modern" respondents, first by way of asking what their current reactions to the event are; and second, by asking them to respond to it as they feel they would have at the time. Such a procedure would pick up any self-recognized drift of opinion, a phenomenon that would be of interest in its own right. And while we would not expect that respondents would be capable of recalling their actual early positions with any accuracy, we would be able to compare "remembered opinion" with distributions actually generated at the time of the event, which would provide some assessment of a more unconscious drift of opinion retrospectively. Not only could gross distributions be compared, but various background correlations of stable demographic attributes with original and remembered opinion as well. Such a tactic would only be available with respect to a very few events that were controversial at the time and well monitored by opinion polls with standard items, but it would seem important to seize upon such opportunities in this fashion.

### **Secular Societal Trends**

Up to this point, we have dealt with major "events" that have some staccato quality, even though certain of them cannot be very clearly demarcated in time. We are also interested in the felt impacts of some of the major social and technological trends that have characterized the past century. We would proceed in much the same way, asking first in unstructured form for the kinds of things that the respondents feel have been most important in these regards, but then subsequently covering the most central possibilities in more structured or standardized form.

The questions that we have in mind require more time perspective than our younger respondents could be expected to have. Thus we would ask persons 50 or older what they would see as the most important differences in daily life when they were growing up and life today, and why they make the selections they do.

For trends such as these, we would like to know how the respondent feels each has affected his or her life, and whether the change is generally beneficial or undesirable. This line of questioning would naturally conclude with a global question as to whether, all things considered, the quality of everyday life has improved or deteriorated in the respondent's lifetime.



## Background Questions

In addition to standard background questions on age, sex, education, occupation, income, place of origin, marital status, and ethnicity, it would be useful to obtain some factual measures of the respondents and their close relatives' involvement in past events.

## Contributions of the Research

To our knowledge no similar information on what might be called the collective memory has been previously gathered, and we believe the descriptive data yielded by the study will have some intrinsic interest for many social scientists. However, our interest is also analytic, and we will attempt to account for why historical memories vary over the population and how they influence judgments of the future. To give one partial example, our pilot study suggested that there is a gender difference in personal connection to events and changes, with women more likely than men to explain the importance of an event or change in terms of some personal impact on their own lives (e.g., loss of a parent in a war). It is possible that this connection plays a role in the frequently reported tendency for women to be more reluctant than men to support new military actions.

The research can also be viewed as a study of social memory that parallels laboratory studies of memory. Those studies document both primacy and recency effects in single-session tests of memory (Martindale, 1981). We expect also to find at least traces of primacy and recency effects in these long-term natural memories. Primacy effects should come about because of sheer rehearsal, much as in laboratory studies. Recency effects, on the other hand, cannot be due to short-term memory, as is assumed in laboratory studies, but to the temporary rehearsal of events that have recently occurred. In our pilot study, for example, mention was made of the Korean Airlines plane that had recently been shot down, but we doubt that that would occur even a month or so later. Since our design calls for gathering our interviews over at least six months, we will be able to go some distance in separating recency effects as such from the intrinsic importance of recent events.

The hypothesis of primacy effects connects this research in a broader sense with the hypothesis that major events occurring early in life have the largest impact (Mannheim, 1927)--other things equal. Without necessarily subscribing to any single critical age, we will be able to examine the connection between recalled events and the ages at which they occurred. Do people in their 70s disproportionately cite the Depression, people in their 60s World War II, people in their 50s the Cold War or associated events, etc? The survey data will allow examination of this issue as well as of the effects of such experience on broader political and social judgments about the interpretation of the events and their implications for future events (Lang and Lang, 1978).

Finally, it is true that a good deal of professional history attempts to address the longer-run impacts of both specific events and larger technological and social trends on the populations experiencing them. Some of these impacts are important whether salient, merely cognized, or

neither. At the same time, attention is frequently addressed to the perceived long-run significance of such events and trends on those who have lived through them, typically based on documents and testimony from individuals who are obviously unrepresentative survivors. We think it would be illuminating to assemble the most salient recollections of a group more representative of the nation as a whole.

### References

- Lang, K., and Lang, G.E. 1978 Experience and ideology: the influence of the sixties on the intellectual elite. Pp. 197-230 in Social Movements, Conflicts and Change, Vol. 1. Greenwich, Conn.: JAI Press.
- Mannheim, K. 1927 The problems of generations. Reprinted in Essays on Sociology of Knowledge, 1952. New York: Oxford University Press.
- Martindale, C. 1981 Cognition and Consciousness. Homewood, Ill.: The Dorsey Press.

## A PROPOSAL FOR THE DEVELOPMENT OF A NATIONAL MEMORY INVENTORY

Endel Tulving and S. James Press

This is a proposal for the development of a national memory inventory. Memory capabilities of the population, together with other cognitive abilities, represent an important part of the nation's intellectual resources. No general and systematic information regarding the quality of these resources is available at the present time. For this reason alone, it would be desirable to have an estimate of memory, as well as other cognitive abilities, on a nationwide basis. The inventory would not only provide an objective picture of the current state of these abilities, but also make it possible to monitor changes in memory and cognitive abilities over time, and to relate such changes to the changing age composition of the nation. Age-related memory and cognitive functions would be of particular interest in this context.

The inventory would provide a set of national norms against which memory performance of individuals or groups of individuals and any impairments in such performance can be evaluated. This kind of evaluation is becoming increasingly critical in a society in which a sizeable proportion of older people suffer from various forms of senile dementia, such as Alzheimer's disease. Impairment of memory and cognitive functions is among the earliest symptoms of these dimensions; the detection of such impairment is, therefore, of considerable importance. The existence of national norms of memory, as well as other cognitive abilities, with which the test scores of individuals can be compared, may significantly facilitate clinical assessment of impairment of such ability.

A national inventory of memory could be used for classification of subjects with memory impairments. It could also be used to scale the quality of memory functions of individuals in situations in which such functions play an important role. For instance, it might be possible to "calibrate" eye-witnesses in court trials on the basis of a battery of suitable tests and to develop "weights" for responses given by individual respondents to recall-type questions on sample surveys of different kinds.

This proposal is one of the products of the CASM seminar at St. Michaels and its follow-up meeting in Baltimore. The two authors of the proposal do not wish to claim any proprietary rights to, or special interest in, the further development of the ideas contained herein, or its eventual implementation, although they are willing to collaborate with other interested individuals and agencies in any further development. The major purpose of the proposal is to stimulate and encourage further thought and possible action along the general lines discussed herein.

The proposal consists of two main parts. The first part contains a short background statement about memory and "testing" of memory; a short general description of the proposed battery, together with a listing of criteria used in selecting individual components of the battery; and a summary of the procedures to be used in the collection and analysis of

the data. The second main part of the proposal consists of a description of memory tasks constituting the battery, together with instructions and examples of the kinds of materials that might be used.

### Memory and Memory “Tests”

Psychological study of memory can be approached from two different vantage points. One is that of cognitive psychology. In this view, memory is a set of interrelated cognitive processes that allow a person to acquire, store, and subsequently retrieve information about the world. These processes are described with reference to a typical individual, the “standard rememberer.” The other vantage point is that of psychometrics. According to this view, memory is an ability or skill that individuals possess and with respect to which different individuals vary.

Memory is conceptualized as a unitary entity in neither the cognitive psychology nor the psychometric approach. Rather, both assume that the concept of memory covers a number of different forms or kinds of acquiring and using knowledge and information. These different kinds of memory operate according to somewhat different principles and, at the level of psychometric analyses, show different correlational patterns of individual differences. Thus, students of memory have talked about visual versus auditory memory, verbal memory versus pictorial memory, rote memory versus meaningful memory, voluntary versus involuntary memory, as well as about facial memory, spatial memory, recall memory, recognition memory, and many other sorts of specialized memories.

Two major distinctions concerning different kinds of memory that are useful to make in the present context are those (1) between episodic and semantic memory (Tulving, 1972, 1983) and (2) between primary and secondary memory (Waugh and Norman, 1965; Craik and Levy, 1976). Episodic memory refers to memory for concrete, personally experienced events; semantic memory refers to a person's abstract knowledge of the world. For example, if a person sees and later recalls a familiar word or a drawing of a common object in a memory test or memory experiment, that person's episodic memory is being tested. If, on the other hand, the person is shown the picture of a public figure and asked to name the figure, it is semantic memory that is being assessed. Primary memory (sometimes also labelled short-term memory) refers to memory for perceived stimuli within a few seconds of their presentation, before the representation of the stimuli has completely left the individual's consciousness; secondary memory (sometimes also labelled long-term memory) refers to memory for information that has left the person's consciousness and has to be brought back into it through particular retrieval queries or cues.

Given the complexity of processes and abilities that the term memory covers, it is generally accepted that there is no simple way of measuring people's memory. Certainly there does not exist a single convenient memory test that could be used to assess the memory abilities of a group of individuals. Instead, a battery of tests is necessary to capture different forms and kinds of memory.

The term memory test is somewhat ambiguous. Its meaning can be clarified by drawing the distinction between a memory task and a memory test. Although the term test is frequently used in the psychometric tradition, referring to the whole operation that permits the attachment of a numerical value to a person's performance on a memory task, it actually, or more precisely, refers to just one component of such an operation. A memory test for the material that the person has learned in a particular situation (an episodic memory test) constitutes only the final stage of a memory task that consists of the following sequence of events: (1) a person examines (observes, studies) some material; (2) there is an interpolated interval of variable duration, usually filled with mental activity involving material other than that studied in the first part of the task; and (3) the person is given a test of what he or she remembers from the initial, study phase of the task. (Note that in the realm of semantic memory, there is usually no need to distinguish between the memory task and the memory test: a semantic memory task can consist of nothing else but a test.)

No generally accepted standard battery of memory tasks exists. When a person's memory has to be assessed for clinical purposes--as in cases of known or suspected brain damage--various instruments have been used. The most popular of these is the Wechsler Memory Scale. It consists of seven subtests, some of which are concerned with questions as to the respondent's awareness of and orientation in space and time ("How old are you?" "What day of the month is this?"), and only some of which tap the respondent's memory for newly presented information. But all of these true memory tests measure the respondent's short-term memory only. This characteristic seriously limits the usefulness of the scale. The Wechsler scale was standardized in the 1930s and 1940s (Wechsler, 1945) on approximately 200 haphazardly selected adults between the ages of 25 and 50. A person's overall score on the battery can be evaluated against the distribution of scores from the standardization group, and on the basis of this evaluation it can be expressed as his or her MQ (memory quotient), whose meaning or interpretation is roughly comparable to that of IQ. The clinical condition known as amnesia can be operationally defined in terms of an abnormally large difference between a person's IQ and MQ.

Testing of larger groups of subjects on various memory tasks has been undertaken only in factor-analytic studies of memory (e.g., Kelley, 1964; Underwood, Boruch, and Malmi, 1978). In these studies, several hundred subjects are typically given a large number of memory tasks, and the scores from the tests are used to derive the factor structure of the tests employed. These studies, too, suffer from the limitation of employing almost exclusively short-term memory tasks in which the subjects are tested immediately after the presentation of the to-be-remembered material.

### **A Short Description of the Battery**

It is probably impossible to construct a completely adequate battery of memory tasks at the present time. Because of lack of appropriate

empirical evidence, there is no general agreement as to how many different kinds of tasks would be necessary to assess most of the important aspects of people's memory performance. Different materials and different conditions under which the materials are studied, retained, and tested are known to influence the performance of a given individual relative to that of others. Thus, there exists a potentially very large set of memory tasks, that is, combinations of materials and conditions of their study and test. The exact constitution of the battery that eventually would be used represents one of the many sub-problems that would have to be solved in the course of the project. The scope and organization of such a battery would necessarily have to reflect a compromise between what is scientifically desirable and what is practically feasible. To make the battery suitable for use in a large, heterogeneous population, the materials for each task would need to be carefully screened and tested to minimize biases favoring the performance of one cultural, ethnic, or socioeconomic group over another. The battery described in what follows constitutes only one of many possibilities.

The sample battery consists of two major parts, A and B. **Part A** consists of the study stages, and in some cases immediate (short-term) testing, of six tasks, together with a test for the memory of the order of the six tasks. **Part B** consists of the delayed (long-term) tests of five of the six tasks.

The sequence of events constituting the battery is summarized next. A more complete description of the tasks--materials, instructions, and test forms--will be found in the second major section of the proposal.

### **Part A**

- (1) Low-frequency words--32 words (such as HYDRANT and BLUEBIRD) presented for study. Delayed recognition and word-fragment completion tests are given in **Part B**.
- (2) Line drawings of common objects--24 line drawings of common objects (such as a basket, a glove, and a lion) are presented for naming by the respondent. Delayed recognition test is given in **Part B**.
- (3) Paired associates--12 pairs of words (such as CLAMP-VALET and RURAL-HEAVE) are presented for study and immediate paired-associate test. This study-test procedure is repeated on the second trial, with the same 12 pairs. Delayed paired-associate test is also given in **Part B**.
- (4) Faces--16 faces of unknown people presented for study. Delayed recognition test is given in **Part B**.
- (5) Categorized words--3 familiar words (such as CARP, MINNOW, BARRACUDA, or SPRUCE, POPULAR, WILLOW) from each of 6 different conceptual categories (a total of 18 words) are presented to subjects initially for identification of category membership (fish or trees) and subsequently for an immediate free-recall test of the 18 words. Delayed cued-recall test is also given in **Part B**.

- (6) Short, high-frequency words--15 three-letter words (such as GUN, ART, ILL, BAY) presented for study and immediate free recall. This study-test procedure is repeated on the second trial with the same 15 words.
- (7) Order of tasks--given a descriptive listing of the six tasks of [Part A](#), respondents are asked to reproduce the order in which they encountered the six tasks.

### Part B

- (1) Categorized words--cued recall test. Respondents are given the names of the six categories of words they saw in [Part A\(5\)](#), and they try to recall the three instances presented in each category. The maximum score is 18.
- (2) Faces--two-alternative forced-choice recognition test. Respondents are shown 16 pairs of faces, one pair at a time. Each pair contains one of the faces seen in [Part A\(4\)](#) and a new face. The respondent has to choose one of the faces in each pair as the one he or she saw earlier. The maximum score is 16.
- (3) Paired associates--cued recall test. Respondents are given the left-hand members of each of the 12 pairs of words seen in [Part A\(3\)](#), one word at a time, and their task is to produce the name of the corresponding right-hand member of the pair. The maximum score is 12.
- (4) Line drawings of common objects--four-alternative forced-choice recognition test. Respondents are shown 24 sets of four different line drawings of common objects, one set at a time. Each set depicts an object (such as a basket, or a glove, or a lion) in four different ways. One of these was seen by the respondent in [Part A\(2\)](#), the other three are new. The respondent's task is to select the one he or she saw before. The maximum score is 24.
- (5) Low-frequency words--yes/no recognition test. Respondents are shown 32 words (such as HYDRANT and COPYCAT), one word at a time. Half of these test words appeared in [Part A\(1\)](#), half are new. The respondent's task is to identify each word as old or new. The maximum score for the old test words is 16, for the new test words, 16.
- (6) Low-frequency words--word-fragment completion test. Respondents are given 32 word fragments (such as \_ \_ U \_ B \_ RD and \_ O \_ O \_ UT). Half of these fragments correspond to words they saw in [Part A](#) (such as BLUEBIRD), while the other half belong to words not previously seen in the session. The respondent's task is to complete the fragment by replacing dashes with letters and thus converting the fragment into a word. Note that respondents are not asked to produce words that they saw in [Part A](#), their task is to produce the word that fits the fragment. Fragments are so constructed that they fit only one word in English. The maximum score for each of the two subsets of test words is 16.

### Criteria for the Selection of Tasks

The criteria governing the selection of tasks for the instrument (the battery) proposed here include the following:



- (1) Emphasis on long-term episodic memory. The tasks in the battery are primarily concerned with “memory proper,” that is, long-term episodic memory.
- (2) Multiple materials. The tasks in the battery tap memory for both verbal and nonverbal material.
- (3) Multiple test types. The tasks in the battery assess both recall memory and recognition memory.
- (4) Sensitivity to age differences. The tasks in the battery include those that are expected to be sensitive to age differences in the population as well as those that are not, or are less sensitive (Craik, 1977).
- (5) Length of the testing period. It should be possible to administer the whole battery in a single session of approximately one hour's duration.
- (6) No special equipment. The battery could be administered under less-than-perfect laboratory conditions without any special equipment.
- (7) Group testing. The administration of the battery could be modified to make it possible to test small groups of respondents simultaneously, if such a procedure has certain practical advantages.
- (8) Alternative response modes. The battery consists of tasks in which either oral or written responses could be given by the respondents without greatly biasing the results.
- (9) Range of scores. The tasks in the battery can be fine-tuned in pretesting to minimize “ceiling effects” in performance while permitting a very large majority of respondents to perform in a way that would justify the examiner to provide occasional positive encouragement to the respondent.
- (10) Alternative forms. It is possible to construct alternative forms of the battery, entailing different versions of the same tasks, that would yield comparable normative data from the population.
- (11) Clinical use. In addition to assessment of memory abilities of samples of the general population, the battery can be used for clinical evaluation of individuals with milder forms of memory impairment.

### General Characteristics of the Battery

The battery is designed to measure both short-term and long-term memory. In one of the tasks (A6), only short-term tests are given, albeit on two separate learning trials. In four of the tasks (B2, B4, B5, and B6), only long-term tests are given. In four other tasks (A3, A5, B1, and B3), both short-term and long-term tests are given. Delayed (long-term) tests are given for tasks in which the respondents are unlikely to be confused as to exactly what it is that they have to try to remember in any given test. For this reason, no delayed free-recall tests are included in the battery.

The test for the order of the six tasks of Part A (A7) is included as an attempt to assess “pure” episodic memory: memory for the temporal sequence of otherwise easily remembered personal events. Other tests tap only the respondent's knowledge of the semantic contents of these events. It is the only test that is likely to produce ceiling effects,

but the data from it may be revealing in cases where subjects do make errors on the test. The test takes only a little time to administer.

Three tasks in the instrument probe people's recognition memory for verbal materials (B5), appearances of objects (B4), and faces (B2). For both of the two latter kinds of materials (B4 and B2), verbal mediation in remembering is precluded: remembering just the name of the originally perceived object will not permit the subject to choose the correct alternative in the test, since all four test items in a set have the same name, and faces are difficult to code verbally to begin with. Yet two different recognition tasks tapping purely visual memory (line drawings and faces) are included because it is quite possible that the ability to remember faces is not highly correlated with the ability to remember the appearances of other visually perceived objects. All recognition tests are given after longer retention intervals: immediate tests of relatively small sets of to-be-remembered materials would be subject to ceiling effects.

Recall tests in the battery are of two kinds, free recall (A5 and A6) and cued recall (A3, B1, and B3). In free recall, the respondents' task is to produce as many studied items as possible, in any order, in response to general instructions to do so. In cued recall, subjects are provided with specific cues for recall of individual items. Two kinds of relations between cues and items to be recalled (or, between cues and targets) are studied in cued recall tests. In the paired-associate task, the cue consists of a word semantically unrelated to the target that was paired with the to-be-recalled word at study; in the categorized words task, each cue, associated with three target words, is represented by the semantically meaningful category name.

The fragment completion test (B6) is known to show the effects of memory in a fashion uncorrelated with other measures of memory (Tulving, Schacter, and Stark, 1982) and is therefore of especial interest. Respondents' performance with the "new" fragments of this test provides a measure of one aspect of their semantic memory.

### **Procedures of Data Collection and Analysis**

The idea is to administer such a battery of memory tasks to a large probability sample of the U.S. population. Before a large national sampling of the country is carried out, it would of course be necessary to do pilot testing with small samples consisting of several hundred subjects. In addition, some preliminary studies would be necessary to establish and test the feasibility of the statistical procedures for analysis of the data collected.

A number of procedural problems must be solved in the course of the pursuit of the ultimate objective of this proposal. Some of these involve improving the currently available fundamental understanding of the determinants of people's performance on memory tasks, and, therefore, involve cognitive psychology as well as other aspects of psychology. Others involve computer science, psychometrics, sociology, and statistics. The procedures to be developed would involve at least the steps indicated in the following paragraphs.

### Identification of Indicators

A set of  $p$ -variables that can be measured by individual testing are derived from the tasks that constitute the battery. These variables are indicators of the dimensions that characterize the diverse aspects of human memory. The measured performance of a subject on various tasks is referred to as the raw scores. Some psychometric “unfolding” techniques could be used here to determine an appropriate space of  $q$  (lower than  $p$ ) dimensions. (For some earlier work in this area, see Underwood et al., 1978.) Techniques such as factor analysis, principal components analysis, multidimensional scaling, etc., could be used based upon  $p$ -tests that probably are not orthogonal in terms of memory characterization. The results of such unfolding analysis is a  $q$ -vector of “reduced scores” (factor scores) for any tested individual,  $x:(q \times 1)$ . That is, a low dimensional space of  $q$  dimensions will be determined that is sufficient to characterize the determinants of memory performance of interest. Each subject will then be scored in the  $q$ -space.

### Data Collection

A pilot survey should be carried out to collect  $p$ -vectors of raw scores for  $N$  people. (The raw scores can later be converted into reduced scores.) The  $N$  people should be selected by stratified random sampling, stratification to take place on variables known to be correlates of memory variations (for example, age). Such stratification variables need to be identified; a sampling frame needs to be constructed; an appropriate sampling procedure must be established; and a final survey instrument needs to be constructed. This should all be done on a small scale before anything major is attempted. The survey instrument would include a questionnaire designed to provide background information on each subject appropriate to relating reduced memory scores ( $x$ ) to background variables (memory correlates).

### Population Determination

A important part of the procedure entails establishing what might be called a typical or normal memory. There will be a multidimensional distribution of such typical memories. Subsequently, a clustering analysis should be undertaken of the  $N$  reduced-score vectors, each defined in  $q$ -space, in order to try to separate the  $N$  points into two populations, typical and atypical, or perhaps several well-defined groups, instead of just two. Scores in the typical group (if there were only two) could be used to establish the distribution of a “standard” population, as the forerunner of national norms. The atypical group (or groups) would have its own distribution and norms.

## Distributional Analysis

This aspect of the procedure involves an analysis to determine the empirical form of the distributions of the various clusters. An attempt should be made to establish theoretical population functional forms and fit parameters.

For example, suppose, just for simplicity of explication, it turned out that for each population cluster (i) of interest, the vectors of reduced scores for subjects,  $x:x_{xl}$ , were all normally distributed  $N(\theta_{i_1}, \sigma_{i_1})$ . That is, the vectors of reduced scores for people with typical memories each followed a multivariate normal distribution with mean vector  $\theta_{i_1}$  and covariance matrix  $\sigma_{i_1}$ . People with memories that are in some atypical cluster have reduced-score vectors that also follow a multivariate normal distribution, but with some other mean vector and covariance matrix. If there were two such atypical clusters the population parameters might be  $(\theta_{i_2}, \sigma_{i_2})$  and  $(\theta_{i_3}, \sigma_{i_3})$ , respectively, for each cluster. Once these distributions were established and the population parameters estimated, any individual person's total memory performance could readily be classified by conventional statistical classification techniques into one of the available memory clusters. It will be necessary to define a sample of individuals who will be used to establish the basic populations that will then be used for future classification purposes.

### The Memory Battery: Instructions, Materials, and Tests

A somewhat more complete description of a possible battery of memory tasks is given next. For each task, instructions to participants, examples of materials to be used, and the nature of the test(s) are provided.

#### Part A

- (1) Low-Frequency Words Instructions: "For your first task of the day, I am going to show you a number of words that you will be asked to remember later on. Since you are going to see many other words later on today, and since we want to be able to talk about the lot I am going to show you now, we will call these words 7- and 8-letter words, because they all contain either 7 or 8 letters. Please pay close attention to each word as it appears as you will have only a few seconds to study it. Are you ready? Here we go!"

**Presentation: 32 words presented at the rate of 3 seconds/word.**

**Sample materials:**

HYDRANT	BLUEBIRD	SMALLPOX	MOONBEAM
NICKNAME	DAFFODIL	CLARINET	PACIFIST
BATHROBE	OCTOPUS	BAGPIPE	PENDULUM
ALMANAC	MACKEREL	COCONUT	MOLECULE

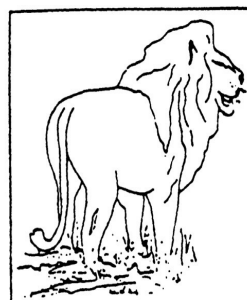
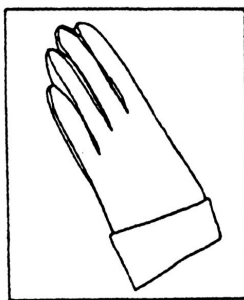
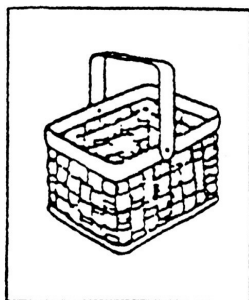
Transition: "This is the end of those 7- and 8-letter words. As I said, we will come back to them later. We now go on to the next task."

- (2) Line Drawings of Common Objects Instructions: "This time I am going to show you a number of pictures of common objects, one picture at a time. When you see a picture, you should name the object that the picture depicts. Later on, I will ask you to remember these pictures. But when you first see each picture, just call its name out aloud, using a short, general description. (Demonstrate: tree) For instance, you would call this picture a tree. Here is another picture. (Demonstrate: jug) What would you call this one? That's correct. Jug, or water jug, is just right. You've got the idea. (Or correct the subject.) So, let us begin with the real pictures. Just call out their names now."

**Presentation:** Present 24 pictures at the rate of 8 seconds/picture.

**Materials:** 24 line drawings of common objects (sample objects shown below):

BASKET	GLOVE	LION	LEAF
MOUSE	BELT	COATHANGER	FROG
TIE	BED	BATH	CHAIR
CANDLE	SNOWMAN	BOTTLE	BUTTERFLY
BELL	ZEBRA	HAMMER	AIRPLANE
RIFLE	SAILBOAT	FLAG	KEY



Transition: "That's all for the pictures. Again, we will return to them later."

- (3) Paired Associates Instructions: "Your next task involves memory for pairs of words. Each pair contains two five-letter words that will be shown together for a few seconds. After you have seen all the pairs, I will test your memory for them by showing you the first word of each pair and asking you to recall the second word that was paired with OK? Here we go!"

**Presentation--Trial 1: Present 12 pairs of words at the rate of 3 seconds/pair.**

**Materials:**

CLAMP-VALET	RURAL-HEAVE	CONIC-ABOUT	SPICE-DUMPY
ULCER-CHIME	STORE-HITCH	RAPID-BLUNT	STALK-PORCH
ALIVE-GLORY	QUASH-FIBER	STAGE-SHADE	HEADY-FINAL

Immediate test 1: "Here comes the test. I will show you the first word of each pair and you will try to recall the second. Do not worry if you do not get too many of them right, it is a difficult test. Here we go!"

Presentation and immediate test--Trial 2: Repeat procedure of Trial 1: Present 12 pairs, and test them as on Trial 1.

Transition: "You are doing all right. We change the pace again, and for the next task I will be showing you some photographs of people's faces."

- (4) Faces Instructions: "This task, as I said involves memory for photographs of people's faces. I will show you a number of faces and later on ask you to recognize them. Pay close attention to each face for you'll see it only once, very briefly. Are you ready? Here we go!"

Presentation: Show 16 photographs at the rate of 2 seconds/photograph.

Materials: A selection of black-and-white photographs of faces of people. The photographs are similar to those that might be used in a yearbook for a large school: small portraits with little other than facial features to distinguish one person from another (no examples shown).

Transition: "I will test you for these faces later on today. We go on now to the next task."

- (5) Categorized Words Instructions: "For this task you will see groups of words, 3 words at a time. The three words in each group belong to a particular category. For instance, if the words were 'London; San Francisco; Tokyo,' the category would be 'cities.' Your task is to identify and tell me the category to which the words in each group of three belong. Just give me a brief label of each category as I show the words to you. If you cannot think of a suitable common name for the three words, just say so, and I will give it to you myself. After you have seen a number of these categorized words, I will ask you to recall them. So, pay close attention to all three words when you study each group. So, look carefully at all words in each group, name the category, and later on recall the words. Is this clear? If so, let us proceed."



**Materials:**

WASP	LAWYER	CANARY	ASPARAGUS	ORCHID	BADMINTON
MOTH	ACCOUNTANT	HAWK	CELERY	AZALEA	WRESTLING
COCKROACH	FARMER	ORIOLE	TURNIP	ZINNIA	VOLLEYBALL
(Insects)	(Jobs)	(Birds)	(Vegetables)	(Flowers)	(Sports)

Immediate test: "All right. Let us see now how many of these words that I showed you, you can remember. Tell me as many words now as you remember. Tell me only the words that I showed you, and NOT the category names that you yourself provided. Go ahead."

Give the subject 60 seconds for recall.

Transition: "That's fine. That's all for the recall of these categorized words. Let's go on to the next task."

- (6) Short Words Instructions: "For the next task, I will show you a number of very short words. Each word consists only of three letters. Look at each word carefully and try to remember it. After you have seen the lot, I will ask you to recall them. You do not have to remember the order in which the words appear; when you recall them, you can recall them in any order that they occur to you. I will ask you to begin as soon as you have seen the last word in the lot, so be ready. Any questions? If not, get ready for the first word."

**Materials:**

GUN	ART	ILL	BAY	LID
BIT	OWN	CUT	ROB	DIM
SET	EAR	TRY	FEE	WIN

Immediate test--Trial 1: "Go ahead, tell me all the words you remember."

Give the subject 60 seconds for recall.

Instructions and test--Trial 2: "We will try this list once more. I will show you the same set of short words again, and again, when you've seen the last one, you try to recall as many of them as you can, in any order in which they occur to you. When you recall the words the second time, recall everything that you can from the whole list, including those words that you already got the first time around. OK? Here we go!"



Present and test the list in the same way as in Trial 1.

Transition: "That's fine. You are doing all right."

## Part B

- (1) Order of Tasks Instructions: "Now we are done with studying and looking at different materials. In the second part of the session, I will ask you to remember the materials that you saw in the first part. The first thing I would like you to do is to tell me the order in which you saw different kinds of materials earlier today. There were altogether SIX things you did. They are briefly described on these six cards. Look at these descriptions and then tell me which of these came first, which one second, and so on to the one that you did last. Take a moment or two to refresh your memory for the tasks, and then order them in the way in which they were presented to you earlier."

Test: Six cards presented to the participant with the following descriptions:

- a. Looking at photographs of faces
- b. Looking at line drawings of common objects
- c. Recalling short 3-letter words
- d. Studying 7- and 8-letter words
- e. Categorizing and recalling groups of 3 words
- f. Studying and recalling pairs of five-letter words

Participant is given up to 2 minutes to order the cards; the order is recorded.

- (2) Categorized Words Instructions: "I am now going to test your memory for the words belonging to categories that you studied and recalled earlier. To help you recall the words, I am going to name all the categories that you saw. These category names are printed on this sheet. You remember that there were three words in each category that you saw. (Give subject the cardboard with the category names.) Go ahead now and tell me the words from these categories that I showed you earlier. Keep thinking and recalling words until I ask you to stop. Do not guess wildly. If you are ready, go ahead!"

Delayed cued recall test of categorized words: The subject is given a cardboard with the following names of categories typed on it:

- a. Insects
- b. Jobs or professions
- c. Birds
- d. Vegetables
- e. Flowers
- f. Sports

The participant is allowed 2 minutes to recall the previously presented words.

- (3) Faces Instructions: "The next thing I have for you to do is a test of recognition memory involving the faces you saw earlier today. I will show you two faces at a time. One of the faces in each pair is one that you saw earlier and the other one is new. Your task is to say which of the two faces you saw earlier. This time you should guess, if necessary. Simply say 'left' or 'right' to indicate which of the two you think you saw earlier. In addition to choosing one of the two faces in each pair as the one you think you saw earlier, you should also tell me each time whether you actually remember seeing the face or whether you are only guessing. You saw 16 faces earlier, thus you will see 16 test pairs. Any questions? If not, let's begin."

Two-alternative forced-choice recognition test for faces: Present to the participant 16 test pairs of faces, at the rate of 8 seconds/pair. Record the participants' choice, and whether the participant reports remembering or guessing. In each test pair, only one of the faces included will have been seen earlier in the set of faces presented for study in Task A(4).

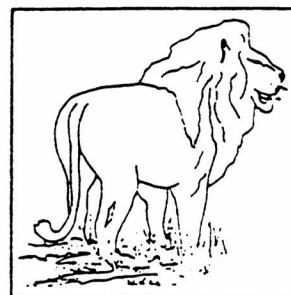
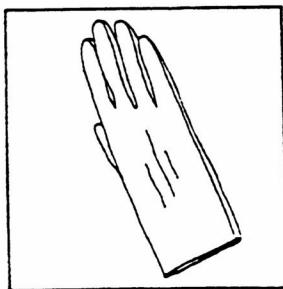
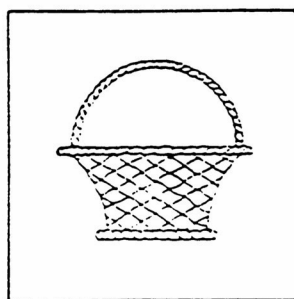
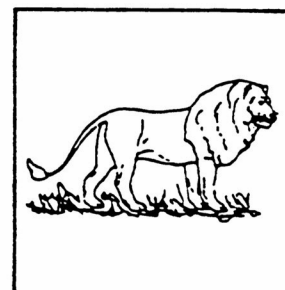
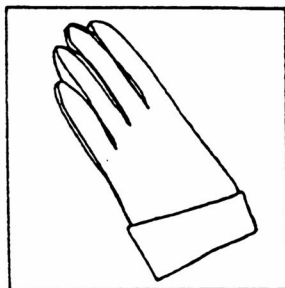
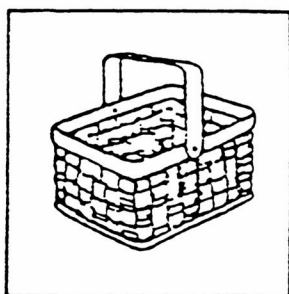
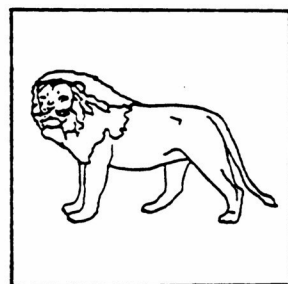
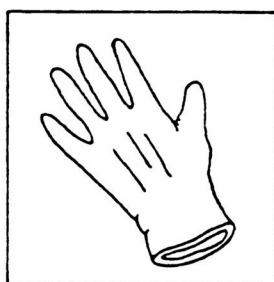
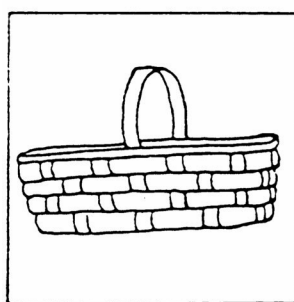
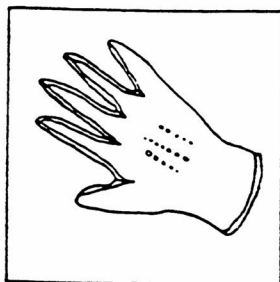
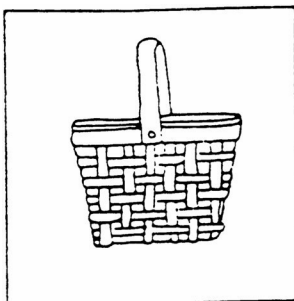
- (4) Paired Associates Instructions: "Next, I am going to test your memory once more for those pairs of five-letter words that you saw and recalled earlier today. This test is exactly like the one that you already took earlier. As before, I will show you the first word of a pair and your task is to try to recall the second member. Don't guess wildly. We will go through the 12 pairs one at a time. Ready?"

Delayed paired-associates test: The subject is given the 12 left-hand members of the pairs, as in the immediate tests, one at a time, and allowed up to 10 seconds to produce the right-hand member of the pair.

- (5) Line-Drawings of Common Objects Instructions: "Your next task is to recognize pictures of objects that you saw earlier today. Each object will be tested by showing you four pictures of the object. Try your best to pick out the exact same picture that you saw earlier. You will note that the four test pictures are labelled A, B, C, and D. Look carefully at all four of them and then tell me the letter of the one that you know or think that you saw earlier. Again, you should choose one picture out of each set of four, guessing if necessary, and you should tell me each time whether you remember the picture you chose or whether you are guessing. Ready? Here we go!"

Four-alternative forced-choice recognition test of pictures of objects: Present to the subject 24 sets of test drawings and allow the subject up to 10 seconds/set to choose one of the alternatives. Record

subject up to 10 seconds/set to choose one of the alternatives. Record the choice and whether the subject reports remembering or guessing.



About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.

- (6) Low-Frequency Words: Recognition Instructions: "Remember those 7- and 8-letter words that you saw way back at the start of today's activities? We are now ready to test your memory for them. It is a recognition memory test. I will show you one word at a time and you tell me whether you remember seeing the word earlier today or not. For each word that I show you, you make the decision and say 'yes' if you remember it and 'no' if you do not. Again, as before, you should also tell me whether you are guessing or whether you are reasonably sure of your decision. So say either 'yes' or 'no' to each test word, and also whether you are reasonably certain of your decision or whether you are guessing. Any questions? Here we go!"

Yes/no recognition test of low-frequency words: Show the subject 16 test words, one at a time, allowing up to 5 seconds per word for the subject to make the yes/no decision and three seconds to report the confidence judgment. Record the decision and whether the subject says he or she is remembering or guessing.

**Sample test words:**

DUCKLING	HYDRANT	SMALLPOX	OMELETTE
NICKNAME	BLIZZARD	MEMBRANE	DAFFODIL
MOSQUITO	BLADDER	PACIFIST	PENDULUM
LETTUCE	MACKEREL	MOLECULE	APRICOT

**Transition: "You are doing fine. We are almost finished."**

- (7) Low-Frequency Words: Fragment Completion Instructions: "This is the last task I am asking you to do today. It involves completing of words from which some letters have been deleted. I will show you a number of such incomplete words and you will try to guess what the word is by mentally filling the missing letters. (Demonstrate) For instance, look at this card: (CH\_PM\_NK). What is the word? Good. (Or: It's CHIPMUNK; do you see it?) Try another one. (\_EMOC\_AT). What's this one? Right. DEMOCRAT. Got the idea? All right, let's start then. Do the best you can, and do not worry if you do not get too many of them."

Word-fragment completion test of low-frequency words: Present to the subject 16 word fragments, one at a time, and allow a maximum of 15 seconds/word for completion of the fragment.

**Sample fragments are presented below:**

_ _ U _ B _ RD	P _ _ _ FF _ N	MO _ _ B _ _ M	AN _ _ _ MY
_ L _ R _ _ ET	B _ GP _ _ E	FL _ _ _ EL	ASB _ _ _ O _
O _ T _ _ US	KN _ P _ _ K	_ _ TH _ OB _	C _ TL _ R _
RA _ _ B _ W	BA _ _ E _ OR	AL _ _ N _ C	_ O _ O _ UT

**Final word: "That's all. Thank you very much."**

### References

- Craik, F.I.M. 1977 Age differences in human memory. In J.E. Birren and K.W. Schaie, eds., Handbook of the Psychology of Aging. New York: Van Nostrand Reinhold.
- Craik, F.I.M, and Levy, B.A. 1976 The concept of primary memory. In W.K. Estes, ed., Handbook of Learning and Cognitive Processes, Vol 4. Hillsdale, N.J.: Lawrence Erlbaum Associates.
- Kelley, H.P. 1964 Memory abilities: a factor analysis. Psychometric Monographs 11.
- Tulving, E. 1972 Episodic and semantic memory. In E. Tulving and W. Donaldson, eds., Organization of Memory. New York: Academic Press.
- 1983 Elements of Episodic Memory. New York: Oxford University Press.
- Tulving, E., Schacter, D.L., and Stark, H.A. 1982 Priming effects in word-fragment completion are independent of recognition memory. Journal of Experimental Psychology: Learning, Memory, and Cognition 8:336-342.
- Underwood, B.J., Boruch, R.F., and Malmi, R.A. 1978 Composition of episodic memory. Journal of Experimental Psychology: General 107:393-419.
- Waugh, N.C., and Norman, D.A. 1965 Primary memory. Psychological Review 72:89-104.
- Wechsler, D. 1945 A standardized memory scale for clinical use. Journal of Psychology 19:87-95.

## PROTOCOL ANALYSIS OF RESPONSES TO SURVEY RECALL QUESTIONS

Elizabeth Loftus

One idea that was received enthusiastically at the St. Michaels seminar was the suggestion for using protocol analysis to study how survey respondents retrieve information from memory to respond to questions about past events. The first section of this item, prepared a few weeks after the seminar, expands on this idea and presents some results from five pilot interviews conducted to explore the potential utility of this technique.

Subsequently, David Fathi, a student of the author, has conducted research in this same area for an honors thesis. The second section of this item describes some results from 23 protocol analyses of responses to two questions; one on use of health care facilities and one on deposits to a credit account. One of the questions arising in this research was how to relate different methods of retrieval to the validity of responses. Since it proved difficult to obtain verification data for the health and credit deposit questions, some subsequent work by Fathi has related to questions asking students in an undergraduate psychology course to recall the exact dates of examinations given in the course.

### Experimental Project: Protocol Analysis

In many national surveys, respondents are asked to recall personal events from their lives. For example, in the National Health Interview Survey, respondents are asked, "During the past 12 months, about how many times did (you) see or talk to a medical doctor?" In the National Crime Survey, respondents are asked, "In the last six months, did anyone beat you up, attack you, or hit you with something, such as a rock or bottle?" Very little is known about the precise strategies for retrieving personal information of this sort.

One method for learning about cognitive strategies is through the use of protocols (Ericsson and Simon, 1980). In the protocol technique, people are asked to think aloud as they answer specific questions. The verbalizations produced are called protocols, and they can subsequently be transcribed and analyzed. This method has an advantage over the similar technique of asking people after the fact to describe how they arrived at a particular answer or estimate. The "after-the-fact" technique has the disadvantage that people often provide reasons or rationalizations for their behavior that are not the true reasons but rather are strategies that subjects believe should have been appropriate (Nisbett and Ross, 1980).

To explore the feasibility of a protocol analysis approach to the problem of how people retrieve personal experiences of the type required on, say, the National Health Interview Survey, we asked five pilot subjects to think aloud while answering specific questions. We first gave subjects some practice questions so they could gain experience in verbalizing their thought processes. These were questions such as, "In



the last 12 months, have you eaten lobster?" Then we asked some health-related questions and some crime victimization questions. For example, we asked subjects, "In the last 12 months, how many times have you gone to a doctor, or a dentist, or a hospital, or utilized any health care specialist or facility?" and, "In the last 12 months, have you been the victim of a crime?"

Many specific questions can be answered by examining the protocols produced by subjects. For example, one specific question is this: Do people answer the health question by starting from the beginning of the 12-month period and moving toward the present (the past to present approach), or do they start from the most recent event and move backward (the present to past approach)? One might predict that respondents would begin with the most recent events, since these might be more "available" in memory (Tversky and Kahneman, 1973). While our results must be considered preliminary, they indicate that, to the contrary, the past to present approach is the favored one. For example, one female respondent answered the health question by saying, "Let's see . . . six . . . six months ago I went to the dentist. Last month I went to the doctor. I think that's it."

If this tendency to prefer the past to present retrieval sequence for the health question were to be documented in a full study, it would suggest that people might be most efficient at retrieving information if prompted to do so by cues that allowed them to start in the past and work toward the present. Of course, this hypothesis would need to be explicitly tested since we know that simply because most people perform acts in a particular way does not necessarily mean this is the most efficient way to do so.

Another specific issue that could be addressed by analysis of the protocols is the extent to which respondents produce new information when asked further questions that relate to ones that were asked earlier. For example, the response of the female quoted above indicates two health related contacts. However, later this respondent was asked, "In the last 12 months have you been to a dentist?" Her answer:

"Let's see . . . I had my teeth cleaned six months ago, and so . . . and then I had them checked three months ago, and I had a tooth . . . yeah, I had a toothache about March . . . yeah. So, yeah, I have." (Interview conducted in July 1983.) This protocol again indicates a preference for the past-to-present retrieval sequence, but also indicates the production of two additional dentist visits that were not provided earlier to the more general question.

Although it is well suspected that additional questions will produce additional instances, it is not known why. Protocols could shed light on this issue. Furthermore, it is not known whether beginning an interview with a general question (e.g., "Have you been to a specialist?") is the optimal technique. It is possible that after having said, "no" to the general questions, subjects may be less likely to search memory in an effort to answer the specific question than they might have had they not been asked the general question to begin with. We simply do not know whether this is the case. However, an examination of protocols given to specific questions that either are or are not preceded by general ones would be more informative.



One interesting observation from the five pilot protocols is the large number of instances in which people change their answer as they are in the midst of speaking. For example, one female subject who was asked the crime question answered: “No, not that I can think of, unless . . . oh, I had two dollars stolen at work, but that's it.” Another said: “No, I haven't, that I can remember . . . Yes, I was--I was thinking about my car, and I had some tapes stolen from my car, in Montlake, about six months ago.”

We could speculate that if subjects had been responding using a more formal checklist technique in which they simply had to say “yes” or “no” that these two instances might never have been reported. Under the more leisurely approach provided by the protocol technique, the instances emerged from memory. One question that naturally comes to mind is whether we can improve on current interviewing techniques to take advantage of this possible discovery. For example, if respondents were asked to think for a minute, and then answer the question, would we be able to accomplish the same benefits within the context of the more typical interviewing procedures?

In short, many interesting issues can be explored through the use of protocols. Specific hypotheses can be tested concerning how personal information is retrieved by people. Moreover, methods for improving the interview process can be tested in this fashion.

### **Order of Retrieval in Free Recall of Autobiographical Material**

Six subjects were asked, “In the last 12 months, how many times have you gone to a doctor, or a dentist, or a hospital, or utilized any health care specialist or facility?” Three subjects were asked, “In the last 12 months, I'd like you to try and recall all the times that you deposited money in your A La Card account\*, and for each time, try and give me the date as accurately as possible,” or a slight variation on this question. Seventeen subjects were asked both these questions, with the A La Card question coming first. Subjects were instructed to “think out loud” as they responded, and their remarks were taped and transcribed verbatim, yielding 23 protocols in response to the health care question and 20 in response to the A La Card question.

Sixteen of the 23 health care protocols and 8 of the 20 A La Card protocols contained fewer than two instances of the behavior in question, thus yielding no information about order of retrieval. For those protocols containing 2 or more instances of the behavior in question, a “+” was assigned for each time a subject went from a temporally more distant instance to one more recent, and a “-” for each time the subject moved from a more recent instance to one in the more distant past. This was regarded as a rough way to quantify the degree to which subjects tended spontaneously to retrieve these autobiographical memories in one direction or another.

---

\*A system under which costs of meals eaten by students are charged against an account to which periodic deposits are made.

When all the + and ` signs were counted up across all the protocols, the results were as follows:

---

Health Care	+	7
	`	3
A La Card	+	23
	`	6

---

When the protocols were classified according to whether the direction of recall was consistently forward (i.e., all +'s), consistently backward (all `s), not clearly in one direction or another (both +'s and `s), or there was no information about direction of recall available (fewer than 2 instances of the behavior in question produced), the results were as follows:

---

Health Care:	all +	4
	all `	1
	both + and `	2
	no information	<u>16</u>
		23
A La Card:	all +	7
	all `	2
	both + and `	3
	no information	<u>8</u>
		20

---

These results were taken as evidence that, at least in response to these questions, subjects tend to retrieve autobiographical memories in a predominantly past-to-present, or forward, direction.

### References

- Ericsson, K.A., and Simon, H.A. 1980 Verbal reports as data. *Psychological Review* 87:215-251.
- Nisbett, R., and Ross, L. 1980 *Human Inference: Strategies and Shortcomings of Social Judgment*. Englewood Cliffs, N.J.: Prentice-Hall.
- Tversky, A., and Kahneman, D. 1973 Availability: a heuristic for judging frequency and probability. *Cognitive Psychology* 5:207-232.

## THOUGHTS AND RESEARCH ON ESTIMATES ABOUT PAST AND FUTURE BEHAVIOR

Lee Ross

### **The Strategy and Tactics of Survey Methodology as They Pertain to Determining Past Actions and Outcomes**

A basic issue raised at the St. Michaels conference concerned the special status of surveys, such as the National Health Interview Survey, which focuses primarily on relatively objective past actions and outcomes rather than opinions, preferences, fears, intentions, or other largely subjective responses. A great deal of the traditional science and art of the survey methodologist has focused on the problem of potential instrument or interviewer bias. Many strategic decisions about instructions, wording of items, and the role of the interviewer are designed to minimize such potential for bias by minimizing the role of both the instrument and interviewer in defining terms, suggesting response strategies, and especially in providing feedback about the responses themselves. While such precautions may be entirely appropriate in the context of political surveying or other attempts to ascertain attitudes, beliefs, or other subjective states of the respondent, they may be less necessary when the responses in question deal with specific concrete past actions and outcomes by the respondent. Furthermore, in attempting to determine the respondent's past actions and outcomes through measurement of his or her recollections or estimates, a rather different set of potential sources of error or bias come into play--i.e., the types of factors with which cognitive psychologists have long been concerned in their study of human memory and judgment.

Cognitive psychologists have identified many factors that impair performance or introduce error in recall or judgment and, perhaps even more pertinent to present concerns, they have identified factors or strategies that lead to improvement in performance. If such psychologists were confronted with the problem of designing instruments and interviewer protocols for facilitating accurate recall and estimation of past actions and outcomes, I suspect that they would worry relatively little about traditional instrument or interviewer effects and a great deal about how to help the respondent remember the events in question or estimate the relevant magnitudes or frequencies associated with such events. They would worry about the impact of the respondents' general theories, schemas, or expectations on their recall or estimates. They would worry about the effect of the respondents' concerns with self-presentation (and perhaps self-perception and evaluation as well). Most importantly, perhaps, they would design instruments and procedures to overcome such obstacles through trial and error testing which measured accuracy of recall or estimation--i.e., compared recollections and estimations to direct measures of the actions and outcomes in question.

I could offer many "radical" suggestions about techniques that could enhance accuracy of recollections or estimates. One might make heavy use of models--i.e., letting the respondent see someone doing a good job of being systematic and complete, or using specific memory aids or mnemonics

in recalling their medical history, consumer behavior, crime victimization experiences, or whatever topic is the focus of the survey. One might encourage the interviewer to offer suggestions about how to remember or estimate the responses or outcomes in question (and train that interviewer in how to prompt such recall). One might explicitly warn the respondent about common biases or errors. One might encourage the respondent to describe events in his or her own words before attempting to answer specific questionnaire items. Perhaps the most extreme possibility would be furnishing the respondent with informative “anchors”—e.g., mean judgments, typical responses, or “base-rates” for people in general or people who are like them in terms of pertinent demographic characteristics.

### Research on “Past” and “Future” Behavioral Estimates for Self and Other

One preliminary piece of research has been undertaken in my laboratory that was prompted by the foregoing comments, even though its direct relevance may not be immediately apparent. It was our thesis that recollection of past events, at least in cases where “episodic recall” is likely to be imperfect at best, or even nonexistent, is closely akin to other related judgment tasks. Specifically, we sought to compare estimates of specific past performances to parallel predictions about future response and to compare estimates and predictions about one's own responses with parallel estimates and predictions about the responses of a peer whom one knows rather well and has ample opportunity to observe on a day-to-day basis.

Ultimately, our concern will be with relative accuracy, and with the relationship between accuracy and confidence, in these four different domains (i.e., self/other x past/future). Pursuing this concern, however, will demand that we choose responses for which we can independently assess actual behavior to which the relevant estimates or predictions can be compared. This will pose significant methodological hurdles and tax our ingenuity. Undoubtedly, it will also restrict our domain of inquiry to responses that are normally recorded (e.g., checks written, purchases made, books checked out of the library, time logged on computer, long distance telephone calls) or at least recordable by an observer (behavior in contrived experimental settings, perhaps television watching, study time, class attendance, etc.)—domains which may or may not be representative of those that figure in survey concerns and domains and which may or may not be typical in terms of difficulty of recall or estimation.

For now, we have chosen to ignore accuracy per se, and to focus on confidence intervals—i.e., to compare subjects' certainty (or, to be more precise, the range of their uncertainty) about the frequency and magnitude of their own past behaviors or outcomes with their certainty about frequencies for parallel future responses. Furthermore, we compare confidence intervals regarding self-estimates and predictions with parallel estimates and predictions for other people. In other words, we are comparing estimates or recollections about one's past behavior, about

which one might have a basis for considerable certainty, with three types of estimate that one would expect to be highly uncertain and heavily “theory based” rather than data based.

In our single research effort, Stanford students were asked to make estimates about a variety of past and future responses either for themselves (self-estimate) or for their roommates (roommate estimates). Items included number of checks written (or to be written) in a 30-day period, money spent in restaurants, hours spent watching television, number of long distance telephone calls, and so forth. For each item they made a best guess and then bracketed that best guess with an upper and lower confidence limit. In fact, they furnished two different confidence limits: 50 percent limits (such that they thought the probability of the upper limit being too low was .25 (or 25 percent) and that the probability of the lower limit being too high was also .25 (or 25 percent) and 80 percent limits (such that the probabilities of the limits being too low or too high were each 10 percent). They also rated the ease or difficulty of making each estimate, and indicated how “surprised” they would be if the “right answer” was not contained within the confidence limits that had been specified (although these items shall not be dealt with in this brief report). The research design made use of both within- and between-subject comparisons--with self versus other a between-subject factor and past versus future a within-subject factor (order of past vs. future was counterbalanced, as was the order of specific items).

The results of this pilot effort (see [Table 1](#)) can be summarized succinctly. Confidence limits for predictions of the future were only

TABLE 1 Relative Width of 50 Percent Confidence Intervals for Estimates of Past Behavior and Predictions of Future Behavior

	Estimate of Past Behavior	Prediction of Future Behavior	Combined
Judgment about self	100	131	115
Judgment about roommate	111	124	117
Combined	105	128	

modestly wider than confidence limits for estimates about the past regardless of whether it was the self or one's roommate who was the target. Furthermore, confidence intervals about one's roommate were no wider, overall, than confidence levels about oneself; for estimates about the past, the confidence interval was slightly wider for roommate than self; for predictions about the future, the reverse was true. None of these main effects or interaction effects, moreover, appear to be statistically significant (although this may change somewhat when data from a second cohort of subjects is added to our analysis). It is particularly noteworthy that the confidence intervals for ostensibly data-based estimates about one's past behavior were only marginally narrower than for the ostensibly theory-based estimates about the future responses of one's roommates. Such data certainly prompt one to wonder exactly how data-based estimates about one's past really are! They also encourage the type of speculation offered earlier--i.e., that the accuracy of any inferences we might want to make about that behavior would be facilitated by procedures that facilitated recall or encouraged more accurate estimation strategies on the part of the respondent. Finally, the data comparing self-estimates and roommate estimates are interesting in their own right, beyond any relevance to concerns of survey methodology. People apparently believe that they can make as accurate and confident estimates about other people's responses as they can about their own, particularly when it is future rather than past responses that are the subject of such estimates. It is obviously tempting to find out whether such relative immodesty regarding one's social predictions, and modesty regarding self-prediction, is justified (just as it is tempting to find out whether one's ability to predict the future is really as good, and one's estimates about the past are really as bad, as suggested by the confidence limits offered in our study). The need for follow-up research is evident--research in which accuracy, and therefore calibration of such confidence intervals, can be assessed directly.

Note: All intervals were transformed to reflect magnitude relative to interval for past behavior of self. Means reported summarize results for 12 behavioral estimates and predictions. A total of 25 subjects offered confidence intervals for self only (both past and future behavior) and an equal number offered confidence intervals for roommate only (both past and future behavior).

## OUTREACH ACTIVITIES

In the months that have elapsed since the January 1984 meeting in Baltimore, several of the seminar participants have been active in publicizing the activities and outputs of CASM and in seeking to encourage others to work in this exciting cross-disciplinary field. The principal method of outreach has been the presentation and publication of papers, but other means are also being used.

On January 26, 1984, Norman Bradburn attended a seminar, "Problems of Measuring Behavior," sponsored by the Economic and Social Research Council in London, England, and presented a paper, "Potential contributions of cognitive sciences to survey questionnaire design." In February Bradburn made presentations on the same subject to three organizations in West Germany: the Zentrum für Umfragen und Methodische Analyse (ZUMA) in Mannheim, the Institut für Demoskopie in Allensbach, and the Max Planck Society in Munich. In July 1984 Roger Tourangeau attended a ZUMA seminar entitled "Social Information Processing and Survey Methodology" and presented a paper, "Question order and context effects."

At the 39th Annual Conference of the American Association for Public Opinion Research held on May 17-20, 1984, at Lake Lawn Lodge, Delevan, Wisconsin, Judith Tanur chaired a session on contributions of cognitive psychology to survey research that included a paper entitled "An information processing approach to recall in surveys" by Norman Bradburn and a paper entitled "Attitude measurement: a cognitive perspective" by Roger Tourangeau. The discussant was Elizabeth Martin of the Bureau of Social Science Research.

At the annual meeting of the American Statistical Association held on August 13-16, 1984, in Philadelphia, Pennsylvania, a session on cognitive aspects of survey methodology was sponsored by the Section on Survey Research Methods and cosponsored by the Social Statistics and Statistical Education Sections. Organized and chaired by Judith Tanur, the session included a paper by Roger Tourangeau entitled "Interchanges between cognitive science and survey methodology" and a paper authored by David C. Fathi, Jonathan W. Schooler, and Elizabeth F. Loftus, presented by Elizabeth Loftus, entitled "Moving survey problems into the cognitive psychology laboratory." The discussants were two of the CASM guests at St. Michaels, Dr. Jacob J. Feldman of the National Center for Health Statistics and Professor Phillip J. Stone of Harvard University.

A paper entitled "Cognitive psychology meets the national survey" by Elizabeth F. Loftus, Stephen E. Fienberg, and Judith M. Tanur has been prepared in response to an invitation from the American Psychologist and is expected to appear in December 1984.

The editor of the Milbank Memorial Fund Quarterly has invited members of the CASM group to prepare papers. A paper entitled "Cognitive aspects of health surveys for public information and policy" is being prepared by Stephen E. Fienberg, Elizabeth F. Loftus, and Judith M. Tanur. A companion piece being prepared by Monroe Sirkin and Judith Lessler stems from their laboratory-based research project at the National Center for Health Statistics. The editor of the Quarterly is inviting several



Health Statistics. The editor of the Quarterly is inviting several discussants for these papers.

In response to a preliminary proposal presented at its board meeting in June 1984, the Social Science Research Council (SSRC) is organizing a working group on cognition and survey research in order to prepare a detailed plan for the activities of a possible SSRC committee that would bear the same name. Cochaired by Robert Abelson and Judith Tanur, the working group includes Roy D'Andrade, Stephen Fienberg, Robert Groves, Robin Hogarth, Don Kinder, and Elizabeth Loftus. The staff person responsible for this effort at SSRC is Robert Pearson, a guest at CASM's Baltimore meeting.

These are the outreach activities that the editors have been able to identify as this report goes to press; there may well be others that have escaped their attention. It seems reasonable to predict, on the basis of this record, that we can look forward to a continuing round of relevant reports and discussions as further results emerge from the cross-disciplinary research programs and activities generated by the CASM project.

## APPENDIX A

### BACKGROUND PAPERS

This appendix contains three papers prepared specifically for the CASM project. The first two--“Cognitive Science and Survey Methods” by Roger Tourangeau and “Potential Contributions of Cognitive Sciences to Survey Questionnaire Design” by Norman M. Bradburn and Catalina Danis--were sent to participants in advance of the St. Michael's seminar. The third paper, “Record Checks for Sample Surveys” by Kent Marquis, is based on a presentation by the author at St. Michaels.

The first two papers focus on the basic questions addressed by the CASM project: What are the potential links between the cognitive sciences and survey research and how can each discipline contribute to progress in the other? However, the authors examine the questions from different perspectives.

Tourangeau identifies four kinds of cognitive operations performed by respondents in survey interviews: understanding questions, retrieving relevant information from memory, making judgments (if called for), and selecting responses. He reviews research in the cognitive sciences on each of the four topics and discusses the possible relevance of research findings to survey interviews. In discussing judgment and response, he gives special attention to the implications of findings from the cognitive sciences for attitude questions in surveys. Tourangeau also discusses the possible uses of surveys as a “laboratory” for cognitive research and explains how using surveys in this way might help to overcome some of the limitations of small-scale laboratory experiments. He describes one instance in which laboratory-based generalizations have been tested in surveys and briefly summarizes the results.

Bradburn and Danis approach the same general subject from the point of view of the survey researcher. The authors first present a conceptual model of response effects, i.e., sources of variation in the quality of data obtained in surveys, and a general model for human information processing. They then contrast the research methods used in the cognitive sciences and in survey research methodological studies. The main part of their paper reviews potential contributions of findings from cognitive research to four of the issues most frequently studied by

survey researchers: question wording, response categories, contextual meaning, and response to survey items that require information on time or frequency of specific events or transactions. For each of these issues, some of the findings from the field of survey research are described, followed by a discussion of possible application of pertinent theories and findings from cognitive research. In the concluding section, the authors caution that the application of findings from cognitive research “to events outside of the laboratory is not direct and requires additional research.”

The paper by Marquis addresses a topic that is of critical importance in attempting to use surveys as a vehicle for cognitive research. Laboratory experiments in the cognitive sciences are generally designed in a manner that permits direct, objective evaluation of the success of subjects in performing specified cognitive tasks. For example, subjects may be exposed to verbal material or other stimuli developed by the experimenter and then asked to recall, recognize, or make judgments about these known stimuli. In survey research, however, the stimuli (in survey research terms, the “truth”) are not in general known to the survey researcher. This makes it difficult, in comparing alternative interview procedures, to determine which ones are most successful in minimizing response effects. Attempts to overcome this difficulty come under the general heading of validity checks, the subject of Marquis's paper. In validity checks, data are sought from external sources, such as official records, that are deemed to be less subject to response effects than are the survey results. These data, either in aggregate or individual form, are compared with data on the same subjects from surveys.

Marquis's paper reviews alternative designs for validity checks and discusses their strengths and weaknesses. He argues that certain designs lead to wrong conclusions about the direction and size of survey reporting bias and gives some examples from record checks associated with health surveys. He concludes that simple forgetting is not necessarily the dominant problem in health surveys.

## COGNITIVE SCIENCES AND SURVEY METHODS

Roger Tourangeau

In the house of cognitive science there are many mansions. In one are the carefully controlled studies of the laboratory; in another, elaborate simulations of cognitive processes by computer scientists. Its subject matter is no less varied than its methods, including topics as diverse as the understanding of language and the inferring of causes, remembering and forgetting, perception, and judgment. It would take a kind of architect of ideas to characterize this rambling structure.

In this review, I shall risk far more undercoverage than any survey researcher. I shall be purposive in my sampling, fitting the selection to my twin aims: my review will focus on the areas of cognitive science that have the most direct bearing on how surveys are conducted, and it will focus on one or two arbitrarily selected areas where the cognitive sciences could benefit most directly from the use of surveys. Much that is relevant will be overlooked. The territory is just too large and too varied for a single foray to include more than a few salient landmarks.

I shall attempt to take two points of view in this survey. In attempting to bring the cognitive sciences to bear on survey methods, I present a cognitive analysis of the task of the respondent. In attempting to bring survey methods to bear on cognitive problems, I shall present the case for adding another mansion to the house of cognitive science. I examine two particular topics--forgetting and "optimism"--that could, I think, benefit from the use of survey samples. My hope is that the reader will fall prey to the amply documented tendency to overgeneralize from a few concrete cases.

### The Respondent's Task

In the usual interview situation, the interviewer reads the respondent a question and a set of response options. The respondent is supposed to attend to and understand the question, recall whatever facts are relevant, make a judgment if the question calls for one, and select a response. This list of cognitive operations is by no means exhaustive: for example, with open-ended questions the respondent must formulate an answer rather than selecting one from the pre-existing options. Sometimes a respondent will short-circuit the process, deciding to refuse to answer rather than to retrieve the relevant facts from memory. Nor is the canonical order--comprehension, retrieval, judgment, and response--invariably followed. The respondent may select the answer before the interviewer even finishes the question. Still, the respondent must usually go through each of these four steps, typically in the prescribed order.

As should be obvious from this description of the respondent's task, there is considerable room for error. Respondents may misunderstand the question or the response categories; they may forget or misremember the crucial information; they may misjudge the information they do recall;

and they may misreport their answer. Consider respondents who are asked whether they have seen a doctor in the past two weeks. First of all, they may misunderstand the reference period: some may think that the question covers the current calendar week plus the preceding one; others may interpret it to mean the period beginning two weeks ago to the day. Even if the respondent does understand which period is being referred to, he or she may find it difficult to determine the exact date bounding the reference period. A study by Elizabeth Loftus (Loftus and Marburger, 1983) indicates that some people may even misreport whether they have had a birthday “within the last six months.” It seems likely that the source of their error is their inability to translate that phrase into a concrete date to which their birthday can be compared.

Even when they understand the question, respondents may forget or misremember the relevant events. Two similar visits to the doctor can be remembered as a single event; Linton's (1982) massive study of her own memory for events in her daily life indicates that the inability to distinguish similar or repeated events is a major source of forgetting.

The respondent who understands the question and who recalls the relevant events may still slip up at the judgment stage. A respondent who recalls a visit to the doctor but who can't quite date it has a judgment to make: Did the visit fall within the reference period or outside of it? Such judgments of recency are affected by a number of factors aside from the actual timing of the event. For example, we tend to judge emotionally salient events as more recent. Even when the respondent can correctly retrieve the relevant facts and correctly judge them, he or she may not report them accurately. We may omit some things to avoid embarrassing ourselves, invent others to avoid disappointing the interviewer. We also resolve ambiguities by recourse to what the situation seems to demand. We ask ourselves whether a telephone call constitutes a consultation with a doctor and our answer may depend on how we perceive the relative demands of completeness versus relevance.

In the next four sections I examine each of these processes--comprehension, recall, judgment, and response--in more detail; I consider the main theoretical approaches used in studying them and attempt to draw out the implications of the theories for the practical problems of survey research.

### Comprehension

Research on comprehension has, for reasons of methodological convenience, concentrated on the comprehension of written material. The same principles, however, are assumed to apply to spoken material as well; I refer to materials of both types simply as “text.”

There are two main approaches to the study of comprehension. One emphasizes the large cognitive structures the reader or listener brings to bear on the text; the other places somewhat more emphasis on the demands of the text itself. The two views are complementary rather than contradictory. I refer to the first as the top-down approach and to the second as the bottom-up approach. Although I try to sharpen the differences between them, the two approaches share several important

notions--the profound impact of context, the use of prior general knowledge by the reader or listener, and the influence of inferential processes during comprehension. The difference in the approaches lies mainly in their views on the nature of the information we use in interpreting a text. The top-down approach emphasizes large pre-existing structures that can organize an entire text; the bottom-up approach emphasizes lower-level structures that can be used piecemeal.

### Top-Down Processing

According to the top-down view, we understand a text by imposing a pre-existing organization on it. Until we impose such a structure, we have considerable difficulty in stitching together successive ideas into a coherent fabric (Bransford and Johnson, 1972):

First you arrange items into groups. Of course one pile may be sufficient depending on how much there is to do. If you have to go somewhere else due to a lack of facilities that is the next step; otherwise you are pretty well set. It is important not to overdo things. That is, it is better to do too few things at once than too many. . . . After the procedure is completed, one arranges the materials, into different groups again. They then can be put into their appropriate places.

Although each sentence in this passage presents no special difficulty, the passage as a whole is nearly impossible to interpret--until, that is, we learn that the passage is about doing the laundry. After that, the sentences fall neatly together into a tidy conceptual package. With most passages we are able to find some conceptual package early on in the text. (The Bransford and Johnson passage is deliberately constructed to prevent this from happening; it uses general terms--"procedure," "items," "different groups," etc.--where particular terms would have tipped us off.) Once the relevant structure is recognized, each succeeding idea finds its niche in the larger edifice. Comprehension is seen by the top-down approach as a process of recognition: first we recognize the general pattern and then we recognize how each piece takes its foreordained place in the pattern. The patterns or structures are of two types. Some theorists (Mandler and Johnson, 1977; Rumelhart, 1975) stress the importance of our knowledge of the general form of texts. We know, for example, that stories consist of settings, beginnings, middles, and endings. We know, further, how each constituent of a story breaks down into smaller constituents (e.g., a setting includes the time, place, and cast of characters) and we know how the constituents relate to each other (e.g., the beginning of the story triggers some goal in the protagonist which he or she then attempts to satisfy in the middle part of the story). Our knowledge of the grammar and logic of stories allows us to fit the events of a story into a coherent structure. Other theorists (Abelson, 1981; Bower et al., 1979; Schank and Abelson, 1977) emphasize more particular knowledge of the reader, knowledge about stereotypical situations. "Scripts" is the label

usually applied to our knowledge about such everyday matters as doing the laundry, going to the doctor, eating at a restaurant. We comprehend a text by finding the pertinent script; then we match the ideas in the text with the prototypical events and roles of the script. The meaning of Bransford and Johnson's passage remains elusive precisely because it evokes no script.

### Bottom-Up Processing

Not everything we read or hear falls neatly into some pre-existing mental cubbyhole (Thoreau, cited in Miller, 1979):

Near the end of March, 1845, I borrowed an axe and went down to the woods by Walden Pond, nearest to where I intended to build my house, and began to cut down some tall, arrowy white pines, still in their youth, for timber . . . . It was a pleasant hillside where I worked, covered with pine woods, through which I looked out on the pond . . . . The ice in the pond was not yet dissolved, though there were some open spaces, and it was all dark-colored and saturated with water.

Miller, taking the bottom-up approach, suggests we interpret this passage by constructing a "memory image" for it (Miller, 1979:204):

My memory image grew piecemeal in roughly the following way. First I read that the time was late March; I formed no image at this point, but filed it away for possible use later. . . . Next, I saw an indistinct Thoreau borrow an axe from an even less distinct somebody and walk, axe in hand, to some woods near a pond.

Most of us have never built a log cabin and have only the sketchiest notion of what it would entail. Despite the absence of anything so well-formed as a script, and despite the discursive, unstorylike nature of the text, we have little difficulty in following the passage. We form some sort of picture of what's going on (Miller's memory image) and the details of the passage fit into that picture. All this is not to deny the importance of prior knowledge or of higher-level cognitive structures in the interpretation process: we may not know much about building log cabins but we understand that Thoreau needs the axe for that purpose. Although Miller does not stress the point, we infer a structure of goals and subgoals: Thoreau borrows the axe to cut down the trees to build his house.

The bottom-up approach (Kintsch and van Dijk, 1978, on "microstructure rules;" Miller, 1979; Miller and Johnson-Laird, 1976) emphasizes the range of the prior knowledge used in interpreting texts. In interpreting Thoreau, we draw upon our knowledge of early spring in New England, of pine trees and axes, of the conditions of 1845, of log cabins, of whatever, in fact, is needed for us to form a coherent mental image for the passage. It is mainly this data-driven character that distinguishes the bottom-up approach. The story grammar and script



theorists tend to focus on a relatively small number of prior structures; when discussing bottom-up processes by contrast, theorists tend to be quite catholic in their tastes, pointing out how the reader uses whatever background knowledge may be relevant to the text at hand.

### **The Importance of Prior Knowledge and Context**

Both approaches to comprehension emphasize our use of prior knowledge--knowledge about the form of texts, knowledge about stereotypical situations, knowledge about concrete details--in interpretation. They also share the notion that context allows us to activate and use the relevant pieces from our vast fund of background information. Without context, we are unable to determine what information is relevant to the passage at hand. With changing contexts, we draw on radically different pieces of information in the comprehension process leading to radically different readings. If Raskolnikov, rather than Thoreau, had borrowed the axe, we would draw rather different conclusions about his purpose.

Prior knowledge is used to connect the ideas in a passage. In Thoreau's account, we saw how different actions are imbedded as goals and subgoals. We fill in the linking connections and the omitted details. Miller's image of Thoreau's passage may be hazy in some respects, but it includes more detail than the passage strictly warrants--he infers a man, for example, perhaps snow on the ground. His reading honors both the claims of the text and the claims of his own knowledge of the world.

### **Implications for Survey Methods**

There are several main themes that both approaches share. We go far beyond the information given in interpreting a text; we fill in gaps and add details, making inferences that our background knowledge and the text at hand seem to call for. The inferences we make depend on the prior knowledge that is activated by the passage and its context; context guides the selection of the relevant prior knowledge.

The work on comprehension has focused on connected text rather than discrete items, on exposition rather than interrogation, on written rather than spoken prose. The relatively little work that has been done on answering questions has also tended to point up the importance of prior structures--such as scripts and story grammars--in guiding the processes by which we seek the information that the question requests (see Bower et al., 1979, and Mandler and Johnson, 1977).

Despite the differences between an interview and a prose passage, there are a few generalizations that probably apply in both situations. First, the inference process can go awry--we may incorrectly infer what is only possible or at best probable. The problem is compounded because we do not sharply distinguish probable inferences (e.g., Thoreau's goal in borrowing an axe) from necessary ones (e.g., that going to the site required the narrator to move). Second, the context of a sentence (or question) will to a large extent determine the nature of the inferences

drawn, the scripts invoked, and the background knowledge brought to bear in interpreting it.

In more practical terms, the literature on comprehension suggests that, other things being equal, related questions that are grouped may have the advantage over questions that are scattered throughout an instrument, because the grouped questions provide a helpful context for interpretation; longer questions may have the advantage over shorter ones, because longer questions create more context; explicit questions may have the advantage over questions that rely on tacit knowledge, because they leave less room for erroneous inferences; questions tailored for different subgroups may have the advantage over more uniformly standardized questions, because different subgroups have different stores of background information that can lead to different interpretations. Each of these generalizations is an oversimplification; they are intended as guidelines rather than hard-and-fast rules. The advantages of long questions, for example, may be offset by the increase in syntactic complexity that usually accompanies increased length.

Survey researchers are hardly unaware of many of these points. Bradburn (1982) describes several methodological studies on context and question order effects. One of the findings apparently illustrates how respondents may draw erroneous inferences about the meaning of a question. In some cases, the answer to a general question (such as “Taken all together, how would you say things are these days--would you say you are very happy, pretty happy, or not too happy?”) may change, depending on whether it comes before or after specific questions (e.g., about one’s marriage). Bradburn interprets this result to mean that when the general question comes last people infer that it covers only what has been left out of the particular questions. Their inference may stem from a general conversational rule that tells us to avoid unnecessary redundancy in our answers (Haviland and Clark, 1974).

Other studies (Belson, 1968; Fee, 1979) show how differences in background can lead to differences in how a question is interpreted. Even terms in common use (such as “energy crisis”) have a wide range of meanings; different kinds of people tend to favor different interpretations of them.

### **Oral Presentation**

Cognitive scientists have concentrated their energies on the comprehension of written materials; it is natural to wonder how far we can apply their conclusions to spoken ones. Much of the work that has been done on the differences between reading and listening has come under the banner of research on attitude change, where the concern has been to compare the effectiveness of oral and written arguments. There is no clear winner. Studies suggest, for example, that oral presentation may be better for simple arguments but worse for complicated ones (Chaiken and Eagly, 1976). Other processes besides comprehension tend to be implicated in these studies--oral presentations may be more effective simply because they are more likely to command our attention, but it may be easier to understand written presentations if we bother to read them.

Mode of presentation may interact with other variables besides complexity. Oral presentation can produce gains in short-term memory performance, but the increases seem to be limited to the last few items in a series. Some phenomena are likelier to arise when material is presented orally--homonyms present little difficulty when read, and channel discrepancies between the verbal and nonverbal messages are unlikely to occur with material presented in writing. On the whole, however, there is little research on how the mode of presentation affects comprehension.

There is a fair amount of evidence that memory for spoken language is generally poor (see, for example, Keenan et al., 1977), although not necessarily worse than memory for written language. (Keenan et al. also suggest that we may remember the irrelevancies, the asides, and the jokes better than the gist of the presentation.) Given the typical length of survey questions, forgetting the question is probably a less common source of respondent error than misunderstanding it.

### Retrieval

Having interpreted--or misinterpreted--the question, the respondent now faces the problem of answering it. Almost all questions require us to rummage about our memories in search of an answer. The fallibility of memory is, of course, more widely appreciated than the fallibility of the comprehension process. This section takes up four issues: What do cognitive scientists say about the structure of memory and the processes used to search for information? What are the sources of forgetting? What can be done about them? When can memory be trusted?

Memory is arguably academic psychology's oldest topic, systematic work having begun more than a century ago with Ebbinghaus. Like all topics, memory has had its ups and downs, but it has never fallen completely out of favor, not even during the behavioristic reign of terror that sent so many mentalistic topics into banishment. It would be pointless to pretend that the vast literature on memory can be summarized in a few broad strokes. It cannot. The best I can hope to do is to highlight a few key ideas that relate most directly to the problems faced by survey researchers.

### Episodic Versus Semantic Memory

Memory researchers distinguish between different stores for information, ranging from the very brief persistence of visual information in iconic storage to the more or less permanent retention of information in long-term memory. Generally, three stores are distinguished: sensory memory (which is thought to be very short lived); short-term memory (which corresponds, in some accounts, roughly to the active contents of consciousness); and long-term memory. Tulving (1972) has introduced a further distinction between semantic and episodic memory. Episodic memory contains our memories for experience, for events that are defined by a particular time and place; semantic memory contains more general

knowledge, such as our knowledge of the meanings of words, that is independent of the setting in which it was learned. Our memory for a particular text generally resides in episodic memory; the background knowledge used in interpreting the text generally resides in semantic memory.

Tulving (1968; Tulving and Thomson, 1973) has argued that the context in which an event is experienced or a word is encountered determines how it is interpreted, stored, and retrieved from episodic memory. His encoding specificity principle states that sometimes we fail even to recognize a previously learned item because at recognition the item fails to reinstate the exact encoding it was given during learning. For example, when one puts eggs on a shopping list, he or she may fail to purchase the right item if candy Easter eggs were intended and hen's eggs were recalled.

Although some psychologists (J. Anderson, 1976) question Tulving's distinction between semantic and episodic memory--arguing that both sorts of memory are retained in a single store and follow the same principles--few question the importance of context in encoding and recall.

### **Memory as an Associative Network**

Our long-term memories are not long lists of all the facts we have learned and the events we have experienced. Memories are connected to each other and the connections can be long and complicated or short and direct. Some theorists (J. Anderson, 1976; Collins and Quillian, 1972) view long-term memory as an associative network. They view individual ideas as nodes or points of intersection in the network; the ideas are connected by links, which correspond to the relations between ideas.

The network analogy is useful for a number of reasons. First, it is compatible with most models of language comprehension; individual sentences or passages of text can be represented using the same network formalism. Second, it allows us to understand how items can be retrieved and forgotten from memory. According to Anderson, for example, we remember by searching relevant portions of the memory network. We forget items that are relatively isolated--few paths in the network lead to the item sought--or whose connections to other items are poorly learned. This perspective implies that good cues for remembering something are those that lead us to search the right part of the network; in most cases, the best cue is the item itself. It comes as no surprise that recognition--recalling whether an item has been presented before--is almost always better than other forms of recall.

### **Retrieval as Reconstruction**

We noted earlier that when we read a passage, we make inferences; we supply connections that are implicit in the text and we add plausible details. There is an abundance of evidence that we encode and store our interpretation of a text rather than the text itself; our memory for gist is far better than our verbatim memory (Bransford et al., 1972; Sachs,

1967). Not only do we lose verbatim details, we also seem to add (and recall) plausible details that fit our interpretation (Bartlett, 1932; Bower et al., 1979; Bransford et al., 1972; Cantor and Mischel, 1977; Mandler and Johnson, 1977). It is not always clear whether we make the inferences when we first encounter the material or when we remember it later. Some theorists argue that we may do much of the “filling in” at the time of recall. What we recall, of course, is likely to be incomplete or spotty; we may round out what we can remember with what we can infer. Once we draw an inference--during our initial encoding or later during an attempt at retrieval--it can be added to our representation of the event. Loftus and her colleagues have conducted several demonstrations on this point. In one (Loftus and Palmer, 1974, Experiment II), the subjects watched a film of a traffic accident. Afterwards, some of them were asked the question, “About how fast were the cars going when they smashed into each other?” The rest were asked a similar question with the word “hit” replacing “smashed into.” The “smashed” subjects were more than twice as likely later on to report incorrectly that they had seen broken glass in the film. The subjects drew such inferences as “smashed into” seemed to warrant; later they could not disentangle what they had seen from what they had inferred.

### Sources of Forgetting

As with research on comprehension, so with research on memory: it is more than a little hazardous to extrapolate from the sort of research conducted by cognitive psychologists to the practical concerns of the survey researcher. The vast bulk of the research on memory has concerned the memory over short periods of time (typically a few minutes) of arbitrary lists generally consisting of nonsense syllables. Although my review has mainly covered research on more meaningful materials, such as text or filmed events, it is still a long way from laboratory studies to everyday memory. Fortunately, Linton's (1982) study of her own everyday memory suggests that the memory problems encountered in the research laboratory are similar to the ones encountered elsewhere.

The research I have attempted to summarize pinpoints several different reasons for forgetting. One reason for forgetting is that we may not have transferred the relevant information into long-term memory. We do not attend to everything and, even when we do, we often attend in a rather mindless fashion. To find out whether one has already purchased the morning newspaper, it may be better to search one's briefcase than one's memory. Another reason that we forget things is that we simply cannot retrieve them. Linton (1982) refers to this type of forgetting as “simple failure to recall.” Some representation of the event made it into long-term memory--in Linton's case, she could often recall such items for months or years--but gradually it becomes inaccessible, lost somewhere in the associative network. Linton describes another type of retrieval failure: “During the fourth and subsequent years [of the study] I began to encounter a few old items that simply did not ‘make sense.’” These items--brief descriptions of events that she had written shortly after the events had taken place--no longer functioned as

interpretable cues for recall. The initial interpretation of the event ceased to have meaning outside of its original context. A number of researchers--beginning with Tulving and his encoding specificity principle--have argued that the value of a recall cue depends on its ability to reinstate the original interpretation (i.e., encoding) of an event. A cue that receives one interpretation in its original context and another in the context of recall will be a poor cue indeed--it leads us to search the wrong part of memory. The converse of this principle is that overlap between features of the original context and the context of recall will facilitate recall. In an intriguing review of the literature, Bower (1981) has shown that memory for lists of words, daily events, and childhood experiences is better when the mood during recall is consistent with the mood during the experience.

Linton's study highlights another source of forgetting--over time, we may lose the ability to distinguish between similar events. It is easy enough to recall one's only trip to St. Michaels, far harder to recall one's trip to the office on March 14. The problem in recalling specific instances of repetitive events is that we forget the particulars and retain the pattern; we can reconstruct the essentials based on the pattern but this strategy offers no basis for recovering the individuating details that would allow us to distinguish one trip to the office from another. In addition, details sometimes seem to wander from one instance of "scripted" event to another (e.g., Bower et al., 1979, Experiment 3.) Details may also be inferred, based on the pattern, even though they were absent from the particular instance (Bower et al., 1979, Experiment 4; cf. Mandler and Johnson, 1977; and Cantor and Mischel, 1977).

In summary, there are several processes of forgetting: (1) The information may never reach long-term memory; (2) it may be irretrievable, particularly when the context of recall differs sharply from the original context of the event; (3) it may be hard to distinguish from related information; (4) it may be tainted with intrusions and inferences made while or after the original information was learned.

### Aids to Memory

It is no news to anyone that memory is fallible. Can anything be done to improve it? A fair amount of research has examined procedures to improve memory that are applied at the initial encoding stage. Although a number of tactics have proved to be effective (such as "deeper" initial processing, the use of mental imagery, and the application of mnemonic tricks), they are, from the point of view of the survey researcher, of not much use--they require the respondent to do something special at the time of encoding and, of course, the researcher is usually not yet on the scene. In longitudinal studies and studies using diaries of course there are more opportunities to train respondents in the use of these mnemonic tactics. Still, tricks that can be applied at the retrieval stage long after the relevant events have taken place are of more general use to survey researchers. Not much is known about this kind of mnemonic trick. Erdelyi and Kleinbard (1976) suggest that when it comes to



retrieval more is better: repeated attempts at recall yield more of the items sought. In addition, the more time we have to recall an item, the more likely we are to retrieve it. For this reason, longer questions, which may have an advantage at the comprehension stage, may also have the edge at the retrieval stage. (On the other hand, these advantages may be dissipated if greater length is accompanied by greater syntactic complexity.) Bower's (1981) work suggests that efforts to reinstate features of the original context will also yield fuller recall.

Left to their own devices, people tend to favor "external" memory aids (such as calendars, diaries, and so on) over the mnemonic tactics usually studied by memory researchers (Harris, 1978). Survey researchers have not ignored the possibilities of using special diaries, logs, and so on for improving recall. Also, some surveys have incorporated existing records as a jog to respondents' memories; a respondent's checkbook may contain the most useful cues to help him or her recall a visit to a doctor. On the whole, people are reasonably aware of the shortcomings of their memories and of how to overcome them (Flavell and Wellman, 1977). Forewarning them about what they will be asked to recall may be the easiest way to improve their memories.

### What Can We Remember?

Not everything is remembered equally well or poorly; it is probably worthwhile to list some kinds of experiences that are likely to be especially well remembered. First of all, a number of researchers have argued that emotional events are unusually well remembered (see, for example, Brown and Kulik, 1977; Holmes, 1970; Sheingold and Tenney, 1982; Colegrove, 1898, is an interesting precursor to the Brown and Kulik work on "flashbulb" memories). Linton (1982) suggests that emotional events are best remembered when they (1) were emotional at the time they occurred, (2) mark transitional points in one's life, and (3) remain relatively unique, to which Bower (1981) might add, (4) they retain their original emotional significance. Research on the reliability of eyewitness testimony, on the other hand, suggests that memory for emotional events may be especially prone to distortion. It is not clear how to reconcile these two lines of research. The key issue may be what happens between the event and the attempt to recall it. Emotional events do not pass unremarked; they are told and retold to friends and relatives, and embellishments that fill in or improve the story may be incorporated into the memory for the event, producing distortions during recall. Owing to the distortions, it may be as hard to remember the details as it is to forget the event.

Other kinds of events have a clearer advantage in memory. Events that stand out against the background of a script or stereotype because they are anomalous in some way tend to be remembered (Bower et al., 1979, Experiment 7; Hastie and Kumar, 1979). In addition, the first few and last few events in a series will be better remembered than events in the middle; drawn-out events will be remembered better than brief ones. We also recall certain aspects of events better than others: "essentials" better than details (this may reflect the ease with which the essentials



can be reconstructed), actions and their outcomes better than their settings or motives, causally connected sequences better than merely temporal ones (Mandler and Johnson, 1977).

### Judgment

Many types of questions require considerable cognitive work from the respondent after the relevant facts have been retrieved. The question may call for a judgment that requires several pieces of information to be combined; it may call for some inference. Hastie's recent review (Hastie, 1983) distinguishes several major approaches to the topic of social inference. After considering two of these approaches, I explore how different kinds of questions will lead to the use of different judgment strategies. I then examine in more detail the process by which attitude questions are answered.

### Information Integration Theory

Although largely a one-man show, information integration theory (N. Anderson, 1974, 1981) has nonetheless had a profound impact in a number of areas within social psychology, ranging from impression formation to equity judgments. The essential idea is simple but powerful: when people must make judgments about something, they combine pieces of information using simple algebraic rules. For example, if we know that someone is a Republican banker, then we judge him to be likeable according to how likeable we find Republicans and how likeable we find bankers. Our judgment about Republican bankers is, according to Anderson, a kind of average of our separate feelings about the two categories. In rendering a judgment--such as a likeability rating--we evaluate the individual pieces of information we have and then we integrate them. The integration is likely to follow some simple rule such as multiplication, addition, or, in this case, averaging. The nuances of the theory involve how we evaluate the individual pieces of information, how we weight them, and when we use one integration rule instead of another. Despite its apparent simplicity, the theory is remarkably adept at explaining a range of phenomena. Jones and Goethals (1972), for example, show that the same information can lead to different judgments depending on the order in which it is presented. (In line with nearly everyone's suspicions, first impressions do seem to carry inordinate weight.) N. Anderson (1974) has argued that this results from waning attention to later items and from the tendency to discount later information when it contradicts what is already given. In the language of integration theory, the weight an item receives depends upon its serial position. Anderson has applied the theory in too many domains to cover in any detail. The key points to bear in mind are that (1) people readily evaluate diverse pieces of information on a common scale, and (2) they seem to combine the information according to simple formulas.

## Judgmental Heuristics

Anderson's work points up the precision with which we can render judgments. The main alternative to information integration theory, the judgmental heuristics approach, presents the contrasting view that we apply loose rules of thumb in rendering many judgments; the rules of thumb ("heuristics") often lead us far astray.

Kahneman and Tversky (1971, 1973; Tversky and Kahneman, 1973) have identified three such heuristics--availability, representativeness, and anchoring and adjustment--and shown how they can lead to systematic errors in judgments of frequency and likelihood.

We use the availability heuristic when we judge the frequency of a class of events based on the number of examples we can generate or the ease with which we generate them. For example, people incorrectly judge that words beginning with the letter R are more frequent than words whose third letter is R--presumably because they can more easily call to mind words that begin with R. Similarly, people may overestimate the likelihood of memorable but rare events because examples are called readily to mind.

Representativeness refers to the tendency for judgments regarding category membership to be rendered solely on the basis of resemblance to some prototype or underlying process. For example, people regard boy-girl-boy-girl-girl-boy as a much likelier sequence of births than three boys followed by three girls--the first sequence seems more representative of the underlying random process. Or they will judge a person described as quiet as likelier to be a librarian than a salesman--even when they are told there are far more salesmen than there are librarians. The representativeness heuristic leads to errors because it fails to incorporate the statistical considerations--such as base rate information--that are also relevant to the judgment. We show little hesitancy about generalizing from a few cases if the cases seem representative of the underlying process or category. Small, biased samples present few problems for nonstatistical inference.

A third judgmental heuristic is anchoring and adjustment. The anchor refers to some starting point or preliminary estimate; this initial judgment is then adjusted to produce the final verdict. Problems arise when the adjustment is insufficient or the anchor misleading. For example, when asked to calculate  $10 \times 9 \times 8 \times \dots \times 2 \times 1$  in their heads, people appear to multiply the first few numbers and, using this as their anchor, adjust the result upward. They come to a consistently higher conclusion than people asked to multiply  $1 \times 2 \times 3 \dots \times 9 \times 10$ , who presumably begin their adjustment from a lower anchor point. Both groups tend to underestimate the actual answer, suggesting that the adjustment is insufficient in either case. A more dramatic example of the inappropriate use of anchoring and adjustment is provided by Ross, Lepper, and Hubbard (1975). Subjects in the study were misled about their performance in distinguishing authentic suicide notes from fake ones. Some subjects were told they had performed much above (or much below) average; later on, the experimenter admitted that the evaluation had been a complete hoax. Despite the discrediting of the evidence on which their opinions were based, "successful" subjects continued to

believe they were better than average at the task and “failure” subjects continued to believe they were worse than average. The subjects apparently underadjust their opinions in the light of the discrediting of the feedback.

### Types of Questions

In the dozen years since Kahneman and Tversky burst upon the psychological scene, the number of heuristics and biases that we have been shown to exhibit has grown at something like an exponential rate. Studies indicate that we persevere in the face of discrediting evidence (Ross, Lepper, and Hubbard, 1975); ignore base-rate information (Kahneman and Tversky, 1971; Nisbett and Ross, 1980, Chapter 7); generalize from small or biased samples (Tversky and Kahneman, 1971); make extreme predictions based on invalid indicators (Kahneman and Tversky, 1973); perceive “illusory” correlations that sustain our prejudices (Chapman and Chapman, 1969; Hamilton and Gifford, 1976); search for evidence in ways that can only confirm our views (Snyder and Swann, 1978); believe in our control over chance events (Langer, 1975); attribute causal significance to whatever happens to be salient during an event (Taylor and Fiske, 1975; McArthur and Post, 1977); ascribe all sorts of personal characteristics to people who are clearly doing only what the situation requires (cf. Ross, 1977, on the “fundamental attribution error”); “recognize” new pieces of information that confirm our stereotypes (Cantor and Mischel, 1977)--and, in spite of all these failings, we remain confident that our judgments are sound (Einhorn and Hogarth, 1978). The list of our shortcomings could be easily be extended; it can be as long as the reviewer is persistent. Still, it is hard to believe that our judgment is as bad as all this might suggest; we must occasionally get things right. Something of a backlash has already set in (see, for example, Loftus and Beach, 1982); in their comprehensive review, Nisbett and Ross (1980) take pains to argue that, despite the thousand natural errors that the mind is heir to, we often do make valid judgments and sound decisions.

### Implications for Survey Research

It would be nice for everyone if we could reduce this long list of biases to one or two fundamental errors from which the rest could be derived, but no one has thus far succeeded in imposing order on this error-filled crew. Failing that, it would be useful if we could identify the particular judgment tasks that typically elicit particular strategies with their characteristic shortcomings, but no one has succeeded even in that endeavor. Nisbett and Ross (1980) do manage to relate different errors to different steps in the judgment process (e.g., gathering evidence, summarizing it, etc.) but particular judgment tasks differ greatly in which steps they require.

Sykes (1982) distinguishes four types of questions commonly used in surveys: those that request factual or behavioral information; those

that assess the respondent's awareness or knowledge of a subject; those that elicit attitudes or opinions; and those that call for a reason or explanation. Questions about what we have done and what we know are probably the least susceptible to judgmental error--although judgments of the frequency or recency of a behavior are known to be distorted in predictable ways. (Most of the Kahneman and Tversky work bears on the issue of frequency judgments.) Questions about the causes or reasons for behavior, on the other hand, require considerable judgment on the part of the respondent, and the judgment process seems to be very flawed (Ross, 1977); what is more, people often seem to use the same process whether they are explaining their own behavior or someone else's (Bem, 1967, 1972; Nisbett and Wilson, 1977)--many of our explanations appear to be based on commonsensical notions of why people behave the way they do rather than on access to some fund of private knowledge. Nisbett and Wilson (1977) argue that if we do know ourselves better than others know us, it is largely in the historical sense that we remember how we have behaved in the past.

### Attitude Questions

Attitude questions are an interesting case because very little is known about how we answer them. The process is likely to be quite complex, involving both the use of judgmental heuristics and the application of integration rules. Attitudes have been shown to be sensitive to a range of influences: beliefs (Ajzen and Fishbein, 1972; Rosenberg, 1956); values (Rokeach, 1971); norms (Ajzen and Fishbein, 1972); feelings (Abelson, Kinder, Peters, and Fiske, 1982); other attitudes (Abelson et al., 1968); behaviors (Freedman, 1965); and arguments. The exact nature of the relationship between attitudes, on the one hand, and beliefs, values, norms, emotions, and behaviors, on the other, is often hard to disentangle; behaviors, for example, are clearly affected by attitudes, but the reverse is also true, and both attitudes and behaviors are affected by still other things, such as norms.

It is tempting to conclude that we must answer attitude questions by retrieving the relevant beliefs, values, behaviors, and so on from memory and then making a judgment based on whatever is brought to mind. Unfortunately, we know that many of the things that affect attitudes--such as persuasive arguments--are quickly forgotten, though their impact lingers on (Hass and Linder, 1972). (A more familiar version of the phenomenon is knowing that we have heard convincing evidence on some issue without being able to recall what the evidence is.)

If we do not retrieve evidence in answering attitude questions, then what do we retrieve? There are probably different answers to this question depending on the nature of the attitudes we hold. Let us first distinguish well-formed attitudes from "nonattitudes" (Converse, 1963). There are many topics on which we have no particular opinion--we may have some beliefs but they will be weakly held and poorly supported; we may have some feelings but they will be inconsistent and changeable.

With nonattitudes, we are likely to retrieve the one or two salient pieces of information that we have--perhaps a cliché we have heard about

the topic or some example of our behavior. Having sampled our memories in this haphazard way, we “scale” what we have retrieved; we ask ourselves “what are the evaluative implications of this cliché, that behavior?” We then combine the evaluations assigned to each piece of information--perhaps we average them or perhaps we use whatever comes to mind first as an anchor and then adjust our judgment as we recall and evaluate additional bits of information. With nonattitudes, then, answering an attitude question has three stages: we sample our views, we scale them, and then we combine them. Since the views we hold are so poorly interconnected, what we retrieve on one occasion may bear little resemblance to what we retrieve on the next. We must expect nonattitudes to be quite unreliable.

There are also topics about which we have passionate feelings or deeply held beliefs. It is probably useful to distinguish several categories of such well-formed attitudes: (1) attitudes that consist of well-developed stereotypes; (2) attitudes that relate to deeply felt, symbolic concerns; (3) attitudes that relate to immediate, practical issues.

Research on stereotypes is hardly new but it has taken on new impetus recently because of the interest in the internal structure of categories (Rosch and Lloyd, 1978; see also the work on schemata by Markus, 1977; Rumelhart and Ortony, 1977). According to the schema theorists (Rumelhart and Ortony, 1977; Ortony, 1979), our conception of a category may consist of a kind of frame (the schema) with a series of “slots” for the details that distinguish one category member from another. The slots may contain “default” values, reflecting our sense of what the prototypical category member is like. With socially significant categories, one of the slots may well contain our evaluation of the instance or a default evaluation that reflects our judgment of the category as a whole. Thus, it is plausible to suppose that when we ask a Marxist how well he or she likes bankers, the question activates a dogeared stereotype or schema that includes an evaluative slot with the information that bankers are capitalists of the worst kind. Similarly, we may store evaluative reactions alongside less emotionally charged information in our memory representation for individual category members such as particular political figures (Abelson et al., 1982).

Schemata do not come into play only with familiar social groups. Attitudes about abstract political issues often take on a schematic cast. For example, research on attitudes towards school busing (Kinder and Sears, 1981) suggests that attitudes on this emotional question are less a function of practical considerations (such as whether one's own children are likely to be bused) than of symbolic ones. People's attitudes on the busing question are most easily predicted from their attitudes about apparently unrelated issues, such as “welfare.” It is as if people saw these issues as manifestations of some larger pattern; in the case of busing, opponents seem to see it as an example of a general schema in which liberal “reforms” pose a threat to traditional American values. In the case of symbolic attitudes, then, what gets retrieved in answering an attitude question is a kind of schema that includes deeply felt emotions.

With some issues, practical considerations do seem to have the upper hand. Surely some of the people who favored the recent tax cuts favored them because they had a good deal to gain from them. People who hold such practical attitudes are probably more open to attitude change than people holding symbolic attitudes or attitudes based on stereotypes. This suggests that questions regarding practical attitudes may be answered through a process quite different from the one used when stereotypes or symbolic issues are involved. Instead of retrieving some schema with its ready-made evaluation, we may retrieve the relevant evidence (beliefs one holds, arguments one has heard), evaluate it, and combine it to render our judgment. This process is quite similar to the one proposed for nonattitudes, only the memories activated are far more detailed.

These four types of attitudes are meant as pure cases. Many attitudes doubtless include elements of both symbolic significance and practical calculation; one man's cliched nonattitude is another man's stereotype. Still, it is worth noting that with some attitude questions we must compute an evaluative judgment; with others, we simply retrieve one. The form of the attitude question may also influence whether we employ the one process or the other. Likert-type items, which require us to rate how much we agree or disagree with some attitude statement, may encourage the retrieval and integration of more detailed information from memory than items calling for a simple evaluation. The best way to assess attitudes may depend therefore on the type of attitude being assessed.

### Response

The respondent's task is not quite finished. Having rendered the judgment that the question demands, he or she must now select (or formulate) a response. This section deals exclusively with the selection of a response from a pre-established set of response categories. Open-ended questions require respondents to formulate their answers, a process too complicated to cover here.

This section deals with two major issues--how respondents select their answer and when they misreport it.

### Response Selection

To select a response is to make a choice. Psychologists have studied a wide variety of decision rules to cover a range of choice situations. One of the most intensively studied rules--the Luce choice rule (Luce, 1959)--underscores the fact that making choices is a chancy business; confronted with an identical set of alternatives, we do not always select the same one. According to the Luce choice rule, this may occur not because we evaluate the options any differently but because the decision process is inherently probabilistic. The Luce choice rule says that we assign a value to each option; the probability of selecting a particular option is the proportion its value represents relative to the total value



assigned to all of the options. Although the Luce choice rule has been successfully applied in a variety of settings (e.g., it has been proposed as the method by which we select the answers to analogy problems, Rumelhart and Abrahamson, 1973), it is not an intuitively appealing model for response selection in surveys. A more intuitive procedure has been suggested by Tversky (1972), who argues that we make choices by eliminating the options that lack some desirable feature or aspect. Like the Luce choice rule, Tversky's model incorporates a chance element: our final decision depends on which aspects we select as a basis for narrowing down the options, and this selection process is partly a matter of chance. Some decision models have been proposed specifically for memory tasks such as deciding whether one recognizes an item. These models assume that whether one accepts an item as "old" depends not only on its level of activation (its feeling of familiarity) but also on one's criterion for accepting an item as old. The criterion will be influenced by the relative probabilities and costs of the two types of possible errors (failing to recognize an old item, falsely recognizing a new one). (See Bower et al., 1979, for another recognition decision model based on the level of activation.)

The Luce rule and Tversky's elimination by aspects model are useful in pointing up the chance component in the choice process. The models of choice applied to memory are useful in underscoring another point--the role of nonfactual considerations in response selection. Respondents with relatively low "criteria" may report events they recall in only the vaguest of terms because they feel that it is better to give some response, however inaccurate, than none. Similarly, respondents may report an attitude simply to avoid the minimal embarrassment of admitting that they do not have one--it is probably not difficult to produce an attitude judgment on demand. Responses to even the most factual of items are influenced by considerations besides the facts.

Many models of choice assume that the respondent has some ideal in mind and then selects the option that is the closest approximation to it. Attitude scaling techniques--such as Thurstone scaling and Coombs's unfolding technique--presuppose the existence of a unitary, underlying attitude dimension on which the respondent's ideal point and the various attitude items can be positioned. These models alert us to the problem of "nonscale" respondents, for whom the dimension is not meaningful, and of extreme ones, who lie outside the range of the options provided.

In many situations, we are willing to settle for something that is considerably less than ideal; often we demand no more than what will suffice. The widespread use of "satisficing" (Simon and Stedry, 1969) rules creates some headaches for the survey researcher. Respondents, like the rest of us, probably do not always wait patiently until all the options have been laid out; they may leap at the first option that seems satisfactory, ignoring the rest. A number of techniques can be used to reduce this tendency--placing normatively less desirable options first in a series, shifting the order around to force the respondent to attend to all the options, announcing the options ahead of time.



## Reporting Biases

Although there is a fair amount of evidence that we describe ourselves in more glowing terms than others do, it is not clear why this is so. We may believe our favorable self-descriptions, or we may merely hope that others believe them. The process may be conscious and deliberate or unconscious and unintended. Besides these “self-serving” biases (Nisbett and Ross, 1980), which are as well-known to survey researchers as to cognitive scientists, there are probably a number of subtler reporting biases. For example, people tend to exhibit consistency among their attitudes (Abelson et al., 1968) and between their attitudes and behaviors (e.g., Salancik and Conway, 1975); pressure toward consistency seems to be induced merely by asking questions (McGuire, 1960). While consistency pressures can result in long-lasting attitude change (Freedman, 1965), they can also produce short-lived and, from the point of view of the survey researcher, artifactual changes.

We strive to present ourselves in a favorable light, and we strive to be consistent. We are prey to other sources of distortion in our responses: we select socially desirable responses (Crowne and Marlowe, 1964), we fulfill the investigator's expectations (Rosenthal and Jacobson, 1966), and we meet the perceived demands of the situation (Orne, 1972). These sources of response error are no news to survey researchers, who probably know more about reducing their effects than cognitive scientists. I shall not dwell on them here.

### Survey Research as Cognitive Laboratory

All of the foregoing has been an attempt to show that the cognitive sciences have something to offer survey methodology. It can point out sources of potential errors and methods for reducing them. This section takes the opposite tack and examines what survey methods have to offer cognitive scientists.

It is perhaps somewhat surprising that survey research is not already among the mansions in the house of cognitive science. Certainly, the problems that cognitive psychologists study can often be studied as readily in the field with survey samples as in the laboratory with college students. Sometimes there is a remarkable convergence, as when Ross et al. (1978) and Fields and Schuman (1976) nearly simultaneously discovered the same phenomenon--our tendency to see our own opinion as relatively common--using laboratory and survey methods.

Investigators are rarely trained in both sets of methods, and they tend to frame their questions in ways that lend themselves to investigation by the methods that they know the best. Researchers are no less susceptible to the unsystematic sampling of alternatives than anyone else, so it is really no surprise that they should tailor their research to the familiar methods that spring to mind.

There are costs to letting the tail wag the dog in this way. In the case of cognitive psychology, critics, such as Neisser (1982), have charged that the research tradition consists primarily of studies involving impractical problems performed in unnatural settings by

unrepresentative samples. In reviewing the literature on comprehension, I noted that the research centers almost exclusively on connected expository text, which makes it difficult to draw conclusions about discrete spoken questions. The same sort of observations can be made about research on memory; the bulk of work concerns memory for arbitrary lists, often composed of nonsense syllables. Neisser (1976, 1978) argues that much of this work lacks “ecological” validity, since the tasks and settings are far removed from anything that exists outside of the laboratory. Of course, in some sense, the whole point of laboratory research is to create artificial situations where the effects of variables can be isolated from each other. It is nonetheless hardly illegitimate to wonder whether laboratory tasks really capture the essence of the real-world problems they are designed to model. The most convincing answer to the critics is to demonstrate that the principles discovered in the laboratory do in fact apply in other settings and with other populations.

Some problems do not easily lend themselves to laboratory methods: it is difficult to study disaster, passionate love, or schizophrenia by attempting to simulate them in the controlled setting of the laboratory. The best that the laboratory researcher can do is to create some weak or “acute” version of the the phenomenon of interest--and to hope that it remains the same phenomenon. The researcher's conflicting aims are to reproduce the phenomena closely enough to be convincing but not so closely as to be unethical. Large survey samples can give access to phenonema and populations that are outside the range ordinarily available to the cognitive psychologist. Two phenomena--forgetting and optimism--illustrate the advantages of adding a little more variety to the techniques of cognitive science.

### **Forgetting**

The issue with forgetting is how far the principles that apply to forgetting in the laboratory generalize to other settings and other types of material. Loftus (1982) lists a number of laboratory-based generalizations that have been tested in survey settings: (1) we forget more as time passes, with a higher rate of forgetting shortly after we have learned the material; (2) our memories get worse as we grow older; (3) we are more likely to forget items in the middle of a sequence than ones at the beginning or the end; (4) we are less likely to forget something the longer we are exposed to it; and (5) we are more likely to recognize something than to recall it unaided. Loftus examines each of these generalizations in the light of results from surveys, with mixed results. The classical negatively accelerated forgetting curve, for example, does not always seem to hold; instead, people sometimes forget at a constant rate (see also Linton, 1982; Sheinwold and Tenney, 1982). On the other hand, the relationship between exposure time and recall appears to be robust. It is unclear why personal experiences seem to follow a different forgetting curve from other memories; it raises the intriguing possibility that different types of material may be forgotten through different processes.

It can be difficult to measure forgetting in a survey setting, but it is hardly impossible. Sometimes records are available as a means of checking the respondents' recollections. Other times, it is only possible to assess overall levels of forgetting by comparing two groups of respondents given similar recall tasks; the group which reports more has presumably remembered more.

### Optimism

With memory research, the main problem is that the range of the stimuli is restricted. With research in other areas, the problem is that the range of the subjects is restricted. Laboratory-based research indicates, for example, that people hold a number of optimistic beliefs: the desirability of an event is seen as being related to its probability; bad events are seen as likelier to happen to others than to ourselves (Weinstein, 1978). Investigators have suggested a number of plausible explanations for this optimistic tendency, ranging from the idea that hope, even false hope, helps to ward off impassivity to the notion that happy outcomes are better attended to and consequently are more likely to be recalled.

Since most of the research on optimism has been conducted with college students, the question naturally arises whether these optimistic beliefs are not, in some sense, realistic--after all, college students in the United States have a fairly cushy existence, sheltered from much of life's storm and stress. Aren't they right to believe that their futures are rosy? A related hypothesis is that optimistic beliefs are the product of the optimistic ideologies that prevail in American society. The Christian belief in personal salvation, the liberal belief in technological progress, the conservative belief in individual efficacy--these form our ideological heritage. It is difficult to get a heterogenous enough sample of respondents to test these ideas unless one has access to a large survey sample. Are optimistic beliefs the residue of generally positive life experiences? We need to talk to the poor, the sick, the lonely to find out whether they, too, believe in their relative invulnerability to life's unpleasantness. Are individual optimistic beliefs just an outcropping of an underlying cultural optimism? We need to talk to new arrivals, to members of cultural outgroups, to the unassimilated to test this hypothesis. And the non-optimists will be of interest for either account. Do the depressed, the anxious, the paranoid show an unpleasant history that justifies their pessimism about the future? Were they somehow immune to the ambient ideological optimism? It is impossible to explore these questions satisfactorily without the sheer variety of people that only a large survey sample can offer.

## References

- Abelson, R. 1981 Psychological status of the script concept. *American Psychologist* 36:715-729.
- Abelson, R., Aronson, E., McGuire, W., Newcomb, T., Rosenberg, M., and Tannenbaum, P. 1968 *Theories of Cognitive Consistency: A Sourcebook*. Chicago: Rand McNally.
- Abelson, R., Kinder, D., Peters, M., and Fiske, S. 1982 Affective and semantic components in political person perception. *Journal of Personality and Social Psychology* 42:619-630.
- Ajzen, I., and Fishbein, M. 1972 Attitudes and normative beliefs as factors influencing behavioral intentions. *Journal of Personality and Social Psychology* 21:1-9.
- Anderson, J. 1976 *Language, Memory, and Thought*. Hillsdale, N.J.: Lawrence Erlbaum Associates.
- Anderson, N. 1974 Cognitive algebra. In L. Berkowitz, ed., *Advances in Experimental Social Psychology*, Vol. 7. New York: Academic Press.
- 1981 *Foundations of Information Integration Theory*. New York: Academic Press.
- Bartlett, F. 1932 *Remembering*. Cambridge, England: Cambridge University Press.
- Belson, W. 1968 Respondent understanding of survey questions. *Polls* 3:1-13.
- Bem, D. 1967 Self-perception: an alternative explanation of cognitive dissonance phenomena. *Psychological Review* 74:183-200.
- 1972 Self-perception theory. In L. Berkowitz, ed., *Advances in Experimental Social Psychology*, Vol. 6. New York: Academic Press.
- Bower, G. 1981 Mood and memory. *American Psychologist* 36:129-148.
- Bower, G., Black, J., and Turner, T. 1979 Scripts in text comprehension and memory. *Cognitive Psychology* 11:177-220.
- Bradburn, N. 1982 Question-wording effects in surveys. In R. Hogarth, ed., *Question Framing and Response Consistency*. San Francisco: Jossey-Bass.
- Bransford, J., Barclay, J., and Franks, J. 1972 Sentence memory: a constructive vs. interpretive approach. *Cognitive Psychology* 3:193-209.

- Bransford, J., and Johnson, M. 1972 Contextual prerequisites for understanding: some investigations of comprehension and recall. Journal of Verbal Learning and Verbal Behavior 11:717-726.
- Brown, R., and Kulik, J. 1977 Flashbulb memories. Cognition 5:73-99.
- Cantor, N., and Mischel, W. 1977 Traits as prototypes: effects on recognition memory. Journal of Personality and Social Psychology 35:38-49.
- Chaiken, S., and Eagly, A. 1976 Communication modality as a determinant of message persuasiveness and message comprehensibility. Journal of Personality and Social Psychology 34:605-614.
- Chapman, L., and Chapman, J. 1969 Illusory correlation as an obstacle to the use of valid psychodiagnostic signs. Journal of Abnormal Psychology 74:271-280.
- Colegrove, T. 1898 Individual memories. American Journal of Psychology 10:228-255.
- Collins A., and Quillian, M. 1972 Experiments on semantic memory and language comprehension. In L. Gregg, ed., Cognition and Learning. New York: John Wiley and Sons.
- Converse, P. 1963 Attitudes and Non-Attitudes: Continuation of a Dialogue. Paper presented at meeting of the International Congress on Psychology, Washington, D.C.
- Crowne, D., and Marlowe, D. 1964 The Approval Motive: Studies in Evaluative Dependence. New York: Wiley.
- Einhorn, H., and Hogarth, R. 1978 Confidence in judgment: Persistence in the illusion of validity. Psychological Review 85:395-416.
- Erdelyi, M., and Kleinbard, J. 1976 Has Ebbinghaus decayed over time? The growth of recall (hypernesia) over days. Journal of Experimental Psychology: Human Learning and Memory 4:275-278.
- Fee, J. 1979 Symbols and Attitudes: How People Think About Politics. Unpublished doctoral dissertation, University of Chicago.
- Fields, J., and Schuman, H. 1976 Public beliefs about the beliefs of the public. Public Opinion Quarterly 40:427-448.
- Fischhoff, B. 1975 Hindsight = foresight: the effect of outcome knowledge on judgment under uncertainty. Journal of Experimental Psychology: Human Perception and Performance 1:288-299.

- Flavell, J., and Wellman, H. 1977 Metamemory. In R. Kail and J. Hagen, eds., Perspectives on the Development of Memory and Cognition. Hillsdale, N.J.: Lawrence Erlbaum Associates.
- Freedman, J. 1965 Long-term behavioral effects of cognitive dissonance. Journal of Experimental Social Psychology 1:145-155.
- Hamilton, D., and Gifford, R. 1976 Illusory correlation in interpersonal perception: a cognitive basis of stereotypic judgments. Journal of Experimental Social Psychology 12:392-407.
- Harris, J. 1978 External memory aids. In M. Gruneberg, P. Morris, and R. Sykes, eds., Practical Aspects of Memory. London: Academic Press.
- Hass, R., and Linder, D. 1972 Counterargument availability and the effect of message structure on persuasion. Journal of Personality and Social Psychology 23:219-233.
- Hastie, R. 1983 Social inference. Annual Review of Psychology 34:511-542.
- Hastie, R., and Kumar, P. 1979 Person memory: personality traits as organizing principles in memory for behavior. Journal of Personality and Social Psychology 37:25-38.
- Haviland, S., and Clark, H. 1974 What's new? Acquiring new information as a process in comprehension. Journal of Verbal Learning and Verbal Behavior 13:512-521.
- Holmes, D. 1970 Differential change in affective intensity and the forgetting of unpleasant personal experiences. Journal of Personality and Social Psychology 15:235-239.
- Jones, E., and Goethals, G. 1972 Order effects in impression formation: attribution context and the nature of the entity. In E. Jones, R. Kanouse, H. Kelley, R. Nisbett, S. Valins, and B. Wiener, eds., Attribution: Perceiving the Causes of Behavior. Morristown, N.J.: General Learning Press.
- Kahneman, D., and Tversky, A. 1971 Subjective probability: a judgment of representativeness. Cognitive Psychology 3:430-454.
- 1973 On the psychology of prediction. Psychological Review 80:237-251.
- Keenan, J., MacWhinney, B., and Mayhew, D. 1977 Pragmatics in memory: a study of natural conversation. Journal of Verbal Learning and Verbal Behavior 16:549-560.
- Kinder, D., and Sears, D. 1981 Prejudice and politics: symbolic racism versus racial threats to "the good life." Journal of Personality and Social Psychology 40:414-431.

- Kintsch, W., and van Dijk, T. 1978 Toward a model of text comprehension and production. *Psychological Review* 85:363-394.
- Langer, E. 1975 The illusion of control. *Journal of Personality and Social Psychology* 32:311-328.
- Linton, M. 1982 Transformations of memory in everyday life. In U. Neisser, ed., *Memory Observed*. San Francisco: U.H. Freeman and Company.
- Loftus, E. 1982 Memory and its distortions. Pp. 119-154 in A.G. Kraut, ed., *G. Stanley Hall Lectures*. Washington, D.C.: American Psychological Association.
- Loftus, E., and Beach, L. 1982 Human inference and judgment: is the glass half empty or half full? *Stanford Law Review* 34:939-956.
- Loftus, E., and Marburger, W. 1983 Since the eruption of Mt. St. Helens, did anyone beat you up? Improving the accuracy of retrospective reports with landmark events. *Memory and Cognition* 11:114-120.
- Loftus, E., and Palmer, J. 1974 Reconstruction of automobile destruction: an example of the interaction between language and memory. *Journal of Verbal Learning and Verbal Behavior* 13:585-589.
- Luce, R. 1959 *Individual Choice Behavior*. New York: John Wiley and Sons.
- Mandler, J., and Johnson, N. 1977 Remembrance of things passed: story structure and recall. *Cognitive Psychology* 9:111-151.
- Markus, H. 1977 Self-schemata and processing information about the self. *Journal of Personality and Social Psychology* 35:63-78.
- McArthur, L., and Post, D. 1977 Figural emphasis and person perception. *Journal of Experimental Social Psychology* 13:520-535.
- McGuire, W. 1960 A syllogistic analysis of cognitive relationships. In M. Rosenberg, C. Hovland, W. McGuire, R. Abelson, and J. Brehm, eds., *Attitude Organization and Change*. New Haven: Yale University Press.
- Miller, G. 1979 Images and models, similes and metaphors. In A. Ortony, ed., *Metaphor and Thought*. Cambridge, England: Cambridge University Press.
- Miller, G., and Johnson-Laird, P. 1976 *Language and Perception*. Cambridge, Mass.: Harvard University Press.



- Neisser, U. 1976 Cognition and Reality: Principles and Implications of Cognitive Psychology. San Francisco: W.H. Freeman and Company.
- 1978 Memory: what are the important questions? In M. Gruneberg, P. Morris, and R. Sykes, eds., Practical Aspects of Memory. London: Academic Press.
- 1982 Memory: what are the important questions? In U. Neisser, ed., Memory Observed. San Francisco: W.H. Freeman and Company.
- Nisbett, R., and Ross, L. 1980 Human Inference: Strategies and Shortcomings of Social Judgment. Englewood Cliffs, N.J.: Prentice-Hall.
- Nisbett, R., and Wilson, T. 1977 Telling more than we can know: verbal reports on mental process. Psychological Review 84:231-259.
- Orne, M. 1962 On the social psychology of the psychological experiment: with particular reference to demand characteristics and their implications. American Psychologist 17:776-783.
- Ortony, A. 1979 Beyond literal similarity. Psychological Review 86:161-180.
- Rokeach, M. 1971 Long-range experimental modification of values, attitudes, and behavior. American Psychologist 26:453-459.
- Rosch, E., and Lloyd, B. 1978 Cognition and Categorization. Hillsdale, N.J.: Lawrence Erlbaum Associates.
- Rosenberg, M. 1956 Cognitive structure and attitudinal affect. Journal of Abnormal Social Psychology 53:367-372.
- Rosenthal, R., and Jacobson, L. 1966 Teachers' expectancies: determinants of pupils' I.Q. gains. Psychological Reports 19:115-118.
- Ross, L. 1977 The intuitive psychologist and his shortcomings. In L. Berkowitz, ed., Advances in Experimental Social Psychology, Vol. 10. New York: Academic Press.
- Ross, L., Greene, D., and House, P. 1978 The false consensus phenomenon: an attributional bias in self-perception and social perception process. Journal of Experimental Social Psychology 13:279-301.
- Ross, L., Lepper, M., and Hubbard, M. 1975 Perservance in self-perception and social perception. Journal of Personality and Social Psychology 32:880-892.
- Rumelhart, D. 1975 Notes on a schema for stories. In D. Bobrow and A. Collins, eds., Representation and Understanding: Studies in Cognitive Science. New York: Academic Press.

- Rumelhart, D., and Abrahamson, A. 1973 A model for analogical reasoning. *Cognitive Psychology* 5:1-28.
- Rumelhart, D., and Ortony, A. 1977 The representation of knowledge in memory. In R. Anderson, R. Spiro, and W. Montague, eds., *Schooling and the Acquisition of Knowledge*. Hillsdale, N.J.: Lawrence Erlbaum Associates.
- Sachs, J. 1967 Recognition memory for syntactic and semantic aspects of connected discourse. *Perception and Psychophysics* 2:437-442.
- Salancik, G., and Conway, C. 1975 Attitude inferences from salient and relevant cognitive content about behavior. *Journal of Personality and Social Psychology* 32:829-840.
- Schank, R., and Abelson, R. 1977 *Scripts, Plans, Goals and Understanding*. Hillsdale, N.J.: Lawrence Erlbaum Associates.
- Sheinwold, K., and Tenney, Y. 1982 Memory for a salient childhood event. In U. Neisser, ed., *Memory Observed*. San Francisco: W.H. Freeman and Company.
- Simon, H., and Stedry, A. 1969 Psychology and economics. In G. Lindzey and E. Aronson, eds., *The Handbook of Social Psychology*, Vol. 5. Reading, Mass.: Addison-Wesley.
- Snyder, M., and Swann, W. 1978 Behavioral confirmation in social interaction: from social perception to social reality. *Journal of Experimental Social Psychology* 14:148-162.
- Sykes, W. 1982 Investigation of the effects of question form. *Survey Methods Newsletter* 9-10.
- Taylor, S., and Crocker, J. 1980 Schematic bases of social information processing. In E. Higgins, P. Herman, and M. Zanna, eds., *The Ontario Symposium on Social Cognition*, Vol. 1. Hillsdale, N.J.: Lawrence Erlbaum Associates.
- Taylor, S., and Fiske, S. 1975 Point of view and perceptions of causality. *Journal of Personality and Social Psychology* 32:439-445.
- Tulving, E. 1968 When is recall higher than recognition? *Psychonomic Science* 10:53-54.
- 1972 Episodic and semantic memory. In E. Tulving and W. Donaldson, eds., *Organization of Memory*. New York: Academic Press.
- Tulving, E., and Thomson, D.M. 1973 Encoding specificity and retrieval processes in episodic memory. *Psychological Review* 80:352-373.

- Tversky, A. 1972 Elimination by aspects: a theory of choice. Psychological Review 79:281-299.
- Tversky, A., and Kahneman, D. 1971 Belief in the law of small numbers. Psychological Bulletin 76:105-110.
- 1973 Availability: a heuristic for judging frequency and probability. Cognitive Psychology 5:207-232.
- Weinstein, N. 1978 Unrealistic optimism about future life events. Journal of Personality and Social Psychology 39:806-820. Chicago: Rand McNally, 1968.

## POTENTIAL CONTRIBUTIONS OF COGNITIVE RESEARCH TO SURVEY QUESTIONNAIRE DESIGN

Norman Bradburn and Catalina Danis

The purpose of this paper is to review what research on cognitive processes can contribute to the understanding of errors in survey questioning. The paper is not meant to be a comprehensive review of the literature, but rather an attempt to discuss some problems that are of concern to survey researchers in the light of research results from the relevant cognitive literature. If the paper is successful, it should further a dialogue between two research traditions that potentially have much to say to one another but so far have not addressed each other. We start by presenting a model of response effects, a model of human information processing, and a discussion of the differences between research traditions in the two fields to indicate some of the difficulties in trying to bring them together. We then proceed to review some of the major response effects discussed in the survey research literature in the light of what appears to be the most relevant cognitive literature.

### Conceptual Model of Response Effects\*

Our model of the survey data collection process conceives of the research interview as a micro-social system. In an idealized form, the system consists of two roles linked by the task of transmitting information from the respondent to the interviewer (and ultimately to the investigator). We distinguish three sources of variation in the quality of the data: that stemming from the characteristics of the task itself, that from the interviewer's performance, and that from the respondent. Much of the research on response effects has focused on interviewer and respondent characteristics: for example, the race of the interviewer or the propensity of the respondents to agree to statements without regard to their content. This concentration of effort is probably misplaced because it is the task itself that gives rise to what Orne (1969) has called the "demand characteristics of the situation." The demand characteristics, in turn, play the predominant role in determining the behavior of the actors in the situation. Thus variables affecting the characteristics of the task are at the heart of a model of response effects. Indeed, the empirical literature suggests that the characteristics of the task are the major source of response effects and are, in general, much larger than effects due to interviewer or respondent characteristics.

The task in surveys of human populations is to obtain information from a sample of respondents about their (or someone else's) behavior, knowledge, or attitudes. The respondent's role is to provide that

---

\*The discussion in this section is drawn from Bradburn (1983).

information; the interviewer's, to obtain the information in the manner prescribed by the investigator (who defines the task by drawing the sample, designing the questionnaire, and specifying the observations to be employed in the research). If respondents are to be "good" respondents, they must provide accurate and complete information. Careful attention must be given to motivating respondents to play such a role and to defining the situation for them so that they know accurately what it is that they are supposed to do. Similarly, through training, supervision, and careful specification of the questionnaire and its mode of administration, the investigator sets the standards by which interviewers will be judged on how well they have performed their role.

Within this general framework, we can see that there are three sources of response effects. The first source is the respondents themselves. While we expect that most of the variance in responses among respondents comes from real differences, it is possible that there are individual differences among respondents that systematically affect their willingness or ability to give accurate responses, particularly to certain kinds of questions, such as those that might affect their self-esteem. In addition, other factors, such as the presence of other people during the interview, events that happened to the respondent before the interview began, or social pressures not to cooperate with strangers, may undermine the willingness of respondents to take the time or make the effort to be "good" respondents.

The interviewer's role may be more or less prescribed. In some surveys, interviewers are given considerable freedom in defining the task for themselves and for respondents, particularly with regard to the formulation of questions or the follow-up of answers to previous questions. Today, however, most large-scale surveys use heavily structured questionnaires that leave little room for independent judgment about what questions to ask, what order to ask them in, or what to do when respondents answer one way rather than another. Interviewers, of course, do not always do what they are supposed to do, and it is impossible to anticipate everything that might happen; some things must be left to the interviewer's discretion. The potential for response effects due to differences in interviewer behavior is real, even in the most tightly structured survey.

The task should be defined carefully by the investigator. Task definition is primarily a matter of what questions are asked; how they are asked, that is, their form and wording; the order in which they are asked; what is said by way of introduction to the survey or to particular questions; and the mode of administration of the questionnaire. It is also the source of the largest response effects (Sudman and Bradburn, 1974).

Let us look at some of the ways in which question wording and context may affect response validity. Consider the following question: "In which of these groups did your total family income, from all sources, fall last year--before taxes, that is?" [Respondent is shown a card with income categories on it from which to choose.]

Several things that might bias the reporting of income occur immediately to us. First, respondents might deliberately omit some types of income that they do not want anyone to know about, such as income from

illegal sources or income not reported on their income tax returns. They may forget about some income or report estimates of income where good records are not readily available. A third problem may arise from misunderstanding the question or not defining terms the same way as the investigator intended. For example, should gifts, inheritances, or insurance payments be reported as income? What about noncash income? The question may not make clear what the investigator has in mind when asking about income. Respondents may also include income from the wrong time period. That error is in remembering when something happened, rather than whether it happened. Finally, some respondents may deliberately misreport their income to impress the interviewer or to make themselves look better off than they are or vice versa.

We can summarize these types of errors by noting that they fall into three classes: (1) deliberate errors in which the respondent adds, omits, or distorts information in order to make a good impression on the interviewer or to prevent the interviewer from finding out something about him; (2) memory errors, which may be about whether something happened or when it happened; and (3) communication errors, that is, either the investigator does not make clear to the respondent what is being asked or the respondent fails to make clear the response to the interviewer so that a wrong answer is recorded.

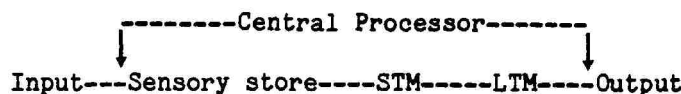
Deliberate errors come from the respondent (and occasionally from the interviewer). Memory errors come from the respondent, but their magnitude may be affected by the way in which the task is defined, that is, by the cues that are given in the questionnaire. Communication errors may come from respondents who do not make their responses clear to the interviewer, from interviewers who do not make clear what they are asking, or from the task, that is, the questions as formulated do not communicate the intended meaning to the respondents.

Research on cognitive processes is most relevant to understanding memory and communication errors. We shall thus concentrate our attention in this paper on them.

### **Conceptual Model for Human Information Processing**

There are a number of models for human information processing; they differ in their details. We shall adopt a very general model that seems to be common to almost all approaches, without attempting to suggest that any particular model is the more nearly correct one.

Information enters the system via one or more sensory modalities. The entire input is held briefly as a sensory store from which information is selected for transfer to short-term memory (STM). The contents of the STM are operated on as appropriate for the cognitive tasks being processed and may be encoded and transferred to long-term memory (LTM) and/or used in further processing, leading eventually to some behavioral output. Presiding over this system is some sort of central processor that performs the operations on the input, analyzes the information, and directs the output. Schematically the system looks something like this:



For the purposes of studying response errors in surveys, the critical parts of the model are the central processor, STM and LTM, and their interrelationships.

The central processor is conceptualized as the part of the system that directs cognitive processes, performs the logical operations, and generally does those things we lump under the rubric of “thinking.” The STM is roughly what we mean when we talk about the focus of attention or what in older terminology might be thought of as consciousness. LTM is here meant to refer to the store of past experience and learning that is somehow mentally represented so that it can be used by the central processor in thinking.

Following Tulving (1972) we will distinguish two LTM subsystems—semantic and episodic memory. Semantic memory is the system concerned with storage and utilization of knowledge about words and concepts, their properties, and their interrelationships. Episodic memory is concerned with temporally dated, spatially located, and personally experienced events or episodes and temporal-spatial relationships between such events. It may be convenient later to distinguish more than two subsystems of LTM, but two is probably the minimum.

Since a survey interview consists of a sequence of questions and answers, we look to cognitive studies for understanding of comprehension of oral or written text on the input side and of information retrieval from episodic memory on the output side.

There are several additional concepts that are useful for thinking about information processing related to survey questions. One of the most important of these is that memories are stored in some organized structured form and not as isolated units. Following Bartlett (1932) we call these structures schemata. Schemata are basic to the organization of both semantic and episodic memory. Words and events may be encoded in many different schemata. Stimuli that activate searches of memory are assimilated to particular schemata that direct the search within their structure. If the wrong schemata are activated, the search may be slowed considerably or fail altogether. For example, if one is involved in moving one's office and someone asks, “Where did that table go?,” one will begin to think about what has happened to pieces of furniture that were previously there. If, however, in the middle of a conversation a research assistant comes in looking for some computer output and asks, “Where did that table go?,” thinking about pieces of furniture may fail to produce the information desired.

Models of the retrieval process are generally couched within the framework of the encoding specificity principle posited by Tulving and associates (e.g., Tulving and Thomson, 1973; Tulving and Wiseman, 1976). In its strong form the principle asserts that retrieval cues will be effective at retrieval only if they provide the context at the time of encoding. This strong version has been challenged (see Burke and Light, 1981, for a discussion), but there is abundant evidence that recall performance can be improved by the reinstatement of cues at retrieval that were present at encoding.

About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.



A distinction is also made between recognition and recall. Recognition requires a judgment that a particular stimulus occurred (e.g., a word was on a list), while recall requires a memory search to reproduce a particular response. Most retrieval models imply that recognition will be better (faster and more accurate) than recall because one has only to recognize the stimulus rather than to generate and then recognize it.

Tulving and Thomson (1973), however, have shown that there are conditions under which recall can be better than recognition (e.g., when the cue in the recognition task activates the wrong association). In the relevant experiments subjects were shown three successive lists of word-pairs. The first two lists were designed to induce subjects to encode each target word with respect to another word. The target words, each paired with its cue word, were shown visually one at a time. Immediately after the end of the presentation of the third list the subjects were provided with 24 haphazardly ordered cue words on a recall sheet and asked to write down the target words. The cue words on the list had a weak association with the target words (e.g., ground/COLD). In a subsequent test phase, words not on the list but having a strong association to the target words (e.g., hot/COLD) were used as retrieval cues. Subjects were then given a variety of tasks that asked them to recall target words, to generate words through associations with the cue words, and to recognize words that they had generated by their associations as belonging to the target list. These experiments demonstrated conditions under which subjects could recall words but did not recognize them as belonging to the target list. The number of words that were recalled but not recognized was about 15 times the number of words that could be recognized but not recalled.

These findings might have some application to understanding the differences between responses to open-ended and closed format questions. Sometimes, questions with lists (e.g., problems facing the U.S., worries) produce different response frequencies than do free-recall questions. The cognitive processes that produce these different response frequencies are not well understood.

Although it is clear that much information that reaches the sensory store is lost and never enters even STM, it also appears that many things enter memory without explicit (conscious) instructions to store them. There is no clear understanding of what determines how well something will be remembered. Rehearsal, labeling, and breadth of relational networks seem to be important in influencing memory, but their relationships are not well worked out. Temporal aspects of memory are important to survey researchers because many surveys involve questions about frequency of events or events that occurred at particular times. The placement of events in time may depend on events being coded with specific markers (e.g., dates) or be inferred from trace strength (e.g., the more vivid the memory trace, the more recent the event). The growing literature on temporal memory will be discussed later.

### **Some Differences Between Survey and Cognitive Research Traditions**

Before proceeding to a more explicit examination of survey methodological problems from a cognitive point of view, we should note some differences in the research traditions of the two fields that make applications of the findings from one field to the other difficult. While these differences are not hard and fast, they do represent emphases that shape the problems addressed and the methods used to study them.

Survey researchers typically focus on properties of the stimulus (questions) that might affect the valid processing of information. These properties can be classified into three categories: (1) wording of the question; (2) structure of the response alternatives; (3) context of question, e.g., other questions asked, the instructions to the respondent, the setting in which the interview is conducted. The concern here is that the meaning of the question be the same to all respondents, and, of course, be the same meaning as that intended by the investigator. We might rephrase this concern by saying that survey researchers are concerned that they are setting the same cognitive task for all of the respondents. There is relatively less attention paid to whether or not a particular cognitive task is possible for the respondent or whether it is an easy or difficult task.

Cognitive researchers, however, typically focus on properties of the processing system that affect the way in which information is handled. These may also be classified into three categories: (1) encoding strategies; (2) retrieval strategies; (3) stores (lexical, semantic, episodic or archival). Experimental cognitive research is oriented toward studying microprocesses of information processing. Survey methodological research is oriented toward studying macroprocesses of comprehension and information retrieval.

Much of experimental cognitive research involves laboratory tasks that are often artificial because of the need to control the effects of other variables. The subjects tend to be college students or others who have a fairly high level of education (and presumably intelligence). There is thus a question about the degree to which findings from the laboratory can be translated to the complexities of a field survey.

There are a considerable number of field experiments in the survey methodological literature. The most common method is to use alternate forms of questions (split ballots) with random assignment of question forms or wording to 1/nth of the sample. These experiments typically have much larger sample sizes than the cognitive laboratory experiments and a more heterogeneous set of respondents on such characteristics as age and education. They have, however, very little control over (and often little knowledge about) the past learning of the respondents.

### **Applying Cognitive Research to the Study of Response Effects**

Because we are reviewing sources of error in survey questioning in the light of cognitive research, we organize our discussion more in line with the questions that are most studied by survey researchers rather than the other way around.

## Question Wording

We noted above that one of the major sources of errors in surveys arises from imperfect communication of meaning from the interviewer to the respondent. Sometimes this miscommunication results from a failure to specify what is to be included in a particular concept (e.g., what income is to be included in a report of family income), and sometimes it arises from the imprecision of the concept itself (e.g., in questions about concepts such as liberalism, defense policy, confidence in a particular institution, etc.). Respondents often recognize that the referent of the question is ambiguous and ask for clarification. In response to the question, "What things do you like best about living in this neighborhood?" a respondent might ask sensibly, "What do you mean by 'this neighborhood'?" Unless the researcher has a specific definition in mind and supplies it to the respondents, respondents are usually told to define the concept for themselves. This, of course, is what people do implicitly with the concepts in all questions. The range of interpretations may be considerably greater than the investigator realizes, however, and differences in interpretation may make responses difficult to analyze. Of course, an analogous problem exists in understanding the responses given by respondents. The researcher needs to know how respondents interpreted the question in order to know how to interpret the answer.

The comprehension process is generally conceptualized as the result of an interaction between the input material (i.e., the question) and previous knowledge. Bransford and Johnson (1972) have shown how previously incomprehensible descriptions become comprehensible when the proper context for understanding them is given. They used the following passage (p.400):

The procedure is actually quite simple. First you arrange items into different groups. Of course, one pile may be sufficient depending upon how much there is to do. If you have to go somewhere else due to the lack of facilities, that is the next step; otherwise, you are pretty well set. It's important not to overdo things. That is, it is better to do too few things at once than too many. In the short run this may not seem important but complications can easily arise. A mistake can be expensive as well. At first, the whole procedure will seem complicated. Soon, however, it will become just another fact of life. It is difficult to foresee any end to the necessity for this task in the immediate future, but then one can never tell. After the procedure is completed, one arranges the materials into different groups again. Then they can be put into their appropriate places. Eventually they will be used once more and the whole cycle will then have to be repeated. However, that is part of life.

College students had a great deal of difficulty understanding this passage unless they were told that it was about washing clothes before they heard it and thus knew the context within which to process the information.

Pichert and Anderson (1977) showed that the events that are recalled from a passage can be affected by manipulating the perspective taken by

the respondent. They used a passage describing two boys playing in a house. One group of respondents reads the passage from the perspective of a home buyer; the other group from the perspective of a burglar. Each group recalled details that were appropriate to the perspective from which they were reading the passage.

Different Comprehension of the Same Term A good example of the differences in comprehension of the same term is given by Fee (1979), who investigated the meanings of some common political terms used in public opinion studies. She adopted a variant of a technique developed by Belson (1968) in which respondents were asked to elaborate on their understanding of terms such as “big government” or “big business.” She found that there were clusters of meanings, which gave rise to different interpretations of questions. For example, she found that the term “big government” had four distinct images: in terms of welfare, socialism, and overspending; in terms of big business and government for the wealthy; in terms of a combination of federal control and diminished states' rights; and in terms of bureaucracy and a lack of due process. Different kinds of people tended to hold different images, which in turn were related to different attitudes. Unfortunately there did not seem to be any way of determining in advance how a respondent would interpret the term.

Different Questions on the Same Topic In the early 1950s, both Gallup and the National Opinion Research Center (NORC) asked questions about support for the Korean War. The Gallup question was: “Do you think the United States made a mistake in deciding to defend Korea, or not?” The NORC question was: “Do you think the United States was right or wrong in sending American troops to stop the Communist invasion of South Korea?” The NORC question consistently drew more support for the war than the Gallup question.

There are actually three important differences between these two questions that may have affected the responses to them. The one to which most attention has been paid is the addition of a reference to stopping the Communist invasion. Researchers have generally found that there is greater approval for American foreign policy decisions when the decision is described as “attempting to stop Communism.”

A second difference is the use of the terms “right or wrong” in the NORC question and the use of the term “mistake” (“or not”) in the Gallup question. We do not know whether “wrong” is conceptually equivalent to “mistake,” particularly when it is paired with “right.”

Finally the Gallup question is asked in a form that implicitly presumes that the policy is a mistake, that is, the respondent has to deny the proposition in the main part of the question in order to support the war. In the NORC question, the respondent has to choose between two alternatives, one positive and one negative.

What can we say about these differences from the cognitive processing point of view? Regarding the use of politically sensitive terms we might hypothesize that the evocation of a new element (reference to stopping Communist invasion) causes more respondents to assimilate the question

Communist invasion) causes more respondents to assimilate the question partially to the schema “stopping Communism” than would be the case without the explicit reference. If the schema “stopping Communism” is highly approved, then we would expect that more approval would be expressed when the phrase is explicitly part of the question than when it is not. One might also ask whether the term “defend” is assimilated to the same schema as “sending American troops?” In order to understand fully the effect of the different wordings, one would need to know the different schemata involved in processing the text of the questions and the underlying evaluation of the different schemata. Schuman and Presser (1981) present further evidence of wording changes that change the apparent meaning of otherwise similar questions.

The problem of synonyms is a difficult one for question wording. Enough is known from methodological studies of question wording to know that the path is fraught with booby traps (Payne, 1951). Similar, if not exactly synonymous, terms that indicate a positive orientation toward an attitude object can have different connotations and yield different responses. For example, the terms “approve” and “disapprove,” and “like” and “dislike” are frequently used in attitude questions, although little attention has been paid to different implications of the two terms. At least one use of the two terms suggested that they are not exactly the same (Murray et al., 1974).

Choosing between two alternatives appears to be a more difficult cognitive task than simply affirming or denying the truth of a proposition. Thus slower and perhaps more thoughtful answers may be given when the respondent has to choose between two or more alternatives. Also the constraints of the interview situation may tend to make it preferable for the respondent to agree rather than disagree with a statement, particularly if it is a proposition that they do not have strong (or any) opinions about. Systematic research on these aspects of question wording could be fruitful.

Antonyms are also troublesome. Approval of a positively worded attitudinal statement is frequently not the same as disapproval of the same attitude when it is expressed negatively using an antonym of the positive wording. One of the best-known examples is that given by Rugg (1941): “Do you think the United States should allow public speeches against democracy?” “Do you think the United States should forbid public speeches against democracy?” To the first question, 21 percent responded that speeches against democracy should be allowed; 39 percent said that such speeches should “not be forbidden.” Similarly, 62 percent said that such speeches should “not be allowed,” but only 46 percent said that they should be forbidden. (Other respondents had no opinion on the matter.) Clearly the concepts “allow” and “forbid” have somewhat different connotations so that their negatives are not equivalent. We need to have more systematic research about how the processing of the negative form of statements differs from the processing of the positive form and how negations differ from statements using an antonym.

The General and the Specific Questions may differ on a generality specificity dimension, and the degree of specificity contained in the question wording appears to affect responses. An example of the effect of increasing specificity can be seen in questions from a Gallup survey in May-June 1945:

Do you think the government should give money to workers who are unemployed for a limited length of time until they can find another job?

[yes, 63 percent; no, 32 percent; don't know, 5 percent]

It has been proposed that unemployed workers with dependents be given up to \$25 per week by the government for as many as 26 weeks during one year while they are out of work and looking for a job. Do you favor or oppose this plan?

[yes, 46 percent; no, 42 percent; don't know, 12 percent]

Would you be willing to pay higher taxes to give unemployed persons up to \$25 a week for twenty-six weeks if they fail to find satisfactory jobs?

[yes, 34 percent; no, 54 percent; don't know, 12 percent]

Since these questions do not differ in single elements but change several levels of specification in each form, it is impossible to say what elements were causing the changes in support. The key elements in these questions appear to be the referent for support (unemployed versus unemployed with dependents); the amount of time for support (a limited length of time versus 26 weeks); the amount of support (an unspecified amount of money versus \$25 per week); the stopping rule (finding another job versus finding a satisfactory job); and the mode of financing the payments (not specified versus higher taxes). A systematic investigation of different elements such as these is typically not done in surveys.

Even with the same question wording, respondents may vary in the degree of generality with which they interpret questions. For example, Smith (1980) reports the following results from the 1973 and 1975 General Social Survey (GSS). Respondents were first asked: "Are there any situations that you can imagine in which you would approve of a man punching an adult male stranger?" Respondents were then asked a series of questions about specific conditions under which they would approve of such an action, such as the stranger hit the man's child after the child had accidentally damaged the stranger's car or the stranger was beating up a woman and the man saw it. Many respondents who said that they would not approve of hitting a stranger "under any situations that they could imagine" went on to express approval for hitting when specific situations were described. Smith suggests that many respondents are not interpreting the general question as literally asked but responding instead to the absolute phrase "Are there any situations that you can imagine" as if it meant "in situations that you can easily think of" or simply "in general."

The question as stated poses a difficult cognitive task for respondents since it is asking for them to imagine a very large range of



situations. Typically the flow of questioning does not give respondents very much time to think about the question, so, as Smith suggests, the effect is to limit responses to the situations that can easily be thought of. The task is also one that is not engaged in every day, so we would expect that the effect would be more pronounced among those who are less “imaginative,” whatever that might be correlated with, perhaps age or education. The specifications ease the cognitive task by providing what imagination was not able to do, that is, some specific situations in which hitting a stranger might be approved.

One way to ease the cognitive task and perhaps equalize the task more nearly for all respondents would be to give a longer introduction to the topic. Such an introduction might cite some of the more common examples where people are apt to approve, or at least tolerate, hitting a stranger. This is also an example of the type of question that might display a substantial order effect depending on whether the general question is asked before or after the more specific questions. (The topic of order effects is discussed below under the section on context effects.)

The mention of specific situations might be thought of as an example of prompts in aid of memory. In questions that require a fairly substantial amount of effort to answer completely, prompts in the form of examples of category instances are often given. For example, in the GSS the question on membership in voluntary organizations is asked as follows: “Now we would like to know something about the groups and organizations to which individuals belong. Here is a list of various kinds of organizations. Could you tell me whether or not you are a member of each type?” Respondents are then handed a card with 16 types of organizations on it, and they are asked about membership in each type of organization. Even this type of prompt may fail to elicit good recall since the respondents still must have the particular organizations that they are members of coded into the categories presented, and the cue of the category must trigger a search that will retrieve the fact that they are members.

This type of prompting is called aided recall by survey researchers, and it produces a higher level of response to recall questions than forms of the question without such aids. The use of aided recall has an unfortunate side effect, however, for items involving time-bounded phenomena, e.g., magazines read in the last month, number of visits to the doctor in the last three months, etc.: it increases “telescoping,” that is, reporting instances of the event from the wrong time period. (Problems related to memory for the timing of events are discussed below.)

How Cognitive Theory Might Help The processing of information on the part of respondents involves two levels of processing. First, the question must be processed according to the rules that we use to understand spoken or written language, that is, the question must be understood. Second, the representation of the question must be processed according to other rules (e.g., Kintsch and van Dijk's, 1978, macrostructures) that retrieve the information necessary to answer the question and perform the operations necessary to make the judgments and produce the answer to the question.



The wording of the question provides the stimulus that sets off a complex set of cognitive processes which finally result in a response. Tulving and Thomson (1973:352) note: "Retrieval operations complete the act of remembering that begins with encoding of information about an event into the memory store. Thus, remembering is regarded as a joint product of information stored in the past and information present in the immediate cognitive environment of the rememberer." They go on to develop the encoding specificity principle (p. 353): "What is stored is determined by what is perceived and how it is coded, and what is stored determines what retrieval cues are effective in providing access to what is stored."

At first look, the encoding specificity principle suggests that we will have great trouble in using standardized questions in surveys since it implies that there could be great individual variation in the way people encode events. Without denying that there may be some individual variation, it is likely that respondents from a similar culture and who speak the same language will encode most everyday events in much the same way. If this were not the case, we would not be able to communicate with one another as well as we do.

The encoding specificity principle is useful in reminding us to pay more attention to the ways in which events are encoded and to develop the wording of questions so as to match those codes. For example, significant subgroups in the population (e.g., those identified by ethnic group, region, education) may have different ways of encoding some types of events so that we might have to alter question wording for different groups. Bradburn et al. (1979) did find somewhat higher (and probably more accurate) reports of use of substances such as drugs and alcohol when respondents were asked questions phrased in their own terms.

The encoding specificity principle may be most important in relation to attitude questions for which the concepts are less well defined and have fewer shared behavioral referents. It suggests that we need to make much greater efforts to map the ways in which respondents encode various attitude objects that we are studying, particularly if they are to be studied over some period of time. Since we frequently want to study change in attitudes over some period of time, it is extremely important that we understand better how to preserve the meaning of questions as events change.

The Focus of Attention in Questions For questions that require recall of past events, question wording may affect responses by directing attention toward or away from the major retrieval categories that are of interest in the question. For example, one of the purposes of an experiment by Vogel (1974), as part of a survey of Wyoming farmers in a Department of Agriculture study, was to study the effect of question wording on estimates of the calf crop on Wyoming cattle ranches in the year 1974. There was reason to believe that the number of calves being raised of less than 500 pounds was underreported. Two forms of the question on calves was used. The question asking about calves weighing less than 500 pounds was changed in the experimental version to emphasize the word "calves" because livestock producers often did not report unweaned calves since they considered the cow and the calf to be one animal unit before

weaning. The question was part of an inventory question whose introduction read:

Please report below all cattle and calves on the land you operate, regardless of ownership (include those now on feed). Also include those owned by this farm or ranch that are now on public grazing land. How many are:

.....

3. beef cows? (Include heifers that have calved)
4. milk cows? (Include milk heifers that have calved)

.....

.....

8. HEIFER, STEER, and BULL CALVES weighing less than 500 pounds.

The alternate version was:

8. CALVES--heifer, steer, and bull calves weighing less than 500 pounds.

In the experimental version the superordinate category of interest (calves) was placed first with the subordinate categories (heifer, steer, and bull calves) given less emphasis by being placed later in smaller type. The experimental version produced about a 10 percent higher estimate for calves weighing less than 500 pounds.

A further complication to the recall process is that remembering is believed to be a reconstructive process. Loftus and her associates (Loftus, 1977, 1980; Loftus and Palmer, 1974) have shown that memory for events can be supplemented and altered through the introduction of post-event information. The incorporation of post-event information can be attenuated to the extent that respondents are more certain about particular memories if recall has been made prior to the introduction of the additional information and if the source of the information is clearly suspect (Loftus, 1982).

### Response Categories

One of the oldest and most puzzling phenomena in survey research concerns the differences between open- and closed-ended questions. In the open-ended question respondents are given no response categories and simply answer the questions using their own terminology. Closed-ended questions give the respondents a set of alternatives from which to choose answers.

The open and closed format difference is similar to the recognition/ recall distinction in cognitive research. In the open format the respondent has to retrieve from memory the material that is asked for in the question. Consider a question such as, "What are the most important problems facing the country today?" Asked in an open-ended format, such a question requires a considerable amount of cognitive work on the part of the respondent to define the limits of the concepts being called for in the question, then initiate a memory search for instances of these

concepts, finally producing the responses. With a precoded question the categories of interest are much more delimited by the response categories offered. The respondent simply has to recognize which particular category is a true instance of behavior or attitude. As with the recognition task, precoded questions appear to be easier for respondents to handle, and in many instances (but not all) they produce as good if not better responses than do open questions. Aided recall generally produces fuller and more accurate responses than do questions without aids, assuming, of course, that the aids are adequate for the recall task. Depending upon the type of list used for aids, the task may be thought of as a pure recognition task or as a recall task with many cues to aid the memory search.

There is one type of behavior report question in which it is clear that the open-ended format does much better than precoded response categories. An experiment by Bradburn et al. (1979) compared the effects of precoded response categories and verbatim recording on behavior reports about topics for which there is a presumption of considerable underreporting, such as alcohol consumption. An example of the type of question used involved those who reported drinking beer during the previous year: "On the average about how many bottles, cans, or glasses of beer did you drink at one time?" In the precoded version the codes were "1," "2," "3," "4," "5," "6 or more." In the open-ended version there were no codes; the interviewer simply wrote down the number given by the respondent.

Estimates of beer consumption based on open-ended questions were about 60 percent higher than those based on the precoded responses. The reasons for this difference are not entirely clear. It is probable that the distribution of alcohol consumption has a long tail on the high side; thus those who report more consumption than allowed for by the highest precoded category tend to increase the average. It is also possible that the precodes are interpreted as an implicit norm and that some high consumers are reluctant to place themselves in the highest category, particularly if they did not consider themselves to be heavy drinkers.

With the type of questions studied by Bradburn and associates in which open-ended questions produced better responses than precoded questions, the questions were quite specific in delimiting the categories of recall and the responses were estimates of quantity or frequency. Thus the response dimension was well specified even though not explicitly given by response categories.

More equivocal results have been found in a series of experiments conducted by Schuman and Presser (1981). They studied questions that involved multiple nominal responses to broad inquiries about values and problems, for example: "What do you think is the most important problem facing this country at present?" "People look for different things in a job; what would you most prefer in a job?" "While we're talking about children, would you please say what you think the most important thing for children to learn to prepare them for life." The results from their series of experiments were complex, but it was clear that there were statistically significant, substantively important differences in marginal distributions between open and closed forms of this type of question. Schuman and Presser also demonstrated that the assumption of

form-resistant correlations, that is, that even though marginal distributions may change between question type, the correlations between responses and significant independent variables would be constant, is not tenable with respect to the open- and closed-ended question comparisons. For example, the correlation between education and work values such as high pay and steady pay varied by the form of the question used.

Comparing responses to open and closed forms may give us some insights into the cognitive processes at work in the two forms of question. A fairly consistent finding, at least for the broad attitudinal questions that Schuman and Presser studied, is that respondents give a greater range of answers in response to the open form than they do when they are given a list of precodes, even when there is a category for “other” responses. Respondents rarely use the “other” category in closed forms. Thus it would appear that one of the important effects of providing response categories is to define much more narrowly the range of categories that is to be included in the response to the question.

It is also likely that there may be some unintended direction in the open forms that can be corrected by the provision of specific alternatives. For example, Schuman and Presser found that the response category “crime” was much more frequently reported as being one of the most important problems facing the country when the question was asked in a closed form rather than an open form. They hypothesized that the reference in the open question to “facing this country” discouraged respondents from thinking of crime as a problem, because crime is perceived by many as a local rather than a national problem. Providing “crime” as an explicit category in the closed form made it clear that the investigator intended to consider it a national problem.

That open versions of questions are cognitively more difficult for respondents is suggested by the fact that in the Schuman and Presser experiments the more educated as compared with less educated respondents gave more responses to the open-ended questions and that there were more missing data in the open forms among the less- as compared to the well-educated. Such differences, however, are less likely to occur when the open questions are extremely focused and require only a simple numerical answer.

Another issue of considerable importance and controversy in attitude measurement is the offering of an explicit “don't know” or “no opinion” category to the respondent. Research on this issue (See Schuman and Presser, 1981:Chapter 4) indicates that providing an explicit “don't know” (DK) option--called a filter--substantially increases the proportion of respondents who give DK responses. Typically, the DK increment is around 20 percent, but this appears to be unrelated to question content or the DK level that prevailed in versions without an explicit DK filter. Schuman and Presser introduced the term “floater” to refer to persons who give a substantive response to an item in a standard unfiltered form and a DK response to a filtered version of the same question. They suggest two models that might account for the nature and identity of floaters. One is a trait model that conceives of floaters as a distinct group who have a high propensity to say “don't know” when offered this choice on a filtered question but will not volunteer such a

response in the absence of an explicit alternative. The second model, which they call “threshold process model,” suggests that there is a process of floating that is created by the form of the question for any given item. This model suggests that there are a number of variables that may influence a respondent's position on a DK propensity dimension for particular item contents. The actual frequency of DK responses depends in part on respondent's position on this propensity dimension and part on the height of the barrier created by the question form. As Schuman and Presser point out (1981:146):

Where the question form strongly discourages DK responses only those very high on the dimension will insist on giving a DK response. Where the question form facilitates or even encourages DK responses (e.g., by providing a full DK filter preliminary to the opinion question) those with both high and moderate positions on the dimension will choose the DK option. In both instances the underlying dimension is the same, but the question form creates different cutting points on it. Moreover, the difficulty of the question content (remote versus familiar issues) will also greatly affect cutting points on the dimension, though, according to the model, equally for both standard and filtered question forms. Thus the DK propensity dimension influences the giving of DK responses on both forms, but there is no special trait that distinguishes floaters as such, nor is there a special group of people to be set apart and described by this term.

Additional questions of interest to survey researchers concern the effects of different response categories, such as the use of middle alternatives as opposed to dichotomous agree/disagree or approve/ disapprove type questions, the effects of using differing numbers of points on rating scales, and the use of familiar analogies such as ladders or thermometers for giving ratings of intensity of favorableness or unfavorableness toward specified attitude objects.

### Contextual Meaning

Questions in surveys are asked in some context. Part of this context is set by the introductory material that the interviewer presents in order to gain the respondent's cooperation. The order in which questions are asked provides a context that will affect the interpretation of the question or provide different cues for retrieval processes. Question-order effects have been widely studied among survey researchers, although until recently there has been relatively little theoretical orientation. A fairly typical example of the types of effects found was reported by Noelle-Neumann (1970), who examined the designation of various foods as particularly “German.” This study was part of an exploration of the image of three basic foodstuffs. In one form of the questionnaire respondents were asked first about potatoes, then about rice; in another form of the same questionnaire this order was reversed. When respondents were asked about potatoes first, 30 percent said potatoes were particularly “German.” However, when respondents were

asked about rice first, 48 percent said that potatoes were particularly “German.” A similar order effect was found for the pair, “noodles/rice.” Findings of this sort suggest that the initial general question, e.g., “Are the following foods particularly German?” activates a general schema relative to the “Germanness” of foods. However, this schema is relatively vague in the beginning and becomes sharpened as different examples are presented. When examples that are particularly atypical (rice) come first, the various attributes of the schema become clearer--a kind of contrast phenomenon.

A persistent mystery is that order effects tend to be asymmetric, that is, they only affect one question or response rather than both of the questions that are rotated. In the above example, order affected the proportion reporting potatoes as “German” but not the proportion reporting rice as “German.”

More recently it has been observed that when a general question and a related more specific question are asked together, the general question can be affected by position while the more specific question is not. For example, in the 1980 GSS, the following two questions were asked:

Taking all things together, how would you describe your marriage? Would you say that your marriage is very happy, pretty happy, or not too happy?

Taken all together, how would you say things are these days--would you say you were very happy, pretty happy, not too happy?

The order of these two questions was rotated in split-ballot experiments. The results indicated that responses to the general question about overall happiness were affected by the order in which the questions were asked, while the question on marriage happiness was not affected by the order. The direction of the shift, however, was not the same in different experiments (see Schuman and Presser, 1981:Chapter 2). One explanation is that when the general question comes first, the question is in fact viewed as including one's whole life, including marriage. However, when the specific question comes first the overall happiness question may be interpreted in a more narrow sense, referring to all aspects of life other than marriage.

Somewhat similar results have been reported by Schuman and Presser (1981) for general and specific attitude items related to abortion. Here higher levels of support were found for a general question about abortion than when the general question came before a series of specific questions about approval of abortion in particular situations, such as in cases of rape or damage to the mother's health. When respondents answer the general question first, it is clear that subsequent specific questions are subsets of the general one; answers to the specific questions may well be different from answers to the general ones. They can be interpreted in their specific sense, independent of what has gone on before them. The general question, however, may be interpreted as really general, i.e., including all of the specific items that are subsets of the general attitude, or it may be interpreted as all the other things



not included in the previously asked specific questions. Unfortunately we do not yet know the range of phenomena to which this type of order effects apply.

With regard to reports of behavior or events, order can also affect responses not only by creating the context for interpretation of the question but also by providing greater cues or time for retrieval processes to take place. In the Wyoming cattle and calf survey mentioned previously (Vogel, 1974), two question orders were used in order to estimate the number of calves born on the land the ranchers operated. In one form of the questionnaire the total number of calves born during the year 1974 was asked before a series of detailed questions asking how many of these calves were still on the farm or ranch and how many had been sold or slaughtered or how many had died. When these detailed questions were reversed and the total was derived by adding up the categories, the estimates for total calves born during the year was about 10 percent higher than when the question on total number was asked first. Similar results occur when respondents are asked a detailed set of questions about whether they get income from different specific sources, such as wages and salaries, savings accounts, transfer payments, etc., before they are asked their total income. The detailed questions appear to act as reminders and facilitate the memory search in order to come up with the requested information.

Marquis et al. (1972) and Bradburn et al. (1979) have found that lengthening the introduction to questions improves reporting of such things as symptoms, utilization of health care facilities, and reports of alcohol and drug consumption. It seems likely that the longer introductions not only direct respondents' attention toward the information requested and start the search process but also give the respondent more time to do the actual retrieval operations.

Unfortunately for survey practice, many of the recommendations based on research on retrieval processes point in directions that are antithetical to other pressures in conducting surveys. Techniques that improve recall, such as more detailed questions, longer questions, and giving the respondent more time to think about answers, also increase the length of the interview and thus costs. There is also concern on the part of the U.S. Office of Management and Budget, which must grant clearance for surveys conducted under government contract, that respondents not be overly burdened by surveys. Respondent burden has typically been defined in terms of length of interview without consideration for difficulty of the tasks.

The setting within which the interview is conducted may also create a context that facilitates or inhibits accurate reporting. The presence of others is the desired case for the National Health Interview Survey, because they may produce interactions that stimulate memory and produce overall better reporting. For the reporting of some types of behavior, however (either sensitive or embarrassing behavior, for instance), the presence of others may serve as an inhibitory factor and reduce reporting.

In their review of methodological studies on response effects, Sudman and Bradburn (1974) found a slight overall negative response effect (that is, a net underreporting) for surveys in which another adult was present



during the interview. In a fragmentary finding, Bradburn et al. (1979) found that the presence of others did not generally affect level of reporting, but that the presence of children seemed to make respondents uneasy about discussing sensitive behaviors such as drug and alcohol consumption and sexual behavior. The presence of adult third parties seemed to stimulate higher item-refusal rates.

### Time and Frequency

In surveys, one of the most important tasks facing the respondent is recalling the correct time period in which the behavior in question occurred. Errors that occur when respondents misremember the time period during which some event occurred are called telescoping errors. While one might imagine that errors would be systematically distributed around the true time period, the errors are more in the direction of remembering an event as having occurred more recently than it did. Thus the net effect of telescoping is to increase the total number of events reported in the more recent time period. This is partly because factors that affect omissions, such as the salience of events, also affect the perception of time (Brown et al., 1983). The more salient (or more frequent) the event, the more recent it appears to be. The result is that there will be overreporting of more salient or more frequent events.

Some of the best studies of telescoping effects on memory have been done by Neter and Waksberg (1964). They developed a procedure, called “bounded recall,” for reducing the effects of telescoping. Bounded recall procedures require panel designs in which respondents are interviewed several times. At the beginning of the second or later interviews (bounded interviews), respondents are reminded of what they said about their behavior during the previous interview and then asked about additional behavior since that time. For example, in the initial interview, respondents may be asked about the number of visits to physicians in the last three months. At the second interview, three months later, respondents review their responses from the previous interview and then are asked, “Since we last talked, how many times have you seen a physician?” The bounded-recall procedure requires a considerable amount of control over the data in order to provide interviewers with the correct information from the previous interview and thus has not been used as often as it deserves to be.

Respondents in surveys are not only asked to report on their behavior or attitudes, but they are also asked to make judgments about how often something has happened or how frequently they feel some way or did something. Some research has been done on the accuracy of time reporting, but most attention has been focused on increasing the validity of behavioral reports for the number of events that occur within a particular time period, e.g., how many times did you visit a doctor during the past 90 days? Relatively little attention has been paid to the cognitive problems of reporting frequency of subjective states, e.g., how often during the past few weeks did you feel bored?

Characteristics of the Ability to Make Frequency Judgments The general consensus among cognitive psychologists is that people are quite good at making relative frequency judgments concerning a variety of events for which frequency information was acquired under incidental learning conditions, that is, under experimental conditions in which subjects are not explicitly asked to attend to the frequency of events.

Studies that have addressed the issue of people's abilities in the domain of frequency judgments have examined two types of events: (1) experimentally controlled events that occur from one to ten times, which are here referred to as "experimental--low frequency" events (E-LF); and (2) naturally occurring events that are generally of high frequency (e.g., the letter "a" in spoken English), which are here referred to as "naturalistic-high frequency" events (N-HF). For E-LF events, judged frequency increases with actual observed frequency (Hasher and Chromiak, 1977; Hintzman, 1969; Hintzman and Stern, 1978; Johnson et al., 1977; Warren and Mitchell, 1980). Most of these studies used printed or spoken English words as stimuli, although pictures have also been used and give similar results. The same basic result is found with N-HF events.

It has been argued that the encoding and updating of information concerning the occurrence of events is obligatory and automatic (e.g., Hasher and Zacks, 1979). Some evidence for this claim is provided by the fact that there is a high positive correlation between people's judgments and the actual occurrence of a variety of naturally occurring "events:" letter frequency (Attneave, 1953), words (Shapiro, 1969) and lethal events (Lichtenstein et al., 1978). More direct experimental evidence on this topic comes from studies that have compared the reporting of frequency information acquired under incidental and explicit learning conditions (Howell, 1973; Attig and Hasher, 1980; Kausler et al., 1981). In general these studies show no or only slight effects of explicit instructions to attend to frequency judgments.

The automatic nature of frequency encoding also apparently extends to the activation of higher order information. Alba et al. (1980) presented subjects with lists of words that were chosen so as to represent 3, 6, or 9 instances of various semantic categories. They found that subjects were able to judge categorical frequencies in a surprise posttest: judgments of category frequency increased as actual frequency increased, and mean judgments for items occurring with different frequencies were significantly different from each other.

It has further been argued that the ability to make frequency judgments does not involve a learned component. Developmental studies that compare the performance of children as young as kindergarten age with college-age adults generally do not find evidence for an age-by-frequency-of-presentation interaction (Hasher and Chromiak, 1977; Hasher and Zacks, 1979). This seems to indicate that the frequency-of-repetition variable affects frequency judgments during childhood and adulthood in similar ways. Age differences in the absolute frequency of responses have been reported (main effects only), but these are not monotonic with age (Hasher and Chromiak, 1977).

There is a limited amount of evidence that the ability to make frequency judgments declines with age. In a comparison of older adults (60-75 years of age) with young adults (18-30 years of age), Warren and

Mitchell (1980) found an age-by-frequency-of-presentation interaction for absolute frequency judgments. This interaction was due to a slight loss of discrimination among items of differing frequency for the older group. Age interactions, however, have not been found when the task has been one of choosing the more frequent of two items (Attig and Hasher, 1980; Kausler et al., 1981).

In addition, practice effects have not been found for subjects making repeated frequency judgments (Erlich, 1964; Hasher and Chromiak, 1977). This also suggests that learning does not play an important role in frequency judgments.

Constraints on the Ability to Make Frequency Judgments A number of studies have produced results that show that some restrictions exist on people's ability to make frequency judgments.

In the first place it is generally known that people are not very accurate in their absolute frequency judgments. The typical result is that for E-LF events people tend to overestimate low frequencies and underestimate high frequencies (Hasher and Chromiak, 1977; Hintzman, 1969). There is some evidence that whether an event occurring with a particular frequency will be overestimated or underestimated depends on the frequency with which other events under consideration have occurred. That is, subjects seem to use some knowledge of the average frequency of events in the experiment. For example, Hintzman (1969) reports that in one experiment in which items were presented with frequencies ranging from one to ten times, subjects overestimated the frequency of items that occurred twice. In another experiment, however, where the highest frequency of occurrence was two, subjects tended to underestimate the frequency of items occurring twice.

The underestimation of low-frequency events and the overestimation of high-frequency events is also found with N-HF events (Attneave, 1953; Lichtenstein et al., 1978; Shapiro, 1969). Attneave (1953) asked subjects to judge the frequency of letters occurring in written language and gave them a standard of comparison of the occurrence of the average letter per 1,000 letters. Lichtenstein et al. (1970) provided their subjects with a standard of comparison of 50,000 deaths per year due to motor vehicle accidents or of 1,000 deaths per year due to electrocution. Shapiro (1969), however, did not provide his subjects with anchor points.

It is not clear what mechanism is responsible for the over- and underestimation. One possibility is that on some percent of trials subjects do not have adequate frequency information and that on those trials they adopt a strategy of using the mean frequency. This could account for the opposite biases on the two ends of the scale, provided that subjects are aware of an average frequency.

Subjects also tend to be more accurate at estimating low-frequency items than at estimating high-frequency items (Alba et al., 1980; Hasher and Chromiak, 1977), and the variability of responses tends to be much higher for higher-frequency items.

The absolute size and accuracy of frequency judgments has been found to be affected by the spacing of repeated occurrences of an event (i.e., a function of the number of intervening events between repetitions of an

event). The general finding is that for E-LF events increasing the interval between successive presentations leads to higher and more accurate frequency judgments (Hintzman, 1969; Jacoby, 1972; Rose, 1980; Underwood, 1969).

The context in which an event is presented also affects the absolute size of frequency judgments but not the discriminability of events occurring with different frequencies. The general finding is that for items of equal frequency, if repeated items are placed in a different context on each occurrence, then the judged frequency will be lower than if the context is the same on each occasion. For example, Hintzman and Stern (1978) presented names of famous people to be learned in the context of a descriptive statement of the person that was either the same or different on all occasions and had subjects make a truth value judgment about the statements; they found higher reported frequency for low-variability contexts. Jacoby (1972) and Rose (1980) report a similar finding. However, Rose (1980) finds equal discriminability for items of differing frequency under the two context conditions.

It is also possible to affect the absolute size of one's frequency judgment by generating the item silently or explicitly by writing or speaking it (Johnson et al., 1977; Johnson et al., 1979a, 1979b). The basic design of these studies included targets presented a variable number of times in the context of the same paired associate. Following a learning trial, some of the targets were tested by presenting the paired associate cue and having subjects write down the target. (Most subjects could do this on most trials.) The basic result is that test frequency (i.e., occasions for self-generated occurrence of an event) affects judgments of external event frequency and vice versa. This has been replicated for two other college-age samples by Johnson and associates and extended to 8- to 12-year-old children (Johnson et al., 1977; Johnson et al., 1979a, 1979b).

Lichtenstein et al. (1978) tried to determine the cause of the systematic overestimation of the frequency of occurrence of some lethal events. They found two variables that account for most of the variability in their subjects' responses: (1) the frequency with which they report learning about the event as a cause of suffering through the media and (2) the frequency with which they report having the death of a close friend or relative caused by the event.

Both the Johnson et al. studies (1977, 1979a, 1979b) and the Lichtenstein et al. (1978) results may be explainable by the availability heuristic of Tversky and Kahneman (1973). This states that frequency judgments may under certain circumstances (see below) be based on how easily instances of the event become available at the time of making the frequency judgment. To the extent that thinking or hearing about an event or having some personal involvement with an event increases its saliency, such conditions may make information more available, and therefore appear more frequent.

A study by Rowe (1974) suggests that the type or degree of processing given to targets in an incidental learning condition may also affect the absolute value of frequency judgments. Rowe had subjects either process items semantically (e.g., rate each target on how strongly it connotes strength) or nonsemantically (e.g., determine the number of syllables in

a target) and found that semantic processing resulted in higher frequency responses.

It is tempting to speculate that this effect may also be due to semantically processed items being made more salient to the subject and therefore, by the logic of the Tversky and Kahneman (1973) argument, being judged as occurring more frequently than less salient items.

Implications for Survey Methodology A number of the factors affecting frequency judgments may lead, under some circumstances, to an incorrect assessment of the occurrence of behavior by the survey researcher. Age differences in the absolute size of response (e.g., Hasher and Chromiak, 1977; Warren and Mitchell, 1980), whether they are attributable to response biases (Hasher and Chromiak, 1977) or to differences in the discrimination of frequencies (Warren and Mitchell, 1980), must be considered by the researcher investigating a behavior over a wide age span. The systematic under- and overestimation of events occurring with high and low frequency, respectively (e.g., Hintzman, 1969), can also be a problem if one is interested in obtaining data on the absolute rather than the relative frequency of occurrence of events. This will be true to the extent that events of different frequencies are compared, as for example, comparisons between physician visits reported by older and younger respondents who have very different frequencies of doctor visits or for visits to dentists as compared with visits to psychiatrists.

It is not clear whether a remedy exists for this systematic distortion of frequency judgments. The fact that this distortion is found under natural learning conditions when no reference point is given (Shapiro, 1969) suggests that this may be the result of the very process by which people make such judgments. It is interesting to note that an analogous distortion is found in the domain of time dating: Brown et al. (1983) found that subjects asked to place events on a bounded time scale tend to bring forward in time older events and move back in time more recent events, a phenomenon they have dubbed the squish effect. They have further found that this effect can be reduced, but not eliminated, by extending the boundary at the far end of the scale (i.e., the earlier time boundary).

Brown et al. (1983) also show that it is difficult for subjects to determine when salient, public events have taken place. Loftus and Marburger (1983) have undertaken some research to attempt to ameliorate this serious problem. They find that simply asking subjects to search memory for events occurring between two specific dates provides more accurate responding than using a general dating procedure such as "in the last six months . . . ." A more important contribution of the Loftus and Marburger (1983) paper is their use of landmarks, public or private, to mark the boundaries of the recall period. They found that the use of a public event, such as the eruption of Mt. St. Helens or New Year's Day, or a private event, such as the respondent's birthday, to mark the beginning of the recall period resulted in more accurate frequency reports than even the use of specific dates. It appears that landmark events are more useful than calendar dates in allowing respondents to make contact with other events in their own lives. More work is needed

to determine what exactly can be considered a landmark and how to circumvent certain inherent problems with the use of landmarks: generality for all groups and applicability of landmark for length of survey (for large surveys interviewing may continue for weeks or months).

Two general models of representation have been presented as the basis for frequency judgments. The first, the strength model, maintains that subjects judge the frequency with which an event occurred by consulting a unidimensional trace of the event and determining the “strength” of the trace rather like the vividness of an impression. All other factors being equal, a more frequent event will be associated with a stronger trace than a less frequent event. By this account, information concerning individual instances (e.g., nuances of context) is not part of the memory representation. The alternative model, multitrace theory, has been proposed in various forms (e.g., Anderson and Bower, 1972; Hintzman and Block, 1971). The different versions share the position that each occurrence of the same nominal event is represented separately in memory. The form of the representation may be in terms of list markers indicating individual occurrences attached to the permanent address of the event in memory (Anderson and Bower, 1972) or as separate episodes with some index to the effect that they are all instances of the same nominal event (Hintzman and Block, 1971).

The current position among experimental psychologists is that strength theory cannot adequately explain the relevant experimental results (Crowder, 1976). Perhaps the best evidence against strength theory is provided by two experiments by Hintzman and Block (1971, Experiments II and III). Results similar to those of Hintzman and Block (1971) have been obtained by Andersen and Bower (1972) for subjects also judging local frequencies of events.

Memory for details corresponding to the unique aspects of nominally identical events also is contrary to the predictions of the strength model. Clearly, people are able to retrieve information which differentiates among various instances of an event (e.g., Hintzman and Block, 1971; Linton, 1982). The distinction between semantic and episodic memory (e.g., Tulving and Thomson, 1973) is in part a reflection of this fact.

Strategies for Making Frequency Judgments Two general strategies have been described in the literature as methods of accessing memory traces for frequency judgments: (1) counting of individual instances and (2) use of the accessibility heuristic. The first requires accessing individual traces by some means (discussed below), determining whether the instance is of the type being searched, and keeping an accurate count of these instances. If subjects do use this strategy, data on accuracy of absolute frequency counts attests to the fact that it is an error-prone process. Some of the possible error sites are: the process of matching a retrieved instance to a description of the required one, tagging of already-counted instances to avoid double consideration, and updating of the counter. The second strategy, the availability heuristic (Tversky and Kahneman, 1973:208) is described as making “. . . estimates [of] frequency or probability by the ease with which instances or



associations could be brought to mind.” This method does not require actual operations of retrieval, only a judgment about the ease with which particular instances could be processed. Tversky and Kahneman (1973, Experiments I-IV) find high, positive correlations between subjects' judgments about the number of instances of various categories they would be able to recall and the number actually recalled by another group of subjects. These authors think that availability can be used to make frequency judgments, especially when individual instances are not distinct, are too numerous to permit counting, or occur at a rate that is too fast for accurate encoding. They caution, however, that availability can be affected by factors such as recency of occurrence and salience, which will result in imprecise judgments of frequency. Several of the experimental findings that have been reviewed here may have resulted from the use of the availability heuristic, for example, Lichtenstein et al.'s (1978) finding of overestimates due to salience of the event to the respondent and all of the other judgments of naturalistic events (e.g., Attneave, 1953; Erlich, 1964; Shapiro, 1969).

The research on frequency judgments can make important contributions to research on errors in surveys, since many surveys ask respondents to report on the frequency of some past behavior. We need to know more about the conditions in survey reporting under which respondents use the two strategies for making frequency judgments and more about what might affect their use of the availability heuristic. Given that people appear to be better at making relative than absolute judgments of frequency, we need to know more about the limits of accuracy in making absolute judgments and the types of reporting where it would be better to be content with relative judgments than to try to get absolute, but inaccurate reports of frequencies.

### Conclusion

In this paper we have reviewed the kinds of problems that are of concern to those working in the field of response effects in surveys and related them to some of the pertinent theories and findings of cognitive research.

A number of the factors shown to have effects in the experimental literature may not have analogous results outside of the laboratory. In the case of a factor such as the similarity of contexts at the time of encoding, the relevant data on what constitutes similarity are not available. Similarly, the effect of spacing on frequency reports may be eliminated when interval between repetition is measured in terms of days rather than seconds and when the number of intervening “events” may be in the hundreds rather than less than ten. The same general critique may be applied to the factors of type of processing given to an event (e.g., Rowe, 1974) and to thinking about an event (e.g., Johnson et al., 1977). The argument here is not that these experimentally identified factors do not have implications for the design of surveys, but rather that the application to events outside of the laboratory is not direct and requires additional research.



The purpose of the paper is to structure a discussion among survey researchers, cognitive scientists, and statisticians to explore the ways in which work in the different fields can be brought together to enrich our understanding of errors in information processing, be they in the survey context, in other contexts of interest to researchers, or in ordinary discourse. From this discussion we hope that an agenda for future research will emerge.

### References

- Alba, J.W., Chromiak, W., Hasher, L., and Attig, M.S. 1980 Automatic encoding of category size information. Journal of Experimental Psychology: Human Learning and Memory 6:370-378.
- Anderson, J.R., and Bower, G.R. 1972 Recognition and recall processes in free recall. Psychological Review 79(2):97-123.
- Attig, M., and Hasher, L. 1980 The processing of occurrence information by adults. Journal of Gerontology 35:66-69
- Attneave, F. 1953 Psychological probability as a function of experienced frequency. Journal of Experimental Psychology 46:81-86.
- Bartlett, F.C. 1932 Remembering: A Study in Experimental and Social Psychology. Cambridge, England: University Press.
- Belson, W.A. 1968 Respondent understanding of survey questions. Polls 3:1-13.
- Bradburn, N.M. 1983 Response effects. In P.H. Rossi and J.D. Wright, eds., The Handbook of Survey Research. New York: Academic Press.
- Bradburn, N.M., Sudman, S., and associates 1979 Improving Interviewing Method and Questionnaire Design. San Francisco: Jossey-Bass.
- Bransford, J.D., and Johnson M. 1972 Contextual prerequisites for understanding: some investigations of comprehension and recall. Journal of Verbal Learning and Verbal Behavior 11:717-726.
- Brown, N., Rips, L.J., and Shevell, S.K. 1983 Temporal Judgments About Public Events. Unpublished manuscript, University of Chicago.
- Burke, D.M., and Light, L.L. 1981 Memory and aging: the role of retrieval processes. Psychological Bulletin 90(2):513-546.
- Crowder, R.G. 1976 Principles of Learning and Memory. Hillsdale, N.J.: Erlbaum.
- Erllich, D.E. 1964 Absolute judgments of discrete quantities randomly distributed over time. Journal of Experimental Psychology 67:475-482.

- Fee, J. 1979 Symbols and Attitudes: How People Think About Politics. Unpublished doctoral dissertation, University of Chicago.
- Hasher, L., and Chromiak, W. 1977 The processing of frequency information. *Journal of Verbal Learning and Behavior* 16:173-184.
- Hasher, L., and Zacks, R.T. 1979 Automatic and effortful processes in memory. *Journal of Experimental Psychology: General* 108:356-388.
- Hintzman, D. 1969 Apparent frequency as a function of frequency and spacing of repetitions. *Journal of Experimental Psychology* 80:139-145.
- Hintzman, D.L., and Block, R.A. 1971 Repetition and memory: evidence for a multiple-trace hypothesis. *Journal of Experimental Psychology* 88(3):297-306.
- Hintzman, D., and Stern, L.D. 1978 Contextual variability and memory for frequency. *Journal of Experimental Psychology: Human Learning and Memory* 4:539-549.
- Howell, W.C. 1973 Storage of event frequencies: a comparison of two paradigms in memory. *Journal of Experimental Psychology* 98:260-263.
- Jacoby, L.L. 1972 Context effects on frequency judgments of word sentences. *Journal of Experimental Psychology* 94:255-260.
- Johnson, J.K., Taylor, T.H., and Raye, C.L. 1977 Fact and fantasy: the effects of internally generated events on the apparent frequency of externally generated events. *Memory and Cognition* 5:116-122.
- Johnson, M.K., Raye, C.L., Hasher, L., and Chromiak, W. 1979a Are there developmental differences in reality monitoring? *Journal of Experimental Child Psychology* 27:120-128.
- Johnson, M.K., Raye, C.L., Wang, A.Y., and Taylor, T.H. 1979b Fact and fantasy: the roles of accuracy and variability in confusing imaginations with perceptual experiences. *Journal of Experimental Psychology: Human Learning and Memory* 5:229-240.
- Kausler, D.H., Wright, R.E., and Hakami, M.K. 1981 Variation in task complexity and adult age differences in frequency of occurrence judgments. *Bulletin of the Psychonomic Society* 18:195-197.
- Kintsch, W., and van Dijk, T.A. 1978 Toward a model of text comprehension and production. *Psychological Review* 85:363-394.
- Lichtenstein, S., Slovic, P., Fischhoff, B., Layman, M., and Combs, B. 1978 Judged frequency of lethal events. *Journal of Experimental Psychology: Human Learning and Memory* 4:551-578.
- Linton, M. 1982 Transformations of memory in everyday life. In U. Neisser, ed., *Memory Observed*. San Francisco: W.H. Freeman and Co.

- Loftus, E.F. 1977 Shifting color memory. *Memory and Cognition* 5:696-699.
- 1982 Interrogating eyewitnesses--good questions and bad. Chapter 4 in R. Hogarth, ed., *Question Framing and Response Consistency: New Directions for Methodology of Social and Behavioral Sciences*, Vol. 4. San Francisco: Jossey-Bass.
- Loftus, E.F., and Marburger, W. 1983 Since the eruption of Mt. St. Helens, did anyone beat you up? Improving the accuracy of retrospective reports with landmark events. *Memory and Cognition* 11:114-120.
- Loftus, E.F., and Palmer, J.C. 1974 Reconstruction of automobile destruction: an example of the interaction between language and memory. *Journal of Verbal Learning and Verbal Behavior* 13:585-589.
- Marquis, K.H., Cannell, C.F., and Laurent, A. 1972 Reporting for health events in household interviews: of reinforcement, question length and reinterviews. *Vital and Health Statistics*. National Center for Health Statistics, Pub. 1000, Series 2, No. 45. Washington, D.C.: U.S. Government Printing Office.
- Murray, J.R., Minor, M.J., Cotterman, R.F., and Bradburn, N.M. 1974 *The Impact of the 1973-1974 Oil Embargo on the American Household*. Report No. 126. Chicago: National Opinion Research Center.
- Neter, J., and Waksberg, J. 1964 A study of response errors in expenditures data from household surveys. *Journal of the American Statistical Association* 59:18-55.
- Noelle-Neumann, E. 1970 Wanted, rules for wording structured questionnaires. *Public Opinion Quarterly* 34:191-201.
- Orne, M.T. 1969 Demand characteristics and the concept of quasi-controls. Pp. 143-179 in R. Rosenthal and R.L. Rosnow, eds., *Artifact in Behavioral Research*. New York: Academic Press.
- Payne, S.L. 1951 *The Art of Asking Questions*. Princeton: Princeton University Press.
- Pinchert, J.W., and Anderson, R.C. 1977 Taking different perspectives on a story. *Journal of Educational Psychology* 69:309-315
- Rose, R.J. 1980 Encoding variability, levels of processing, and the effects of spacing of repetitions upon judgments of frequency. *Memory and Cognition* 8:84-93.
- Rowe, E.J. 1974 Depth of processing in a frequency judgment task. *Journal of Verbal Learning and Verbal Behavior* 13:638-643.
- Rugg, O. 1941 Experiments in wording questions, II. *Public Opinion Quarterly* 5:91-92.

- Schuman, H., and Presser, S. 1981 Questions and Answers in Attitude Surveys: Experiments on Question Form, Wording and Content. New York: Academic Press.
- Shapiro, B.J. 1969 The subjective estimation of relative word frequency. Journal of Verbal Learning and Verbal Behavior 8:248-251.
- Smith, T. 1980 Situational Qualifications to Generalized Absolutes: An Analysis of "Appraisal of Hitting" Questions and the General Social Survey. Technical Report No. 21. Chicago: National Opinion Research Center.
- Sudman, S., and Bradburn, N.M. 1974 Response Effects in Surveys: A Review and Synthesis. Chicago: Aldine.
- Tulving, E. 1972 Episodic and semantic memory. In E. Tulving and W. Donaldson, eds., Organization and Memory. New York: Academic Press.
- Tulving, E., and Thomson, D.M. 1973 Encoding specificity and retrieval processes in episodic memory. Psychological Review 80(5):352-373.
- Tulving, E., and Wiseman, S. 1976 Encoding specificity: relation between recall superiority and recognition failure. Journal of Experimental Psychology: Human Learning and Memory 2:349-361.
- Tversky, A., and Kahneman, D. 1973 Availability: a heuristic for judging frequency and probability. Cognitive Psychology 5:207-232.
- Underwood, B.J. 1969 Some correlates of item repetition in free recall learning. Journal of Verbal Learning and Verbal Behavior 8:83-94.
- Vogel, F.A. 1974 How Questionnaire Design May Affect Survey Data: Wyoming Study. Statistical Reporting Service, U.S. Department of Agriculture, Washington, D.C., Mimeo.
- Warren, L.R., and Mitchell, S.A. 1980 Age differences in judging the frequency of events. Developmental Psychology 16:116-120.

## RECORD CHECKS FOR SAMPLE SURVEYS

Kent Marquis

We need to include validity checks in research on survey cognitive processes for at least two reasons: (1) to help define the measurement problems needing explanation and (2) to know when proposed solutions have improved the answers that respondents give.

Record checks are one kind of validity check. They are a form of “criterion validity” in which survey interview responses are compared with criterion or “true” values on a case-by-case basis.

But the appropriate design of record checks is not so straightforward as the definition might imply. If the design is incomplete, as often happens in the field of health survey methods, we risk drawing the wrong conclusions about the nature of reporting problems, and we risk implementing measurement solutions that may actually do more harm than good.

In this paper I first present the generic record-check designs and their response bias estimators. Then, in the second section, I work through some numerical examples to illustrate the design and interpretation problems of incomplete record checks. From this we will learn that the full record-check design is robust to several important kinds of design and process mistakes that can occur. In the third section I view the empirical results of health survey record checks on the reporting of hospital stays and physician visits. These results tend to confirm the existence of the design-induced interpretation problems that were illustrated in the previous section. These studies also suggest that the net response biases for survey estimates of health service use are close to zero. Some implications of the zero net survey bias are mentioned in the last section.

### Basic Logic of Record Checks

In this section we view the underlying logic of record-check designs and the inferences about survey response bias. For didactic purposes I restrict our attention to the binary variable case: we will learn about three “pure” record-check design possibilities and the bias estimators that go with each of the designs.

Table 1 shows the cross-classified record-check outcome possibilities for a dichotomous variable such as whether a person was hospitalized, visited a physician, or has a particular chronic health condition. The record outcomes (yes, no) are arrayed along the top; the survey outcomes are arrayed along the left side. The cells contain the cross-classified values reflecting both the survey and record observation outcome. On the agreement diagonal, the A cell contains the observations for which both the record and the survey indicate “yes” while the D cell contains agreements about the absence of the characteristic. The off-diagonal cells, B and C, contain the disagreements or errors: it is the

misestimation of these cells (e.g., their frequency or probability) that leads to the common misinterpretation problems.

TABLE 1 Basic Record-Check Matrix for Binary Variables With No Missing Data

Survey	Record		Total
	Yes	No	
Yes	A	B	A + B
No	C	D	
Total	A + C		A + B + C + D

Net survey bias is the discrepancy between the survey estimate of the population characteristic and the true value of the characteristic. In the framework used here, the survey estimate is  $(A + B)/(A + B + C + D)$ , and the true value is  $(A + C)/(A + B + C + D)$ . This discussion deliberately omits the possibilities of nonresponse, processing, and sampling biases, so we may interpret the discrepancy as the response bias.

There are three different ways of conducting a record check, which I call the AC design, the AB design, and the full design. Collectively I refer to the AC and AB designs as partial or incomplete designs.

In the pure AC design, which has been used widely in health survey studies, we first go to the records to select cases with the characteristic of interest present. This might be a sample of people with hospital admissions, people with doctor visits, or people with a chronic condition such as one of the heart diseases. We then interview these people to see if they report the characteristic of interest. Our estimate of survey bias is  $C/(A + C)$ , the underreporting rate. Note that we ignore the B errors.

The pure AB design is a different approach. Here we conduct the survey first and then check records for people who report the presence of the characteristic of interest. If Ms. Smith reports a visit to Dr. Brown, for example, we write to Dr. Brown and ask him to confirm or deny the validity of Ms. Smith's survey report. We estimate the survey bias as  $B/(A + B)$ , which we call the overreporting rate. Note the absence of C-cell information.

We can arrange full designs in several ways. In health survey studies the most common approach is to identify a population of interest, sample from it independently of record or survey values of the characteristic(s) of interest, obtain both survey and record information for each sampled element, and compare the two information sources. The principal feature of the full design is its ability to obtain unbiased

About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.

estimates of at least cells A, B, and C. The full design estimate of survey response bias is  $(B - C)/(A + B + C + D)$ .

### Illustrations of Record-Check Inferences About Net Survey Response Bias

I use numerical examples next to illustrate the interpretive pitfalls inherent in estimates that come from incomplete record-check designs. In an earlier paper (Marquis, 1978a) I provide formal derivations of the conclusions that the AB and AC designs (1) can overestimate the size of the response bias and (2) can misestimate the direction of the bias (e.g., inferring forgetting when the predominant errors are false positive responses).

Using whether hospitalized as the variable of interest and referring to [Table 1](#), we get entries on the agreement diagonal, cells A and D, when the survey and record are either both yes (both report a hospitalization) or both no (both deny a hospitalization).

For convenience, let's assume that the record is always correct<sup>1</sup> and that each person in the population either was or was not hospitalized and reports either yes or no. Then the true rate of persons hospitalized is  $(A + C)/(A + B + C + D)$ , and the survey estimate is  $(A + B)/(A + B + C + D)$ . The net survey or response bias is the difference between the true rate and the survey estimate.

#### Illustration 1: Bias Estimates From Partial Designs are Too Large Because the Estimators Use the Wrong Denominator

One of two problems in the partial design bias estimators is that they use a denominator that excludes some relevant information. This results in an estimate of the response bias that is too large.<sup>2</sup>

For example, assume we have a population of 100 people and that records show that 60 of them stayed in the hospital at least once within the last 12 months. If interviewed, assume 10 of these 60 people would fail to report being hospitalized. Also, if we interviewed the 40 people who were not hospitalized, assume they would all report this fact correctly. The cross-classified observations for this example are in [Table 2](#).

<sup>1</sup>The derivations in Marquis (1978a) address the effects of record errors on the estimates of survey bias.

<sup>2</sup>The partial design estimators are conditional error rates. They are “wrong” only in the sense that they are often misinterpreted as estimates of net bias. The other problem with a partial design estimate is with the numerator and this is discussed later.



TABLE 2 Hypothetical Data Illustrating Denominator Problem

Survey	Truth or Perfect Record		Total
	Yes	No	
Yes	50	0	50
No	10	40	50
Total	60	40	100

$$\text{The true survey bias is } \frac{(A + B) - (A + C)}{(A + B + C + D)} = \frac{50 - 60}{100} = -.10 ;$$

$$\text{the AC design estimate is } - \frac{C}{A + C} = - \frac{10}{60} = -.17 .$$

The correct denominator for an AC design is  $\frac{C}{A + C}$ , but if the denominator  $\frac{C}{A + B + C + D}$  is used, as in this example, it provides the desired estimate of the survey bias,  $-\frac{10}{100} = -.10$ .

### Illustration 2: Estimates From Partial Designs Can Yield the Wrong Sign for the Survey Bias

Because an incomplete design makes it impossible to observe all the survey errors, the bias estimate could have the wrong sign.

For example, assume we have a population of 100 people, 60 of whom stayed in the hospital at least once within the past 12 months. If we interview the hospitalized people, assume 55 of them would report this correctly. If we interview the 40 people who weren't hospitalized, assume 10 would report that they were hospitalized (possibly because they "telescoped" hospital stays occurring more than 12 months ago into the 12-month reference period). The cross-classified observations for this example are in [Table 3](#).

TABLE 3 Hypothetical Data Illustrating the Misestimation of the Sign of the Bias

Survey	Truth or Perfect Record		Total
	Yes	No	
Yes	55	10	65
No	5	30	35
Total	60	40	100

The true survey bias is  $\frac{65-60}{100} = .05$  ;

The true survey bias is

The AC design estimate is  $\frac{5}{60} = .08$  ;

the AC design estimate with correct denominator is  $\frac{5}{100} = .05$  .

The AC design estimates,  $.08$  and  $.05$ , have the wrong sign. Although the true net survey bias is positive, the AC design provides estimates that suggest that the bias is negative.

The reader may wish to create another example to show that the AB design will provide a positive response bias estimate when the true survey bias is negative [viz.,  $(A + B) < (A + C)$ ]. The results will be clearer if you make the B value greater than zero.

### Illustration 3: Partial Design Estimates Overstate the Size of the Net Survey Bias in the Presence of Match Errors

Let us now introduce match errors, that is, mistakes in cross-classifying the survey and record values. Match errors occur either (1) because one of the sources gives incorrect information about the details of an event that precludes a correct match or (2) because someone makes a mistake in carrying out the matching procedures.

For this illustration we assume that all survey and record information is matched, that is, there are no “left over” interviews or records when we finish the matching operation. With this assumption, each time we make a match error, we cause a “compensating” match error to occur somewhere else. The illustration to follow shows that, under the “no nonmatch” assumption, the full design estimate of net response bias is unaffected by match errors (see Neter et al., 1965, for the proof), but the partial design estimates are biased when match errors are present.

It is common practice in record checks to require that the survey and record information agree (within tolerance limits) on a number of dimensions other than the characteristic of interest before an entry is made on the agreement diagonal for the characteristic of interest. In our hospital stay record check, for example, procedures might require agreement on the patient's name, address, sex, the name of the hospital, and possibly the date or length of stay within some range.

An attribute reporting error is a response mistake about one of these match variables that prevents a (correct) cross-classification of the survey and record information.

For example, suppose Mrs. Smith mistakenly reports that her son Tommy was in the hospital to have his tonsils removed when it was really her daughter Suzy who had that experience. This one mistake about the name attribute would generate two classification errors: Tommy's event would show up as a count in the B cell, and Suzy's event would show up as a count in the C cell. Note that Mrs. Smith correctly reported the number of hospital stays for her family, but our special match requirements for the record check cause the report to appear as two offsetting errors. Similarly, a careless clerk could mismatch Tommy's and Suzy's correct interview reports to the hospital information, generating entries in cells B and C instead of in the A and D cells where they belong.

To show the effects of offsetting match errors, let's interview all or part of another population of 100 people and insert only offsetting match errors into the cross-classified observations, as in Table 4, in which 20 match errors have generated 40 off-diagonal entries.

TABLE 4 Hypothetical Data Illustrating Match Errors

Survey	Truth or Perfect Record		Total
	Yes	No	
Yes	50	20	70
No	20	10	30
Total	70	30	100

Both the true bias and the full design bias estimate are zero. The adjusted (using the correct denominator) AB design estimate is .20; the adjusted AC design estimate is .20.

Both the adjusted AB and AC approaches overestimate the (absolute value of the) survey bias by including half of the unsystematic (offsetting) errors in their estimate of the systematic survey bias.

To generalize a little further, any random (offsetting) mistakes made by respondents, interviewers, and processing clerks that do not change the (expected value of the) subject-matter estimate from the survey cause the AB design inference of a systematic positive response bias and cause the AC design inference of a systematic negative (i.e., forgetting) bias. Since random mistakes probably are inevitable, the continued use of incomplete record-check designs will continue to mislead us about the direction and size of net survey response biases.

**Illustration 4: The Combined Effects**

As a final example, I have generated a matrix reflecting the various kinds of systematic and unsystematic errors discussed above to show how full and partial record-check design estimators handle them. Table 5 is based on the following assumptions: a population of 100 people; 60 people were hospitalized; 10 of the 60 would fail to report the hospitalization (i.e., forgetting) if interviewed; 40 people were not hospitalized; 5 of the 40 would falsely report a stay (e.g., telescoping) if interviewed; and 15 match mistakes that generate 30 cross-classification errors.

TABLE 5 Hypothetical Data Illustrating the Combined Effects of Errors

Survey	Truth or Perfect Record		Total
	Yes	No	
Yes	50 true positive reports	5 false positive errors	55
	15 match errors	+15 match errors	
	35 A-cell entries	20 B-cell entries	
No	10 false negative errors	35 true negative reports	45
	+15 match errors	15 match errors	
	25 C-cell entries	20 D-cell entries	
Total	60	40	100

$$\text{The true survey bias is } \frac{(A + B) - (A + C)}{A + B + C + D} = \frac{55}{100} - \frac{60}{100} = -.05 ;$$

$$\text{the full design estimate is } \frac{B - C}{A + B + C + D} = \frac{20 - 25}{100} = -.05 ;$$

$$\text{the AC design estimate is } -\frac{C}{A + C} = -\frac{25}{60} = -.42 ;$$

About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.

$$\text{the adjusted AC design estimate is } - \frac{C}{A + B + C + D} = -.25 ;$$

$$\text{the AB design estimate } = + \frac{B}{A + B} = \frac{20}{55} = .36 ;$$

$$\text{the adjusted AB design estimate is } \frac{B}{A + B + C + D} = .20 .$$

The AC estimate has the correct sign (minus) but is too large. The adjusted AC estimate shrinks in the right direction (toward  $\cdot 05$ ) but not far enough. The adjusted estimate ( $\cdot 25$ ) reflects the 10 false negative survey errors and half (15) of the 30 match-error entries. It, of course, completely ignores the 5 false positive errors and the other half of the match errors.

In this example the AB estimate (+.36) has the wrong sign; the unadjusted estimate is too large; and the adjustment shrinks it in the right direction (toward  $\cdot 05$ ) but not far enough. The adjusted AB estimate (+.20) reflects the five telescoping errors and half of the match errors but is unaffected by the presence of 10 forgetting errors and the “compensating” match errors.

The conclusions I draw from the above are in the nature of cautions. In defining the survey problems to approach using cognitive theories and methods, do not assume that forgetting dominates survey reporting error. Much of the evidence for forgetting comes from AC design record checks whose design and bias estimator guarantee large negative values even in the presence of zero bias or a positive net bias.

The corollary caution concerns future laboratory studies one might undertake. If these involve a criterion validity component, the full design principles are just as applicable in the laboratory as in the field.

### Record-Check Estimates of Reporting Errors in Health Surveys

How well do the logic-based principles of record-check problems hold up in practice?

In this section we look at actual record-check estimates of survey reporting bias for hospital stays and physician visits. We will see that the full design studies tend to estimate very small net reporting biases (close to zero), while the incomplete AC designs produce a negative bias estimate (e.g., forgetting), and the incomplete AB designs yield positive estimates of response bias.

The effect of the incomplete design approaches on the size and sign

of the response bias is as predicted: the idea of a dominant forgetting bias appears to be a methodological artifact.<sup>3</sup>

I show estimates of hospital stay reporting errors first and then estimates of physician visit reporting errors.

### Hospital Stay Reporting Biases

The hospital stay bias estimates from the three kinds of record-check designs (Table 6) show the expected effects with respect to the direction of the response errors. The AC designs imply a net omission (i.e., forgetting) bias and the AB designs imply a net positive (i.e., telescoping) response bias. The two ABC design studies come close to full design procedures; both find approximately equal rates of underreporting and overreporting; one estimates a slightly higher relative underreporting rate while the other estimates a slightly higher overreporting rate.

Although a full design estimate isn't possible (neither study provides enough information to estimate the D-cell), net hospital stay response bias is probably very small across surveys, and most of the incomplete designs miss this conclusion entirely by concluding the existence of a large directional bias.<sup>4</sup>

Some of the studies that used incomplete designs obtained information about the “missing” cell; for example, Cannell et al. (1967) found some apparent overreports in their AC design and Andersen and Anderson found some underreports in their AB design. This can happen if the person reports (or the record contains) more than one hospital stay for the person and the two sources disagree about the extra stay(s). But multiple stays are atypical events and not a good basis for inferring a value for the missing cell. Occasionally researchers have pointed out the low frequency in the “missing” cell of incomplete designs and offered this as empirical evidence for assuming that those kinds of errors are infrequent enough to ignore. But readers who have stayed with me this far can see the potential fallacy in that argument.

The modified AB designs, which were used in the studies at the bottom of Table 6, unsuccessfully tried to turn an AB design into a full design. The modification was to check the records of all hospitals that each family reported using for all members of the family. For example, if Mr. Jones said he stayed at Metropolitan Hospital,

---

<sup>3</sup>The full design estimates suggest the absence of substantial net reporting biases. This is not something that can be derived from the record check design and estimation characteristics. It is “new” information that I speculate about in the final section.

<sup>4</sup>A reviewer suggests that I caution readers not to generalize this result to response errors in reporting details of hospital admissions. Methods and estimates of the latter kind are given in Marquis (1980, Section III). For example, using “A-cell” cases for which both survey and record information was available, out-of-pocket cost net bias was close to zero while reports of length-of-stay and month of admission showed a small net positive bias.

TABLE 6 Estimated Rates of Overreporting and Underreporting of Hospital Stays in Household Surveys, by Type of Record-Check Design

	Rate of Underreporting (percent)	Rate of Overreporting (percent)
ABC Designs		
Belloc (1954)	` 14	11
Feather (1972)	` 11	14
AC Designs		
Balamuth (1961)--MED-10 Study	` 13	--
Cannell et al. (1961)	` 13	3
Cannell and Fowler (1963)--Procedure A	` 17	4
Kirchner et al. (1969)--Hospital Check	` 13	--
AB Designs		
Andersen and Anderson (1967)	` 1	10
Balamuth (1961)--Prospective Check	--	3
Barlow et al. (1960)--Blue Cross Check	` 5	13
Barlow et al. (1960)--Hospital Record Check	` 1	4-8
Kirchner et al. (1969)--Prospective Check	--	15
Loewenstein (1969)--Prospective Hospital Record Check	--	12
AB Designs (Modified)		
Anderson et al. (1979)	` 7	10
Marquis (1980)	` 7	26

Note: Most of the data in this table are from Marquis (1978b), which contains a discussion of each of the studies except those listed under "AB Designs (Modified)."

About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.



the researchers also asked Metropolitan Hospital to check its records for stays of Mrs. Jones and for each of the children. While the modified procedure does yield more C-cell entries, the C-cell estimate is still not an unbiased estimate of this population parameter. The resulting estimate of the net response bias remains unsatisfactory.

### Physician Visit Reporting Biases

Record-check estimates of errors in reporting physician visits also illustrate the design biases of the incomplete approaches. These data illustrate two other things: (1) the effect of a bias in the record information (it can lead to an interpretation of a reporting bias with the opposite sign) and (2) the length of the survey reference (recall) period apparently does not have the systematic effect on omissions that a memory decay theory would predict (the omission bias does not seem to increase as the length of the average recall period gets larger). These three points are discussed next using the estimates in Table 7.

The incomplete designs yield physician visit reporting bias estimates with the predicted signs. The three studies using AC designs (asking about doctor visits known from records) produce large underreporting bias estimates and suggest that either about a quarter, a third, or a half of known visits are not reported in surveys. On the other hand, the two AB design studies yield positive (overreporting) estimates of the response bias. Within each of the ABCD design studies, the estimates of over- and underreporting rates are approximately equal, raising the possibility that the net reporting bias is not significantly different from zero.<sup>5</sup>

The ABC design estimates from Feather (1972) show how record bias can distort estimates of the survey reporting bias, even when using something close to a full design. Her overreport estimate (46 percent) is substantially larger than her underreport estimate (14 percent), suggesting that, on the average, her respondents reported a lot of physician visits that actually did not occur within the 2-week reference period of the survey. Feather's records of doctor visits were summary "fee submissions," which are bills for complete outpatient services, each of which may represent more than one office visit (e.g., a visit for the complaint, a visit for diagnostic testing, several visits for treatment, and possibly one or two follow-up visits to monitor the treatment's effectiveness). Thus the records "underreport" visits. When they are cross-classified with the survey reports, there are a large number of B-cell entries (respondent reports of visits that cannot be matched to record reports of visits). The effect is to inflate the survey overreporting estimate,  $B/(A + B)$ . Had Feather calculated a full design estimate,  $(B - C)/(A + B + C + D)$ , it would have had a large positive value also. This is an illustration of the general principle that record biases can show up as response biases with the opposite sign.

---

<sup>5</sup>Indeed, Marquis et al. (1979) used the information in these studies to provide full design estimates of the net response bias. These estimates are not significantly different from zero.

TABLE 7 Doctor Visits: Estimated Rates of Overreporting and Underreporting, by Type of Record-Check Design

Design Type	Study	Visit Measure	Reference Period	Rate of Underreporting (percent)	Rate of Overreporting (percent)
ABCD (full design)	Loewenstein (1969, Prepaid Health)	Any Visit	12 months	11	14
	Madow (1973)	Any visit	12 months	6	5
	Cartwright (1963)	Any visit	4 weeks	3	4
ABC	Sudman et al. (1974) Combined interview and diary	Number	3 months	17	24
	Feather (1972)	Number	2 weeks	14	46 <sup>a</sup>
AC	Loewenstein (1969) Health Dept. records	Any visit	7 months	52	
	Cannell and Fowler (1963)	Number	2 weeks	23	
AB	Balamuth et al. (1961)	Any visit	2 weeks	36	
	Andersen et al. (1975)	Number	12 months		13
	Loewenstein (1969)	Any visit	12 months		14

Note: Adapted from Marquis et al. (1979).

<sup>a</sup>Records contain a large negative bias (unit is multivisit fee submission).

The last point about [Table 7](#) concerns memory decay forgetting bias. The studies use different reference period lengths, ranging from 2 weeks to 12 months. (They also use different visit measures but the visit measure and the reference period length are not completely confounded.) If the memory decay hypothesis is correct, we should observe higher rates of omission (underreports) in studies that use longer recall (reference) periods. But the underreporting rates are not always larger for the longer reference periods. For example, restricting our attention to the “any visit” measure, the 36 percent underreport estimate for 2 weeks in the Balamuth et al. (1965) study is larger than the 11 percent underreport estimate for 12 months in the Loewenstein (1969) study.<sup>6</sup> I return to the memory decay issue in the next section.

### Reporting Biases for Other Variables

The principles illustrated above also apply to record checks for other “objective” subject matter variables. Of particular interest is a recent summary of “sensitive topic” record checks (Marquis et al., 1981). The topics include receipt of welfare, wage or salary income, illegal drug use, alcohol consumption, arrests and convictions for crimes, and embarrassing chronic conditions (e.g., hemorrhoids and mental illness). Not surprisingly, there have been incompletely designed record checks in these areas, and the conclusions drawn (e.g., people won't report socially undesirable information) are a function of the type of design used. Surprisingly (at least to me), the distribution of bias estimates from the full design record checks center on zero or possibly a small positive value for most of the sensitive topics, suggesting that the errors respondents (or records, or match procedures) make are largely “offsetting” on the average (rather than being mostly in one direction) as would be the case for forgetting or lying about sensitive information.

### Implications for Future Research

In this section I speculate about the implications for research of a net reporting bias close to zero, mention an example of a model that might fit the data, and suggest that perhaps a different conceptualization of offsetting errors is needed for health service use reports.

The empirical findings from fully designed record checks of hospital stay and doctor visit reports suggest that the net response bias is approximately zero, at least for surveys designed as those involved in the cited record checks. I see three implications of these findings.

One implication is that forgetting (and its possible underlying causes, such as failure to acquire the information, failure to retrieve

---

<sup>6</sup>However, within a single reference period length, record checks usually show a positive correlation of underreport probability and elapsed time between the event and the interview. See Cannell et al. (1977) or [Table 3](#) in this paper for examples.

it, or decisions not to report it) is not necessarily the dominant response problem in health surveys as we normally design them. Thus, cognitive research that focuses only on forgetting hospital stays and doctor visits may not have much applied value for contemporary health surveys.

A second implication for cognitive research on health surveys is that it should contain criterion validity features such as fully designed record checks or carefully thought-out strategies of construct validity. Over the past 10-15 years, survey methodologists have sometimes substituted the assumption that “more is better” for empirical validity studies, inferring that a survey procedure that produces more reports of something yields better estimates than a less productive procedure. But the more-the-better assumption is unwarranted when fully designed record checks show that the net response bias for current procedures is approximately zero.

A third implication is that survey methodologists may want to reexamine the effects of recent changes in health survey designs that were incorporated to reduce forgetting biases (see, e.g., Cannell et al., 1977, Appendix). Such changes may be producing a net positive response bias.

A zero net response bias does not mean that survey responses are given without error, only that the errors are offsetting when calculating the survey mean (proportion, or other first-moment statistic). Offsetting errors are of concern to survey designers and analysts because they place limits on the precision of some population estimates and cause biased estimates of other population parameters (e.g., coefficients of association). Cognitive science can make an extremely important contribution to survey design by describing (e.g., modeling) these errors and discovering what causes them (especially if the causes can be influenced by features of the survey design such as the construction of question sequences, the length of recall periods, or the behavior of the interviewer). The important point here is that the errors to be described, understood, and controlled can be offsetting or possibly random. They are not the product of a single cognitive process that produces only omissions. Whatever cognitive processes are operating can produce just as many false positive reports as false negative reports, so our explanations need to expand to take this phenomenon into account.

Sudman and Bradburn (1973) provide an example of a relevant response error model and show that it can make useful predictions in household expenditure surveys. They assume that forgetting and telescoping are two separate cognitive processes that are affected by the length of the recall interval. The proportion of forgetting response errors increases as the recall interval increases while the proportion of telescoping errors decreases with elapsed time. The effects of the two types of errors will be approximately equal, then, for one particular reference period length. So far, the full design record-check results do not contradict this formulation if one assumes that the omission and telescoping tendencies have approximately balanced out to create a net reporting bias close to zero in all of the full design studies examined. In addition, Cannell et al. (1977) cite several AC design studies that show an apparent increase in omissions with an increase in elapsed time.

But is this another methodological artifact or are we, indeed, observing half of the compensating phenomenon suggested by Sudman and Bradburn? If, for example, we could find AB design studies that show overreporting rates inversely proportional to elapsed time, the compensating forgetting and telescoping explanation would receive support in the health survey context. Unfortunately, none of the AB design record-check studies have published elapsed time data.

Feather's ABC design study (Feather, 1972) does not provide elapsed time data for both hospital stay underreporting and overreporting. I have adapted these data in [Table 8](#), using the record version of date of admission (where possible) to make the elapsed time classification. Since date of admission was asked only for the most recent (or only) hospital stay, we must confine the analysis to these stays (somewhere between 75 and 80 percent of all reported and of all recorded stays).

The overreporting and underreporting trends in [Table 8](#) do not support the Sudman and Bradburn model. Both overreporting and underreporting increase with elapsed time. Although not shown, the underreporting trends are similar to those cited by Cannell et al. (1977), suggesting that the reporting dynamics in Feather's research are similar to those operating in the other record-check surveys that ask about hospital stays. Thus, a different error model may be needed for reporting of hospital stays and physician visits.

We do observe, in the Feather data, an increase in both kinds of response errors with greater elapsed times between the survey and the event occurrence. It is this phenomenon (that some survey methodologists label random response error or simple response variance) that cognitive science might address. For example, is this apparent increase in "carelessness" with elapsed time characteristic of reporting of other kinds of events? Is it due to problems recalling the relevant event attributes correctly? Are the reporting problems due to faulty decisions of the respondent about whether to report a recalled event or in the retrieval of incorrect information about the event followed by a logically correct decision about reporting? Understanding these kinds of phenomena is not going to be gained by traditional survey evaluation methods such as record checks. They need the kind of creative hypothesis formulation and carefully controlled laboratory testing that cognitive science can provide.

TABLE 8 Effect of Recall Interval on Hospital Stay Reporting in a Full Design Record-Check Study

Recall Interval: Elapsed Time Between the Discharge and the Survey (weeks)	Number of Most Recent and Only Hospital Episodes (12 months)		
	In Record	In Survey	Percent Overreported
1-18	173	184	1
19-36	157	159	8
37-52	121	123	15
Not Reported	0	7	16

Note: Adapted from Feather (1972). The recall interval classification is determined by the record value unless the survey report could not be matched to a record report.

About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.

## References

- Andersen, R., and Anderson, O. 1967 A Decade of Health Services. Chicago: University of Chicago Press.
- Andersen, R., Kravits, J., and Anderson O., eds. 1975 Equity in Health Services: Empirical Analyses in Social Policy. Cambridge, Mass.: Ballinger.
- Andersen, R., Kasper, J., and Frankel, M. 1979 Total Survey Error. San Francisco: Jossey-Bass.
- Balamuth, E., Shapiro, S., and Densen, P.M. 1961 Health interview responses compared with medical records. Health Statistics, Publication No. 584-D5, Ser. D, No. 5. Washington, D.C.: U.S. Public Health Service.
- Barlow, R., Morgan, J., and Wirick, G. 1960 A study of validity in reporting medical care in Michigan. Pp. 54-65 in Proceedings of the Social Statistics Section. Washington, D.C.: American Statistical Association.
- Belloc, N.B. 1954 Validation of morbidity survey data by comparison with hospital records. Journal of the American Statistical Association 49:832-846.
- Cannell, C.F., and Fowler, F.J. 1963 Comparison of hospitalization reporting in three survey procedures. Health Statistics, Ser. D., No. 8. Washington, D.C.: U.S. Public Health Service. Reprinted in Vital and Health Statistics, Ser. 2, No. 8, July 1965.
- Cannell, C.F., Fisher, G., and Bakker, T. 1961 Reporting of hospitalization in the health interview survey. Health Statistics, Ser. D., No. 4. Washington, D.C.: U.S. Public Health Service. Reprinted in Vital and Health Statistics, Ser. 2, No. 6, July 1965.
- Cannell, C., Marquis, K., and Laurent, A. 1977 A summary of studies of interviewing methodology. Vital and Health Statistics, DHEW Publication No. (HRA) 77-1348, Ser. 2, No. 69. Washington, D.C.: U.S. Government Printing Office.
- Cartwright, A. 1963 Memory errors in morbidity survey. Millbank Memorial Fund Quarterly 41:5-24.
- Feather, J. 1972 A Response/Record Discrepancy Study. University of Saskatchewan, Saskatoon, November 1972. Available from the National Technical Information Service, Springfield, Virginia.
- Kirchner, C., Lerner, R.C., and Clavery, O. 1969 The reported use of medical care sources by low-income inpatients and outpatients. Public Health Reports 84:107-117.



- Lowenstein, R. 1969 Two Approaches to Health Interview Surveys. School of Public Health and Administrative Medicine, Columbia University.
- Madow, W.G. 1973 Net differences in interview data on chronic conditions and information derived from medical records. Vital and Health Statistics. DHEW Publication No. (HRA) 75-1331, Ser. 2, No. 57. Washington, D.C.: U.S. Department of Health, Education and Welfare.
- Marquis, K.H. 1978a Inferring health interview response bias from imperfect record checks. Pp. 265-270 in Proceedings of the the Section on Survey Research Methods. Washington, D.C.: American Statistical Association.
- 1978b Record Check Validity of Survey Responses: A Reassessment of Bias in Reports of Hospitalizations, R-2319-HEW. Santa Monica, Calif.: The Rand Corporation.
- 1980 Hospital Stay Response Error Estimates for the Health Insurance Study's Dayton Baseline Survey, R-2555-HEW. Santa Monica, Calif.: The Rand Corporation.
- Marquis, K., et al. 1979 Appendix B to Summary Report: An Evaluation of Published Measures of Diabetic Self-Care Variables, N-1152-HEW. Santa Monica, Calif.: The Rand Corporation.
- Marquis, K.H., et al. 1981 Response Errors in Sensitive Topic Surveys: Estimates, Effects, and Correction Options, R-2710/2-HHS. Santa Monica, Calif.: The Rand Corporation.
- Neter, J., Maynes, F.S., and Ramanathan, R. 1965 The effect of mismatching on the measurement of response errors. Journal of the American Statistical Association 60:1005-1027.
- Sudman, S., and Bradburn, N.M. 1973 Effects of time and memory factors on responses in surveys. Journal of the American Statistical Association 68:805-815.
- Sudman, S., Wallace, W., and Ferber, R. 1974 The Cost-Effectiveness of Using the Diary as an Instrument for Collecting Health Data in Household Surveys. Survey Research Laboratory, University of Illinois.

About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.

## APPENDIX B

# DESIGNING AND BUILDING THE BRIDGE

The body of this report and [Appendix A](#) cover the substantive aspects of the CASM project. This appendix describes the process: the organization and structure of the June 1983 seminar at St. Michaels, the January 1984 follow-up meeting in Baltimore, and related project activities. It can be thought of as a case study of an effort to foster interdisciplinary collaboration. Readers who have found the products of the seminar to be of some value may want to consider some of the same approaches for similar undertakings.

### PREPARING FOR THE SEMINAR

Plans for the CASM seminar were developed by the chair, Judith Tanur, with assistance from staff of the Committee on National Statistics (CNSTAT), based on a proposal developed earlier by Stephen Fienberg and Miron Straf. Three specific objectives were set for the St. Michaels meeting: (1) to review recent work in the cognitive sciences and its potential applications to survey research; (2) to propose specific research and experimentation that might lead to improvements in the questionnaire and interviewing procedures used in the National Health Interview Survey (NHIS) and other surveys; and (3) to generate ideas for basic research, including research in the cognitive sciences using surveys as vehicles for experimentation and the collection of relevant data.

In planning the seminar, the organizers made three basic decisions. First, because the number of participants would be relatively small--about 20--and they would meet for approximately one week, it was important to find an informal and isolated setting where the participants' full attention could be given to the work of the seminar. A suitable location was found at St. Michaels, Maryland, on the Eastern Shore of the Chesapeake Bay, and the seminar was scheduled to convene there for six days, from June 15 to June 21 (Wednesday through Monday), 1983.

Second, while the scope of the seminar was to be broader than that of the 1980 Bureau of Social Science Research Workshop, experience from the

earlier meeting showed the value of structuring the discussions around specific problems encountered in one or more complex, large-scale surveys. The NHIS, which is conducted for the National Center for Health Statistics (NCHS) by the Census Bureau, was selected for primary focus at St. Michaels for several reasons. First, the subject matter of the NHIS is complex, requiring performance of difficult cognitive tasks by interviewers and respondents and thus providing a rich source of examples for analysis. Second, the error structure of the survey results is reasonably well understood as the result of extensive and continuing methodological research that began just prior to the inception of the NHIS in 1957. Last--and very important--the two agencies responsible for the NHIS, the NCHS and the Census Bureau, expressed a strong interest in the objectives of the seminar and a willingness to assist in familiarizing participants with the NHIS. In addition, the organizers selected two other surveys for less intensive concentration: the General Social Survey, conducted by the National Opinion Research Center, and the National Election Survey, conducted by the Institute for Social Research of the University of Michigan.

Third, it was decided that the participants would be asked not only to attend the seminar, but also to give some of their time to the project both before and after the seminar. Prior to the seminar, all the participants would be asked to read relevant background materials and to be interviewed in the NHIS. Some participants would also be asked to provide background materials for distribution in advance of briefings at the seminar. Following the seminar, all participants would be urged to submit specific research proposals and would be asked to attend a shorter follow-up meeting if the group decided that such a meeting would be valuable.

Preparation for the St. Michaels meeting included--in addition to the obvious (though crucial) aspects such as participant selection, agenda development, and logistics--four special undertakings detailed below.

- (1) The preparation and distribution to the participants of a set of background materials. The complete set of background materials, listed in [Appendix C](#), was mailed to all participants early in May 1983. Of special significance were two papers prepared expressly for the CASM project: "Cognitive Science and Survey Methods" by Roger Tourangeau and "Potential Contributions of Cognitive Sciences to Survey Questionnaire Design" by Norman Bradburn and Catalina Danis (presented in [Appendix A](#)). Both papers are about existing and potential links between the cognitive sciences and survey research, but the authors approached the subject in different ways.

The 16 other background materials fell into several different categories: the relationships between cognitive science and survey research (items 1, 2, and 11); the NHIS content and procedures and related methodological research (items 8-10); other surveys to be discussed at the seminar (items 12-14); important concepts and research in the cognitive sciences (items 3-5); two pertinent tools used in survey research--computer-assisted telephone interviewing (CATI) and validity checks (items 6, 7, and 15); and a list of selected additional readings in cognitive science and survey methods (item 16).

- (2) Arranging for seminar participants to be interviewed in the NHIS. Thanks to excellent cooperation by the Census Bureau, nearly all of the participants were interviewed by regular Census Bureau interviewers for the NHIS. The experience of being respondents in NHIS personal interviews gave the seminar participants a first-hand understanding of the nature and difficulty of the cognitive tasks required of respondents to that survey. Although these interviews were not part of the regular NHIS sample and were not included in the survey estimates, standard interview procedures were used except for a requirement that the seminar participant be one of the respondents to the core questions and that he or she be selected to respond to the supplemental questions on use of alcohol and tobacco. After being reviewed for interviewer errors at the Census Bureau, the completed questionnaires were returned to the participants for their use.
- (3) Arranging for two NHIS interviews to be videotaped for showing at the seminar. Interviews with two volunteer respondents were conducted by experienced Census Bureau interviewers in the respondents' homes and were videotaped by a specialist in documentary videography. The interviewers were paid by the Committee on National Statistics to conduct the interviews; they were not acting as Census Bureau employees. To protect the rights of the interviewers and respondents, informed consent procedures were developed and were reviewed and approved by the National Academy of Sciences' Committee to Review Human Studies, and the videotapes were copyrighted.
- (4) Adaptation of the NHIS questionnaire to a computer-assisted telephone interviewing (CATI) system for demonstration at the seminar. The adaptation of the questionnaire was done by Albert Madansky, one of the seminar participants, using a CATI system that he had developed. The objective of demonstrating a CATI system was to familiarize participants with a new interviewing technique that is likely to be increasingly used in surveys and that may involve significant changes in the cognitive tasks to be performed by interviewers and respondents.

#### THE ST. MICHAELS MEETING

The agenda for the St. Michaels meeting was divided into four principal phases, intended to proceed more or less in sequence with some overlap. The first phase was intended to facilitate getting acquainted. (The first step had been taken before the meeting by circulating curriculum vitae of all participants.) At the initial group session, each participant was asked to describe his or her research interests and put on the table rough ideas for research relevant to the CASM project objectives.

The second phase consisted of background presentations and discussions. These sessions were intended to give all participants a common information base from which they could work together to develop

research proposals. The topics of these discussions were essentially the same as those covered by the background materials; cognitive sciences for survey researchers; survey methods (including CATI) for cognitive scientists; specific surveys, with major focus on the NHIS; and reactions to participant and videotaped interviews. The background presentations and discussion occupied most of the second and third days of the seminar.

The third phase was devoted to working group meetings. In order to allow substantial time for informal, intensive discussions in smaller groups, two sets of three working groups were established. The first set of the three groups covered the major cognitive aspects of survey interviews (comprehension, retrieval, and judgment and response), and the second set covered the principal topics included in the core portion of the NHIS questionnaire (utilization of health services, health conditions, and restricted activity).

Prior to the seminar, participants had been asked to rank their choices for the topics in each set of working groups. Assignments were then made on the basis of several criteria: participant choices (virtually all designations were either first or second choices within each set); formation of heterogeneous groups, with more or less proportional representation of cognitive scientists, survey researchers, agency representatives, and project staff; and minimization of overlap between working groups in the two sets, so that each participant would work with the largest possible number of the other participants in the working group setting. For each working group, the chair appointed a convener and a member of the project staff to serve as rapporteur.

In response to suggestions by several participants during the seminar, a third set of working groups--which came to be known as the brainstorming groups--was organized. These brainstorming groups had two objectives: (1) to allow and encourage all participants to present and get reactions to additional proposals for research and (2) to allow participants to work in a small group setting with some of the seminar participants with whom they had not been associated in either of the other two sets of working groups. The designation of these groups took into account requests from some participants to be in the same group with specific individuals.

The fourth and final phase was feedback and integration. General feedback sessions were scheduled for the evenings of the third and fifth days in order to permit discussion of possible agenda changes and improvements in procedures. Plenary sessions were scheduled toward the close of the seminar for working group reports. In addition to the plenary sessions, a final session was scheduled at which the participants were asked to summarize their thoughts about relevant research that they would like to undertake or participate in.

With minor exceptions, all seminar participants attended the entire seminar. Several other people were invited to be guests at the seminar: they included representatives from the National Science Foundation, the National Center for Health Statistics, the Bureau of the Census, the Committee on National Statistics, and a university (see pp. v-vi). Most of the visitors came to the seminar site on Saturday or Sunday (the fourth and fifth days of the seminar) and stayed through noon on Monday. This allowed them to attend the plenary sessions at which the reports of

the first two sets of working groups were presented and to view one of the videotaped NHIS interviews. In addition, each guest was invited, at a plenary session, to react to the conclusions and research ideas presented by the working groups.

All plenary sessions were tape recorded, and most of them were transcribed after the seminar. Members of the project staff took detailed notes at each session. These materials, plus the written working group reports, form the basis of the proceedings of the CASM project presented in [Chapter 1](#) of this report.

At the final session, participants agreed that a report of their observations, conclusions, and research ideas should be prepared and that the CASM participants should reconvene for about two days in late 1983 or early 1984 to review a draft report and discuss ideas and plans for further research.

### FOLLOW-UP: THE BALTIMORE MEETING

After the St. Michaels seminar, the CASM chair and staff developed the agenda for a two-day meeting in Baltimore in January 1984 and prepared a partial draft of the project report, which was mailed to participants for review about five weeks before that meeting.

More important, initial results from the St. Michaels seminar began to take shape. Several of the participants, working individually or in small groups, began new research activities or developed proposals for research along the lines that had been discussed at the seminar. It soon became clear that part of the Baltimore meeting could be usefully devoted to presentation and discussion of these activities and proposals. Brief write-ups of several of them were included in the draft report sent out before the meeting.

In order to involve other researchers, some steps were taken to begin publicizing the interim results of the CASM project. One of the videotaped NHIS interviews was shown to and discussed with a small group of Census Bureau employees in November 1983. Plans were made for a session on the CASM project at the annual meeting of the American Association for Public Opinion Research in May 1984. The CASM chair, Judith Tanur, in response to a request, organized an invited paper session on cognitive aspects of survey methodology for the annual meeting of the American Statistical Association in August 1984.

All of the regular participants in the St. Michaels seminar returned for the January 1984 meeting in Baltimore. This fact may indicate their perceptions of the value of the CASM project and their commitment to it. Three guests were invited for part of the meeting: two (who had also been guests at St. Michaels) from the National Science Foundation and one from the Social Science Research Council.

The agenda for the Baltimore meeting retained some of the features of the St. Michaels seminar, adapted to the shorter time available. At the start of the meeting, to stimulate further thinking on the interview process, parts of two videotaped interviews from the pretest for the 1984 round of the National Opinion Research Center's General Social Survey were shown and discussed. Unlike the NHIS, the General Social Survey has



many questions on respondent attitudes and perceptions and thus involves cognitive processes different from those required to respond to the basically factual questions of the NHIS.

Part of the time was spent in small group discussions. The participants were again divided into three working groups, using essentially the same criterion of heterogeneity within groups used to form the groups earlier at St. Michaels. In this case, all groups were to discuss the same two topics: (1) further ideas for relevant research and (2) proposals for ways of continuing and expanding collaboration between cognitive scientists and survey researchers. On the final day of the meeting, the working groups reported back to the full group.

A key feature of the Baltimore meeting, to which considerable time was allocated, was the presentation and discussion of the several research activities and ideas that had been developed by individuals and small groups of participants following the St. Michaels seminar. Prior written reports on those activities and ideas were updated, and several of them evoked strong interest and constructive comments from other participants.

The final agenda item for the Baltimore meeting was to “achieve closure” on the CASM project. The participants agreed on the form of this report and discussed means of publicizing the activities and output of CASM so that other researchers might become involved in cross-disciplinary collaboration. Several other kinds of outreach were suggested in these discussions, including:

- publishing relevant articles in journals read by cognitive scientists and presenting papers at meetings of associations to which they belong;
- extending the dissemination process to international journals and conferences, such as the 1985 meeting of the International Association of Survey Statisticians;
- arranging for appropriate journals to publish special issues devoted to relevant themes, e.g., cognitive studies and survey research methods;
- when sufficient results are available from collaborative research studies, holding symposia in university settings to present and discuss them;
- organizing a short course on cognitive aspects of surveys in conjunction with an annual meeting of the American Statistical Association;
- encouraging cognitive scientists interested in surveys as vehicles for research to attend courses in survey methods, e.g., the summer course presented annually at the Survey Research Center at the University of Michigan;
- preparing short annotated bibliographies (1) for survey researchers interested in learning about relevant aspects of the cognitive sciences and (2) for cognitive scientists who want to become familiar with survey research methods;
- asking cognitive scientists to participate in proposing and planning specific investigations in the Census Bureau's planned three-year program of research on telephone survey methodology.

The question of how and by whom such outreach activities might be undertaken led to the broader question of the future of the CASM project in a more formal sense. The consensus of the participants in Baltimore was that CASM had met and perhaps exceeded its original goals: it had generated several promising interdisciplinary research activities and plans, and it had established an informal network of scientists who appreciate the benefits of collaboration between cognitive scientists and survey researchers. The participants agreed that the necessary momentum had been established for the development of collaborative research and that the important thing now was to proceed with the research. The network would continue to function and expand without requiring formal identification or support at this time. Support for specific research plans or, when more research findings are available, for symposia, would, of course, be necessary.

Thus, the formal aspects of the CASM project have been completed. Evaluation of the results, however, will continue as the organizers and participants monitor the progress of relevant research, the application of cognitive research findings in surveys, and the use of surveys as vehicles for cognitive research.

About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.

## APPENDIX C

### BACKGROUND MATERIALS FOR THE SEMINAR

The first 15 items listed in this appendix were mailed to all participants prior to the St. Michaels seminar. The items are listed in the order of the conference sessions for which they provided background information; notes explain their relevance to the session topic. The last item is a reading list the staff provided for the CASM participants.

Session Topic: 1980 Workshop on Applying Cognitive Psychology to Recall Problems of the National Crime Survey

1. Report of the workshop, Jeffrey C. Moore, rapporteur. Although our seminar will have a broader focus and will look at applications largely in health surveys rather than in crime surveys, in a sense we will be picking up where this workshop left off. Several of the seminar participants took part in the 1980 workshop.

Session Topic: Background Paper No. 1

2. Tourangeau, R. (1983) Cognitive Science and Survey Methods. This paper was written specifically for the seminar. Items 3, 4, and 5 were recommended by Tourangeau as further background on areas of cognitive science most relevant to survey design. (The revised version of Tourangeau's paper appears in [Appendix A](#).)
3. Linton, M. (1982) Transformations of memory in everyday life. In U. Neisser, ed., Memory Observed 77-91. The author describes an experiment in memory for everyday events; the experiment (on the author's own personal memories) covers several years.
4. Abelson, R. (1981) The psychological status of the script concept. American Psychologist 36:715-729. Summary, by one of the main proponents of the "script" concept, of the main evidence for the influence of scripts on memory and comprehension.

5. Nisbett, R., and Ross, L. (1980) Human Inference: Strategies and Shortcomings in Social Judgment. Englewood Cliffs, N.J.: Prentice Hall. Chapter 2. An introduction to work on the “representativeness” and “availability” heuristics.

Session Topic: An Introduction to CATI

6. Roshwalb et al. (1979) New Methods of Telephone Interviewing: A & S/CATI. Paper presented at XXIII ESOMAR Congress. Describes the CATI (computer-assisted telephone interviewing) system that will be available for demonstration and use at our conference.
7. Rustemeyer et al. (1978) Computer-Assisted Telephone Interviewing: Design Considerations. Presented at the annual meeting of the American Statistical Association. A general introduction to CATI systems: their advantages, the precise role of computer assistance, the basic elements of a system, and some unresolved issues. A particular system is also described.

Session Topic: Introduction to HIS (Health Interview Survey)

8. National Center for Health Statistics (1982) Current Estimates from the National Health Interview Survey: United States, 1981, Ser. 10, No. 141. A note attached to the cover of this publication identifies the sections that should be of particular interest to seminar participants.
9. Bureau of the Census. National Health Interview Survey Interviewer's Manual (excerpts). The excerpts are: (1) the table of contents for Parts A, D, and E of the manual (other parts are not relevant to the interview), and (2) Part D, Chapter 2, General Instructions for Using the HIS Questionnaires. Most participants should already have the copies of the HIS core questionnaire and supplement that were completed when they were interviewed; additional blank copies will be brought to the meeting. One or two complete copies of Parts A, D, and E of the manual will also be available.
10. National Center for Health Statistics (1977) A Summary of Studies of Interviewing Methodology, Ser. 2, No. 69. Describes the design and results of several methodological studies, most of which were carried out by the Survey Research Center of the University of Michigan for the NCHS in the 1960s. The studies were “. . . designed to test the effectiveness of certain questionnaire design and interviewing techniques used in the collection of data on health events in household interviews and to investigate the role of behaviors, attitudes, perceptions and information levels of both the respondent and the interviewer.”

Session Topic: Background Paper No. 2

11. Bradburn, N., and Danis, C. (1983) Potential Contributions of Cognitive Sciences to Survey Questionnaire Design. This paper was written specifically for the seminar. (The revised version appears in [Appendix A](#).)

Session Topic: Other NCHS Surveys

12. National Center for Health Statistics (1981) Data Systems of the National Center for Health Statistics, Ser. 1, No. 16. See note on cover that identifies relevant parts of this report.

Session Topic: The General Social Survey

13. National Opinion Research Center, University of Chicago (1982) A short description of the survey. We will have available at the seminar a copy of the Cumulative Codebook (1972-1982), which includes all of the survey questions used during this period and marginal totals for each item in each round of the survey.

Session Topic: The National Election Survey

14. Three items are included:
  - a. Institute for Social Research, University of Michigan. American National Election Studies, 1952-1982. A one-page description of the program.
  - b. The questionnaire for the 1982 fall-winter cross-section interview.
  - c. Respondent booklet used in conjunction with the questionnaire.

Session Topic: Validity Checks

15. Marquis, K. (1978) Inferring health interview response bias from imperfect record checks. Proceedings of the Section of Survey Research Methods. American Statistical Association. Discusses some problems in determining "truth" in a survey context. A paper based on a presentation by Marquis at the St. Michaels seminar appears in [Appendix A](#).

### Supplemental Reading List

16. The following list of selected readings in cognitive sciences and survey methods, prepared by the CASM staff, contains items that were considered for mailing to all seminar participants, but had to be omitted in order to stay within reasonable size limits. It was not meant to be comprehensive in any sense; there are more complete listings of relevant items in the reference lists for the two background papers prepared for the seminar.

#### A. Books

- Bradburn, N.M., Sudman, S. and Associates (1979) Improving Interviewing Methods and Questionnaire Design. San Francisco: Jossey-Bass.
- Dijkstra, W., and van der Zouwen, J. (1982) Response Behavior in the Survey Interview. New York: Academic Press.
- Hogarth, R.M., ed. (1982) Question Framing and Response Consistency. New Directions for the Methodology of Social and Behavioral Science, No. 11. San Francisco: Jossey-Bass.
- Moss, L., and Goldstein, H., eds. (1979) The Recall Method in Social Surveys. Studies in Education (new series) 9. Windsor, Ontario: NFER Publishing Company.
- Payne, S.L. (1951) The Art of Asking Questions. Princeton, N.J.: Princeton University Press.
- Schuman, H., and Presser, S. (1981) Questions and Answers in Attitude Surveys: Experiments in Question Form, Wording, and Context. New York: Academic Press.
- Sudman, S., and Bradburn, N.M. (1974) Response Effects in Surveys. Chicago: Aldine.
- Sudman, S., and Bradburn, N.M. (1982) Asking Questions: A Practical Guide to Questionnaire Construction. San Francisco: Jossey-Bass.

#### B. Publications of the National Center for Health Statistics

Series 2, Data Evaluation and Methods Research. The following publications in this series describe various methodological studies related to the National Health Interview Survey.

- No. 6 (1965) Reporting of Hospitalization in the Health Interview Survey.
- No. 7 (1965) Health Interview Responses Compared with Medical Records.
- No. 8 (1965) Comparison of Hospitalization Reporting in Three Survey Procedures.
- No. 16 (1966) Identifying Problem Drinkers in a Household Health Survey.



No. 18 (1966) Interview Responses on Health Insurance Compared with Insurance Records.

No. 23 (1967) Interview Data on Chronic Conditions Compared with Information Derived from Medical Records.

No. 26 (1968) The Influence of Interviewer and Respondent Psychological and Behavioral Variables on the Reporting in Household Interviews.

No. 41 (1971) Effects of Some Experimental Interviewing Techniques on Reporting.

No. 45 (1972) Reporting of Health Events in Household Interviews: Effects of Reinforcement, Question Length, and Reinterviews.

No. 48 (1972) Interviewing Methods in the Health Interview Survey. (Reports on a split-panel test with two substantially different versions of the HIS questionnaire.)

No. 49 (1972) Reporting Health Events in Household Interviews: Effects of an Extensive Questionnaire and a Diary Procedure.

No. 50 (1972) Optimum Recall Period for Reporting Persons Injured in Motor Vehicle Accidents.

No. 54 (1973) Quality Control and Measurement of Nonsampling Error in the Health Interview Survey.

No. 57 (1973) Net Differences in Interview Data on Chronic Conditions and Information Derived from Medical Records.

Series 1, Programs and Collection Procedures. The following publications in the series are relevant to the Health Interview Survey or other NCHS household surveys.

No. 1 (1965) Origin, Program, and Operation of the U.S. National Health Survey. Reprint of earlier publication.

No. 2 (1964) Health Survey Procedure: Concepts, Questionnaire Development, and Definitions in the Health Interview Survey.

No. 11 (1975) Health Interview Survey Procedure: Concepts, Questionnaire Development, and Definitions in the Health Interview Survey.

No. 15 (1981) Plan and Operation of the Second National Health and Nutrition Examination Survey: 1976-80.

Other

National Health Interview Survey: Report of the National Committee on Vital and Health Statistics (1980). This report of a Technical Consultant Panel on the Health Interview Survey includes recommendations for changes in the content of the HIS questionnaire.

### C. Other Published Reports

- Loftus, E. (1982) Memory and its distortions. In A.G. Kraut, ed., The G. Stanley Hall Lecture Series. Washington, D.C.: American Psychological Association. Contains some ideas about how to study memory through survey research.
- Marquis, K. (1978) Record Check Validity of Survey Responses: A Reassessment of Bias in Reports of Hospitalizations. R-2319-HEW. Santa Monica, Calif.: The Rand Corporation.
- Schuman, H., and Kalton, G. (1985) Survey methods and interviewing. Chapter in G. Lindzey and E. Aronson, eds., The Handbook of Social Psychology, 3rd edition. New York: Random House. [Available in prepublication form.]
- Schuman, H., Smith, T., and Turner, C. (1984) Variability in survey measurements of subjective phenomena: empirical results. Chapter 5 in Surveying Subjective Phenomena. Report of the Panel on Survey Measurement of Subjective Phenomena, Committee on National Statistics. New York: Russell Sage Foundation. [Available in prepublication form.]
- Skogan, W. (1981) Issues in the Measurement of Victimization, NCJ-74682. Washington, D.C.: Bureau of Justice Statistics. An extensive overview of 15 years of methodological development in refining the methods by which criminal victimization can be measured through survey interviews.

### D. Journal Articles: Cognitive Sciences

- Bower, G., Black, J., and Turner, T. (1979) Scripts in memory for text. Cognitive Psychology 11:177-220. Provides the main empirical support for script theory.
- Einhorn, H., and Hogarth, R. (1978) Confidence in judgment: persistence of the illusion of validity. Psychological Review 85(5):395-416. Describes the tendency for people to be overconfident in performing a range of judgment tasks.
- Erdelyi, M., and Kleinbard, J. (1978) Has Ebbinghaus decayed with time?: the growth of recall (hypermnnesia) over days. Journal of Experimental Psychology: Human Learning and Memory 4(4):275-289. Sometimes memory improves over time, particularly in the face of repeated efforts at recall.
- Loftus, E., and Beach, L. (1982) Human inference and judgment: is the glass half empty or half full? Stanford Law Review 34:939-956. A review of R. Nisbett and L. Ross (1980) Human Inference: Strategies and Shortcomings of Social Judgment. Englewood Cliffs, N.J.: Prentice-Hall.
- Nisbett, R., and Wilson, T. (1977) Telling more than we can know: verbal reports on mental processes. Psychological Review 84(3):231-259. Describes some limits to self-knowledge and claims that self-reports reflect our theories of behavior more than a direct introspective awareness of its causes.

Smith, E., and Kluegel, J. (1982) Cognitive and social bases of emotional experience: outcome, attribution, and affect. Journal of Personality and Social Psychology 43(6):1129-1141. Reports on a study, based on results from a national survey, of cognition-emotion links, taking into consideration the social context of the individual.

#### E. Journal Articles: Health Interview Survey

- Givens, J., and Moss, A. (1981) Redesigning the National Health Interview Survey's Data Collection Instrument. Paper presented to the American Public Health Association.
- Kovar, M., and Wilson, R. (1976) Perceived health status--how good is proxy reporting? Pp. 495-500 in Proceedings of the Social Statistics Section. Washington, D.C.: American Statistical Association.
- Kovar M., and Wright, R. (1973) An experiment with alternate respondent rules in the National Health Interview Survey. Pp. 311-316 in Proceedings of the Social Statistics Section. Washington, D.C.: American Statistical Association.
- Massey, J., and Gonzalez, J. (1976) Optimum recall periods for estimating accidental injuries in the National Health Interview Survey. Pp. 584-588 in Proceedings of the Social Statistics Section. Washington, D.C.: American Statistical Association.
- Massey, J., Marquis, K., and Tortora, R. (1982) Methodological issues related to telephone surveys by federal agencies. Proceedings of the Social Statistics Section, Washington, D.C.: American Statistical Association.
- Monsees, M., and Massey, J. (1979) Adapting procedures for collecting demographic data in a personal interview to a telephone interview. Proceedings of the Social Statistics Section. Washington, D.C.: American Statistical Association.
- Nisselson, H., and Woolsey, T. (1959) Some problems of the household interview design for the National Health Survey. Journal of the American Statistical Association 34(285):69-87. This article, published about two years after the start of the National Health Interview Survey, discusses many of the basic survey design issues that were addressed in subsequent methodological research.
- White, A., and Massey, J. (1981) Selective reduction of proxy response bias in a household interview survey. Pp. 211-216 in Proceedings of the Social Statistics Section. Washington, D.C.: American Statistical Association.

#### F. Journal Articles: Other Survey Research and Methods

- Kalton, G., and Schuman, H. (1982) The effect of the question on survey responses: a review. The Journal of the Royal Statistical Society, Series A, 145, Part I, 42-73. Includes comments by the discussants.

- Loftus, E. and Marburger, W. (1983) Since the eruption of Mt. St. Helens, did anyone beat you up?: improving the accuracy of retrospective reports with landmark events. Memory and Cognition 11:114-120.
- Sykes, W. (1982) Investigation of the effect of question form. Survey Methods Newsletter. Social and Community Planning Research, London. Includes a classification of question forms used on survey questionnaires.

## APPENDIX D

# BIOGRAPHICAL SKETCHES OF PARTICIPANTS

JUDITH TANUR (Chair) is associate professor of sociology at the State University of New York at Stony Brook. She received a B.S. in psychology from Columbia University in 1957, an M.A. in mathematical statistics from Columbia in 1963, and a Ph.D. in sociology from the State University of New York, Stony Brook, in 1972. She is editor of *Statistics: A Guide to the Unknown* and co-editor of *The International Encyclopedia of Statistics*. Her research has centered on interpersonal control processes, the applications of statistics to social science data, and parallels between the methodology of survey research and that of social experiments. She is a member of the Committee on National Statistics.

NORMAN BRADBURN is the Tiffany and Margaret Blake Distinguished Service Professor at the University of Chicago. He is a member of the Department of Behavioral Sciences, a member of the faculties of the Graduate School of Business and the College, and director of the National Opinion Research Center. He received a B.A. in 1952 from the University of Chicago, a B.A. from Magdalen College, Oxford University, in 1955, and a Ph.D. from Harvard University in 1960. His research interests include the study of psychological well-being and assessments of the quality of life, particularly through the use of large-scale sample surveys; nonsampling errors in sample surveys; and social indicators. He is a member of the Commission on Behavioral and Social Sciences and Education of the National Research Council.

PHILIP E. CONVERSE is professor of political science and sociology at the University of Michigan and director of the Center for Political Studies at the Institute for Social Research. He received a B.A. from Denison University in 1949, an M.A. in English literature from Iowa State University in 1950, and an M.A. in sociology and a Ph.D. in social psychology from the University of Michigan in 1956 and 1958, respectively. He has been engaged for 30 years in large-scale survey research in the United States, Europe, and Latin America in areas that include political behavior, mass/elite relationships, the daily allocation of time use, and perceptions of the quality of life.

ROY D'ANDRADE is professor of anthropology at the University of California, San Diego. He received a B.A. from the University of Connecticut in 1957 and a Ph.D. in 1962 from Harvard University in social relations. He has previously taught at Stanford University and Rutgers University and has been a fellow of the Center for Advanced Study in the Behavioral Sciences. His principal research interests concern American culture and cognitive anthropology.

STEPHEN E. FIENBERG is professor of statistics and social science and head of the Department of Statistics at Carnegie-Mellon University. He received a B.Sc. in mathematics and statistics from the University of Toronto in 1964 and an M.A. and a Ph.D. in statistics from Harvard University in 1965 and 1968, respectively. He previously taught at the University of Chicago and at the University of Minnesota. His principal research has been on the development of statistical methodology, especially in connection with the analysis of cross-classified categorical data, and on the application of statistics in such areas as accounting, criminal justice, ecology, federal statistics, law, medicine and public health, neurophysiology, public policy, and sociology. He has served as chair of the Committee on National Statistics and is currently chair of its Panel on Statistical Assessments as Evidence in the Courts.

ROBERT R. FUCHSBERG is director of the Division of Health Interview Statistics at the National Center for Health Statistics in the U.S. Public Health Service. He received a B.S. in economics from City College of New York in 1949 and also did graduate work at the college in statistics. His interest in health surveys started in 1957 when he joined the National Health Survey Program as an analytical statistician. His main professional interests are the development of improved survey methods and training in survey methodology. During the past 10 years he has developed techniques for adapting telephone procedures to the special needs of health surveys. As a consultant to the government of Portugal since 1980 he has assisted in the development and implementation of the Portuguese National Morbidity Survey.

THOMAS B. JABINE is consultant to the Committee on National Statistics and professorial lecturer at George Washington University. He was formerly statistical policy expert for the Energy Information Administration, chief mathematical statistician for the Social Security Administration, and chief of the Statistical Research Division of the Bureau of the Census. He received a B.S. in mathematics and an M.S. in economics and science from the Massachusetts Institute of Technology in 1949. He has provided technical assistance in sampling and survey methods to several developing countries for the United Nations, the Organization of American States, and the U.S. Agency for International Development. His publications are primarily in the areas of sampling and survey methodology.

WILLETT KEMPTON is research associate in the Family Energy Project at Michigan State University and adjunct assistant professor in the Department of Anthropology. He received a B.A. in 1972 from the University of Virginia in sociology and anthropology and a Ph.D. in 1977

from the University of Texas, Austin, in cognitive anthropology. He received postdoctoral training in quantitative anthropology and public policy at the University of California, where he developed his interest in social and cognitive aspects of energy policy. His publications have dealt with color perception by speakers of Tarahumara and English, variation and change in folk classification systems, and home energy use. His current studies of home energy combine ethnographic, survey interview, and behavioral data.

ELIZABETH LOFTUS is professor of psychology at the University of Washington, Seattle. She received a B.A. in mathematics and psychology from the University of California, Los Angeles, in 1966 and an M.A. and a Ph.D. in psychology from Stanford University in 1967 and 1970, respectively. Her principal research has been in the area of human cognition and memory. She is the author of 10 books, including *Eyewitness Testimony* (Harvard University Press), which won a National Media Award from the American Psychological Foundation in 1980.

ALBERT MADANSKY is professor of business administration and director of the Center for the Management of Public and Nonprofit Enterprise in the Graduate School of Business of the University of Chicago. He received a B.A. in liberal arts in 1952, an M.S. in statistics in 1955, and a Ph.D. in statistics in 1958, all from the University of Chicago. He has been a research mathematician at the RAND Corporation, senior vice president of a large advertising agency, president of a computer software and data processing firm, professor and chairman of computer sciences at City College of New York, and fellow of the Center for Advanced Study in the Behavioral Sciences. Among his research interests is computer-assisted telephone interviewing.

ROBERT MANGOLD is the chief of the Health Surveys Branch of the Bureau of the Census. He received a B.S. from Duquesne University in 1963 and an M.B.A. in 1965 from Ohio University. During his career at the Census Bureau, he has worked on various aspects of large national surveys, such as the National Health Interview Survey, the National Hospital Discharge Survey, the Survey of Neurological Disorders, the Surveys of Veterans, the Current Population Survey, the Point of Purchase Survey, the National Longitudinal Surveys of Work Experience, and the Survey of Criminal Justice Employees.

KENT H. MARQUIS is chief of the Center for Survey Methods Research at the Bureau of the Census. He received a B.A. in psychology from Yale University in 1961 and a Ph.D. in social psychology from the University of Michigan in 1967. He has been a study director at Michigan's Survey Research Center, where his methods research emphasized cognitive and behavioral approaches to survey measures of health, employment, and crime victimization, and he has served as associate director of the Statistical Research Division at the Research Triangle Institute in North Carolina. While a senior social scientist at the RAND Corporation in Santa Monica, he designed and evaluated much of the measurement for a large social experiment concerned with the demand for health insurance. His current interests include research on nonsampling errors in surveys conducted by telephone and other methods.



ANDREW ORTONY is professor of psychology and of education at the University of Illinois at Urbana-Champaign, where he has been since 1973. He has an honors degree in philosophy from the University of Edinburgh, and a Ph.D. in computer science from London University's Imperial College of Science and Technology. His primary areas of research are in the comprehension of language, with special interest in figurative language, and in the relation between affect and cognition. He has edited a widely read interdisciplinary volume, *Metaphor and Thought* (Cambridge University Press, 1979) and with Donald A. Norman is the series editor of *Tutorial Essays in Cognitive Science* (to be published by Erlbaum Associates, Hillsdale, N.J.).

S. JAMES PRESS is a professor and chairman of the Department of Statistics at the University of California, Riverside. He received a B.A. in physics in 1950 from New York University, an M.S. in mathematics from the University of Southern California in 1955, and a Ph.D. in statistics from Stanford University in 1964. He has taught and done research at the University of Chicago, Yale University, the University of British Columbia, the University of California, Los Angeles, the London School of Economics and Political Science, and University College, London. He has worked on statistical problems at Brookhaven National Laboratories, Northrop Aircraft Corporation, Douglas Aircraft Corporation, and the RAND Corporation. His research and numerous publications include work on various aspects of Bayesian inference, subjective probability, and group judgment formulation involving sample surveys.

LEE ROSS is a professor of psychology at Stanford University, where he has taught since 1970. He received a B.A. from the University of Toronto in 1965 and a Ph.D. from Columbia University in 1969. His research interests are in the field of cognitive social psychology, and he has published several books and papers, particularly on strategies and sources of bias in lay inference and judgment.

WILLIAM J. SALTER is a research scientist at Bolt Beranek and Newman, Inc., in Cambridge, Massachusetts. He received a B.S. in mathematics from the University of Chicago in 1969 and an M.S. and a Ph.D. in cognitive psychology from Yale University in 1980 and 1984, respectively. During the 1970s he worked at Abt Associates, Inc., in social program evaluation, policy analysis, and survey design and administration. His current research focuses on causal reasoning and the organization of beliefs, particularly in the domain of economics. His methodological interests center on applying the detailed theories of cognitive science within the context of large samples and subgroup differences in sample survey research.

HOWARD SCHUMAN is director of the Survey Research Center, Institute for Social Research, and professor of sociology at the University of Michigan. He received a B.A. in philosophy from Antioch College in 1953, an M.S. in psychology from Trinity University in 1956, and a Ph.D. in

sociology (social relations) from Harvard University in 1962. Before going to Michigan, he spent three years as a research associate for Harvard's Center for International Affairs, half the time as field director for a survey research project in Bangladesh. His publications have dealt mainly with the nature of survey questions and interviewing, with attitude measurement and the attitude/behavior relation, and with their applications to social issues such as race, economic development, sentiments on war and peace.

MONROE G. SIRKEN is associate director for research and methodology at the National Center for Health Statistics. He received a B.A. and an M.A. from the University of California, Los Angeles, in 1946 and 1947, respectively, and a Ph.D. from the University of Washington in 1950. He has taught at the University of California, Los Angeles, the University of Washington, the University of California, Berkeley, and the University of North Carolina. His research interests and publications are primarily in the fields of sampling and measurement errors in large-scale surveys.

MIRON L. STRAF is research director of the Committee on National Statistics. Previously he taught at the University of California, Berkeley, and the London School of Economics and Political Science. He received a B.A. and an M.A. in mathematics from Carnegie-Mellon University in 1964 and 1965, respectively, and a Ph.D. in statistics from the University of Chicago in 1969. His interests and research are in statistical theory and a variety of applications of statistics, including environmental protection, epidemiology of skin cancer and of mental retardation, apportionment of funds by statistical formulae, linguistics, and the use of statistical assessments in the courts.

ANNE M. SPRAGUE is a staff member of the Committee on National Statistics. She has a B.A. degree in anthropology from the University of Maryland. Her interest in surveys derives from graduate studies involving participant interviews, particularly as they pertain to the nomadic Rom. She has conducted numerous consumer-oriented surveys and has designed and implemented site-specific realty surveys.

ROGER TOURANGEAU is the technical director for the New York office of the National Opinion Research Center. He received an A.B. in psychology and English from Cornell University in 1973 and a Ph.D. in psychology from Yale University in 1978. He has conducted research on cognitive and social-psychological topics and has served as a statistical consultant on a number of national surveys. He has also taught at Connecticut College, Yale University, and Columbia University.

ENDEL TULVING is professor of psychology at the University of Toronto. He received a B.A. in 1953 from the University of Toronto and a Ph.D. in 1957 from Harvard University. From 1970 to 1975 he taught at Yale University, and in 1977-1978 he spent a year at the University of Oxford as a Commonwealth Visiting Professor. He has done experimental and theoretical research on human memory for over 25 years. His most recent book is *Elements of Episodic Memory* (Oxford University Press, 1983).

About this PDF file: This new digital representation of the original work has been recomposed from XML files created from the original paper book, not from the original typesetting files. Page breaks are true to the original; line lengths, word breaks, heading styles, and other typesetting-specific formatting, however, cannot be retained, and some typographic errors may have been accidentally inserted. Please use the print version of this publication as the authoritative version for attribution.

## INDEX

- Affective carryover effects, [36-37](#)
- Age      frequency judgments, ability to make, [120-21](#),  
[123](#)      memory ability relationship, [6-7](#), [44](#), [92](#)  
          perceptions of national history, [38](#)
- Aids to memory (cues), [4](#), [15](#), [82-83](#), [104](#), [105](#), [111](#), [114](#), [118](#)
- Alternatives      choosing between, [109](#)      middle  
                         alternatives, [116](#)
- Alzheimer's Syndrome (disease), [7](#), [44](#)
- American Association for Public Opinion Research, [69](#), [153](#)
- American Statistical Association, [69](#), [153](#)
- Anchoring and adjustment, [10](#), [12](#), [85-86](#)
- Ancillary information, [16](#)
- Antonyms, [109](#)
- Attitude questions, [14](#), [16](#), [35](#), [36](#), [87-89](#), [112](#)
- Attitudes, [4](#), [91](#)      nonattitudes, [87-88](#)      symbolic  
                         and instrumental, [37](#)      well-formed, [88-89](#)
- Attitude scaling techniques, [90](#), [116](#)
- Availability heuristic, [85](#), [122](#), [124-25](#)
- Bias, [17](#), [65](#), [86](#), [91](#) See also      Record checks.
- Bottom-up processing, [76-77](#)
- Bounded recall, [32](#), [119](#)
- Bureau of Justice Statistics, [2](#)
- Bureau of Social Science Research, [2](#), [149](#)
- Categorized words task, [50](#)
- Census Bureau, [2](#), [3](#), [21](#), [150](#), [152](#)
- Central processor, [103-4](#)
- Chicago, University of, [35](#), [37](#)
- Choice models, [89-90](#)
- Closed-format questions, [105](#), [113-15](#)
- Cognitive abilities      new investigation methods, [9](#)  
                         proposed survey topics, [6-7](#), [44-60](#) See also  
                         Comprehension, Judgment, Memory.
- Cognitive aspects of survey methodology (CASM) seminar  
                         background materials, [71-72](#), [150](#), [157-64](#)      bio-  
                         graphical sketches of participants, [165-70](#)      follow-  
                         up activities [69-70](#), [153-55](#)      organization and  
                         structure, [1-3](#), [149-55](#)      themes and topics, [3-6](#)
- Cognitive failures, [7](#)
- Cognitive processes, in survey response, [35-37](#), [71](#), [73-74](#)
- Cognitive sciences      improvement of survey methodol-  
                         ogy, [10-21](#)

- laboratory-based research on survey methods, 26-34
  - research methods of, as distinguished from survey methodology, 4, 6, 26-27, 31, 91-92, 106
  - study area of, 1
  - survey methodology, collaboration efforts, 1-2, 24, 69-70
  - surveys as cognitive experiments, 7-8
- Committee on National Statistics, 1, 2, 149, 152
- Communication errors, 103
- Comprehension, 2, 13-14, 74-79, 107, 108
- Computer-assisted telephone interviewing (CATI), 15, 17, 151
- Conditioning, 32-33
- Conditions See Health.
- Confidence measures, 16, 20, 66-68
- Confidentiality of responses, 33-34
- Consistency, 91
- Construct validity, 143
- Context comprehension, 77-78
  - frequency judgments, 122
  - memory, 80, 125
  - question responses, 78, 116-19
- Conversations, and interviews, 11, 19
- Correlations, form-resistant, 115
- Crime victimization, 2, 62-63
- Criterion validity, 143
- Cued recall tests, 50
- Cues (memory aids), 4, 15, 82-83, 104, 105, 111, 114, 118
- Cultural and social differences
  - attitudes, 38-39
  - comprehension, 13
  - illness concepts, 6, 22
  - memory, 47
  - motivation, 18
  - scripts, 12
- Dating events, 35, 37, 105, 123-24
- Debriefing, 20
- Decision models, 89-90
- Deliberate errors, 103
- Denial scales, 16
- Dichotomous response categories, 116
- “Don't know” (DK) filters, 115-16
- Economic and Social Research Council, 69
- Economic beliefs survey, 7
- Education levels, and response quality, 115
- Elimination-by-aspects model, 90
- Emotional experiences survey, 7
- Emotions
  - emotional events, memory of, 83
  - response effects of, 36-37
- Encoding specificity, 80, 104, 112
- Episodic memory, 45, 49, 79-80, 104, 124
- Errors, 4, 11
  - match errors, 134-36
  - offsetting errors, 143-44
  - reduction strategies, 13-20
  - sources of See Response effects.
  - types of, 103
- Estimates
  - over- and underestimation, 121-23
  - past and future behavioral estimates for self and others, 65-68
  - preliminary estimates (anchors), 10, 12, 85-86
  - probability estimates, 16
  - time and frequency judgments, 10, 35, 37, 119-25
- Estimation strategies, 16
- Ethnographic studies, 21
- Examples, to clarify questions, 13-14, 15
- Factual questions, 86, 87
- Fatigue effects, 18
- Filters, 115-16
- Floater, 115-16
- Forewarning of interview content, 15, 83
- Forgetting, 81-82, 92-93, 137, 142-44
- Form-resistant correlations, 115

- Fragment completion test, 50
- Frames (schemata, scripts), 5, 11-12, 14, 36, 65, 75-76, 88, 104, 109
- Free recall tests, 50
- Frequency judgments, 10, 35, 37, 119-25
- Gallup survey, 108
- General Social Survey, 110, 117, 153
- Health  
  identification of health condition, 6, 22-23  
  protocol analysis, 61-64  
  restricted activity, 23-24  
  self-perceptions of, 24  
  utilization of medical services, 21-22
- Health and Nutrition Examination Survey, 23
- Health surveys, 23  
  record-check estimates of reporting errors, 137-42 See also National Health Interview Survey.
- Homonyms, 79
- Human information processing, model, 103-5
- Hypermnesia, 12-13
- Illness concepts. See Health.
- Information integration theory, 84
- Information processing, model, 103-5
- Institute for Social Research, 150
- Institut für Demoskopie, 69
- Instrumental attitudes, 37
- Interfering variables, 9
- Interviewers  
  behavior, improvement strategies for, 17  
  response effects, 101, 102  
  role of, 5, 9, 65, 102
- Interviews  
  computer-assisted telephone interviewing (CATI), 15, 17  
  conceptual model, 101  
  conduct of, 5, 18-20  
  length of (respondent burden), 18, 118  
  matching respondent script, 12, 19  
  ordinary conversations, comparison, 11  
  presence of other persons, 9, 19, 20, 118-19  
  validation interviews, 8, 20
- Judgment, 2, 84-89  
  frequency judgments, 35, 37, 66-68, 85-87, 119-25  
  improvement strategies, 15-17  
  proposed research on, 9-10
- Judgmental heuristics, 12, 85-86, 124-25
- Justice Statistics, Bureau of, 2
- Landmark events, 15, 123-24
- Lie scales, 16, 20
- Long-term memory, 45, 49, 79, 80, 103, 104
- Luce choice rule, 89-90
- Management and Budget, Office of, 118
- Match errors, 134-36
- Max Planck Society, 69
- Medical services  
  reporting biases, 137-42  
  utilization of, 21-22
- Memory  
  aids (cues) 4, 15, 82-83, 104, 105, 111, 114, 118  
  as associative network, 80  
  definitions and types of, 45, 79-80  
  distribution of memory abilities, 5, 44, 52  
  episodic memory, 45, 49, 79-80, 104, 124  
  forgetting, 81-82, 92-93, 137, 142-44  
  memorable events, 83-84  
  national and personal history intersection, 7, 25, 38-43, 123-24

- national memory inventory, 6, 25, 44-60  
    organization of memory, 5, 9  
    proposed research areas, 8-9  
    proposed survey topics, 6-7  
    reconstruction of information, 80-81, 113  
    semantic memory, 45, 49, 79-80, 104, 124  
    written and spoken language, 79
- Memory errors, 103
- Memory tasks, 46-49
- Memory tests, 46  
    national memory inventory, 25, 44-60
- Mental health, 23
- Metamemory, 7
- Michigan, University of Institute for Social Research, 150  
    Survey Research Center, 154
- Middle alternatives, 116
- Milbank Memorial Fund, 69-70
- Minnesota Multiphasic Personality Inventory, 16
- Multitrace theory, 124
- National and personal history relationship, 7, 25, 38-43, 123-24
- National Center for Health Statistics, 3, 25, 26-34, 150, 152
- National Crime Survey, 2, 61
- National Elections Survey, 150
- National Health Interview Survey (NHIS), 6, 10, 14, 18, 19, 21-24, 61, 118  
    as focus of CASM seminar, 3, 149-55  
    as test vehicle for laboratory research on survey methods, 26-34
- National Medical Care Utilization and Expenditure Survey, 23
- National memory inventory, 6, 25, 44-60  
    battery description, 46-50  
    data collection and analysis, 50-52  
    instructions, materials, tests, 52-59  
    memory and memory tests, 45-46  
    organization of memory, 5, 9  
    proposed research areas, 8-9  
    proposed survey topics, 6-7  
    reconstruction of information, 80-81, 113  
    semantic memory, 45, 49, 79-80, 104, 124  
    written and spoken language, 79
- National Opinion Research Center, 35, 36, 108, 153
- National Science Foundation, 1, 152, 153
- National survey on cognitive failures, 7
- Network, memory as, 80
- Network sampling, 33
- Nonattitudes, 87-88
- “No opinion” response category, 115-16
- Office of Management and Budget, 118
- Offsetting errors, 143-44
- Open-ended questions, 89, 105, 113-15
- Opinion (attitude) questions, 14, 16, 35, 36, 87-89, 112
- Optimism, 93
- Oral presentation, 78-79
- Overreporting, 119, 144
- Paired-associate tasks, 50
- Part-set cuing, 13
- Personal history  
    national history relationship, 7, 25, 38-43  
    survey question order, 14
- Practice effects, 121
- Presentation mode, 78-79
- Primacy effects, 42
- Primary memory, 45
- Prior knowledge, 77
- Proactive inhibition, 14
- Probability estimates, 16
- Protocol analysis, 12, 20, 61-64
- Proxy respondents, 33, 66-68
- Question order  
    affective carryover effects, 36-37  
    as interpretive context, 14, 78, 116-18  
    interviewer discretion to determine, 17  
    and organization of memory, 5, 9, 12, 14, 19  
    and proactive inhibition, 14



- Questions      additional questions to produce additional information, 62      attitude questions, 14, 16, 35, 36, 87-89, 112      closed-format questions, 105, 113-15      different interpretations by respondents, 6, 13-14, 77-78, 107-8      different questions on same topic, 108-9      to elicit reasons, 35, 36      filter and floater options, 115-16      focus of attention in, 112-13      generality/specificity dimension, 110-11      lengthening introduction to, 111, 118      open-ended questions, 89-105, 113-15      types of, 35, 86-87
- Question wording and context and response effects, 14, 36, 102-3, 106-13
- Randomized response, 33
- Rating scales, 90, 116
- Reasons, questions to elicit, 35-37, 87
- Recall, 2      bounded recall, 32, 119      improvement strategies, 14-15, 65-66      protocol analysis, 12, 20, 61-64      question order effect on, 5, 9, 12, 14, 19      question wording effect on, 112      recognition distinguished from, 105, 114      as reconstruction, 80-81, 113      retrieval order, 62, 63      telescoping, 9, 31-32, 111, 119, 143-44      validation procedures, 8, 9, 20-21, 130-147
- Recall interval, 143-45
- Recall tests, 50
- Recency effects, 42
- Recognition memory, 50, 80      recall distinguished from, 105, 114
- Reconstruction of information, 80-81, 113
- Record checks, 8, 9, 20, 72, 130-47      basic logic of, 130-32      matching errors, 134-36      net survey response bias, 132-37, 142-45      reporting errors in health surveys, 137-42
- Reference periods, 15, 22
- Reporting biases, 91 *See also*      Record checks.
- Representativeness heuristic, 85
- Research methods      cognitive sciences/survey research methodology compared, 4, 6, 26-27, 31, 91-92, 106      laboratory-based research on survey methods, 26-34      material to be recalled unknown to researcher, 9      tools for methodological research, 20-21
- Respondent burden, 118
- Respondents      choice of respondents, 19-20      cognitive process involved in survey responding, 35-37, 73-74      debriefing, 20      interactions among household members, 9, 19, 20, 118-19      interview behavior, 18-20      interview role of, 101-2      motivation, 18-19      multiple respondents, 19-20      preserving anonymity of, 33-34      proxy respondent, 33, 66-68      response effects, 101, 102
- Response categories, 113-16
- Response effects/errors, 11, 13-17, 71-72, 101-29      conceptual model, 101-3      contextual meaning and, 116-19      of emotions, 36-37      human information processing, model, 103-5      perceptions of confidentiality and, 33-34

- question wording and, [36](#), [102-3](#), [106-13](#)      response  
categories, [113-16](#)      survey and cognitive research  
traditions compared, [106](#)      temporal memory and,  
[119-25](#)
- Response selection, [89-90](#)
- Rest periods, [14](#)
- Restricted activity [See](#)      Health.
- Retrieval [See](#)      Recall.
- Retrieval cues (memory aids), [4](#), [82-83](#), [104](#), [105](#), [111](#), [114](#),  
[118](#)
- Royal Statistical Society, [1](#)
- Saliency, [16](#), [119](#), [123](#), [125](#)
- Satisficing rules, [90](#)
- Schemata, scripts (frames), [5](#), [11-12](#), [14](#), [36](#), [65](#), [75-76](#), [88](#),  
[104](#), [109](#)
- Secondary memory, [45](#)
- Self-enumeration, [34](#)
- Self-report data, [5](#)      compared with proxy reports, [66-68](#)
- Semantic memory, [45](#), [49](#), [79-80](#), [104](#), [124](#)
- Sensitive items, [16](#), [33-34](#), [108](#), [142](#)
- Sensory memory, [79](#)
- Short-term memory, [45](#), [49](#), [79](#), [103](#), [104](#)
- Social Science Research, Bureau of, [2](#), [149](#)
- Social Science Research Council (United Kingdom), [1](#)
- Social Science Research Council (United States), [70](#), [153](#)
- Spacing effects, [121](#), [122](#), [125](#)
- Split-ballot studies, [20](#), [21](#), [36](#), [106](#)
- Squish effects, [123](#)
- Stereotypes, [88](#)
- Stores, [106](#)
- Strength models, [124](#)
- Survey bias [See](#)      Record checks.
- Survey methodology      cognitive sciences,      col-  
laboration efforts, [1-2](#), [24](#), [69-70](#)      improvement  
strategies, [10-21](#)      laboratory-based research on  
[26-34](#)      research methods of, as distinguished  
from cognitive sciences [4](#), [6](#), [26-27](#), [31](#), [91-92](#), [106](#)
- Survey Research Center, [154](#)
- Surveys      as cognitive experiments, [7-8](#)      method-  
ological research component, [11](#)      proposed  
research areas, [8-10](#)      proposed surveys on cogni-  
tive abilities, [6-7](#)      reliability assessments, [11](#) [See](#)  
[also](#)      Interviews.
- Symbolic attitudes, [37](#)
- Synonyms, [109](#)
- Telescoping, [9](#), [31-32](#), [111](#), [119](#), [143-44](#)
- Temporal memory, [10](#), [35](#), [37](#), [105](#), [119-25](#)
- Threshold process models, [116](#)
- Time for retrieval, [143-45](#)
- Top-down processing, [75-76](#)
- Trait models, [115-16](#)
- Underreporting, [6](#), [8-9](#), [12-13](#), [15](#), [21-23](#), [118](#), [144](#)
- United Kingdom, [1](#)
- Utilization [See](#)      Health.
- Validation interviews, [8](#), [20](#)
- Validity checks, [20-21](#) . [See also](#)      Record checks.
- Verbatim memory, [80-81](#)
- Videotaped interviews, [17](#), [20](#), [153](#)
- Visual aids, to clarify questions, [14](#)
- Wechsler Memory Scale, [46](#)
- Yale University, [35](#)
- Zentrum für Umfragen und Methodische Analyse, [69](#)