

Enhancing Collaboration in Computer-Administered Survey Interviews

Michael F. Schober

New School for Social Research
Department of Psychology, AL-340
65 Fifth Avenue
New York, NY 10003
schober@newschool.edu

Frederick G. Conrad

Bureau of Labor Statistics
Room 4915
2 Massachusetts Ave., N.E.
Washington, DC 20212
conrad_f@bls.gov

Jonathan E. Bloom

Dragon Systems, Inc.
320 Nevada St.
Newton, MA 02160
jonathan_bloom@dragonsys.com

Abstract

We investigated the extent to which a collaborative view of human conversation transfers directly to interaction with non-human agents, and we examined how a collaborative view can improve user interface design. In two experiments we contrasted user-initiated and system-initiated clarification in computer-administered surveys. In the first (text-based) study, users who could clarify the interpretations of questions by clicking on highlighted definitions comprehended questions more accurately (in ways that more closely fit the survey designers' intentions) than users who couldn't, and thus they provided more accurate responses. They were far more likely to ask for help when they had been instructed that clarification would be essential than when they were merely told that help was available. In the second (speech-based) Wizard-of-Oz study, users responded more accurately and asked more questions when they received unsolicited clarification about question meaning from the system in response to their linguistic cues of uncertainty (*ums* and *uhs*, restarts, talk other than an answer, etc.) than when they did not. The results suggest that clarification in collaborative systems will be successful only if users recognize that their own conceptions may differ from the system's, and if they are willing to take extra turns to improve their understanding.

Introduction

Saying something doesn't guarantee it will be understood. People engage in dialog to make sure that what the speaker intended has been understood—to *ground* their understanding (e.g., Clark and Brennan 1991; Clark and Schaefer 1987, 1989; Clark and Wilkes-Gibbs 1986; Schober and Clark 1989). People ground their understanding to a criterion sufficient for their current purposes; in casual conversations (e.g., at a cocktail party), people may not need to understand precise details to satisfy their conversational goals, but in other settings (e.g., air traffic control tower conversations, calls to a technical help desk when your computer crashes, or

conversations with your ex-spouse about child visitation) the stakes are higher.

This collaborative view of human conversation differs from traditional accounts of language use (what Akmajian et al. 1990 called the "message model" of communication), where listeners interpret utterances directly and all alone. The traditional view is that the meaning of an utterance is contained within the words themselves, and that the process of comprehension involves looking up those meanings in the mental dictionary and combining them appropriately; a collaborative view argues that accurate comprehension also requires dialog so that people can clarify what is meant (see Clark 1992, 1996; Schober 1999a, 1999b).

In the studies reported here we investigate the extent to which this collaborative view of human conversation transfers directly to interaction with non-human agents, and we examine whether a collaborative view can improve user interface design. We propose that examining collaboration in human-computer interaction forces us to specify details of the collaborative view that can test its limits and refine our theories of collaboration.

We contrast two approaches to designing collaborative systems that support the clarification of word meanings. Under one approach, clarification is *user-initiated*—that is, if the user explicitly requests clarification, the system provides it. This requires users to recognize that they need clarification and to be willing to ask for it. Under the other approach, clarification is *system-initiated*—that is, the system provides (or offers to provide) clarification when it diagnoses misunderstanding, based on user behavior. For example, in a text or speech interface a system could provide clarification when the user takes too long to act; in a speech interface a system could provide clarification when the user's speech is hesitant or disfluent (containing *ums* and *uhs*, restarts, etc.).

We examine these issues in the context of survey interviewing systems, where systems present questions and users answer them. To our knowledge, current dialog systems for surveys (see papers in Couper et al. 1998 on

“computerized self-administered questionnaires”) do not allow either user- or system-initiated clarification of meaning. Rather, they embody strict principles of standardization developed for human-human interviews, where the interpretation of questions should be left entirely up to respondents (e.g., Fowler and Mangione 1990). The argument for standardization is that if interviewers help respondents to interpret questions, they might influence responses, but if interviewers read scripted questions and provide only “neutral” feedback, responses are less likely to be biased. We have demonstrated that in human-human interviews even supposedly nonbiasing feedback by interviewers can affect responses (Schober and Conrad 1999a). More importantly, strict standardization can actually harm data quality because it prevents respondents from grounding their understanding of the questions. This is a problem because people’s interpretations of seemingly straightforward questions like “How many bedrooms are there in your house?” can vary enormously; without grounding their understanding of questions, respondents may conceive of questions in unintended ways, and the resulting data may not fulfill the survey designers’ purposes (Clark and Schober 1991). We have shown that responses in strictly standardized interviews can be less accurate than responses in more interactive interviews where respondents can ground their understanding of questions with the interviewers (Conrad and Schober 1999b; Schober and Conrad 1997, 1999b).

Dialog systems for surveys differ from many human-computer interaction situations. First, in survey systems users provide information to the system rather than retrieving information from the system, and so the task goals are quite different from those in, say, a database query system or a web search interface, where the user extracts information from the system. Second, survey system users’ desire for precise understanding may be lower than when they interact with other systems. Users may care less about precisely understanding the words in survey questions when providing opinions to market researchers (misunderstanding has few consequences for the user) than understanding the words in an on-line job application or an on-line health claims form (where misunderstandings can be costly).

Experimental Methods

In our studies we assess whether systems that enable users to clarify the concepts in survey questions do actually lead to improved comprehension of those questions (and thus improved response accuracy), as a collaborative theory would predict. We examine the effects of clarification on task duration—clarification probably takes more time, and this may offset any benefits of clarification. We also examine the effects of clarification on user satisfaction; even if clarification (user- or system-initiated) improves comprehension, it could be annoying.

Our first study (Conrad and Schober 1999a) uses a text interface, in which the computer displays questions on a screen. The user enters responses and asks for clarification with the keyboard and mouse. Our second study (Bloom 1999; Bloom and Schober 1999) uses a speech interface, in which the computer, using a synthesized voice, asks questions through a headset. The user answers questions and asks for clarification by speaking into the headset microphone.

In both studies, all users were asked the same survey questions, which had been used in earlier studies of human-human survey interviews (Schober and Conrad 1997, 1999b; Schober, Conrad and Fricker 2000). We adapted 12 questions from three ongoing government surveys. Four questions were about employment, from the Current Population Survey (e.g., “Last week, did you do any work for pay?”); four questions were about housing, from the Consumer Price Index Housing survey (e.g., “How many people live in this house?”); four questions were about purchases, from the Current Point of Purchase Survey (e.g., “During the past year, have you purchased or had expenses for household furniture?”). The three question domains (employment, housing, purchases) were randomly ordered for different respondents, although the questions within a domain were always presented in the same order as they appeared in the government surveys from which they were borrowed. For each question, the survey designers had developed official definitions for the key concepts, which clarified whether, for example, a floor lamp should be considered a piece of household furniture, or whether a student away at college should be considered to be living at home.

Users answered these questions on the basis of fictional scenarios, so that we could measure response accuracy—that is, the fit between users’ answers and the survey designers’ official definitions. For each question there were two alternate scenarios, one typical and one atypical. With the typical scenario, the survey question was designed to be easy for users to interpret—to map onto the user’s (fictional) circumstances in a straightforward way. For example, for the question “Has Kelley purchased or had expenses for household furniture?”, the typical scenario was a receipt for an end table, which is clearly a piece of furniture. With the atypical scenario, it was less clear how the survey question should be answered. For example, for the household furniture question the atypical scenario was a receipt for a floor lamp, which is harder to classify without knowing the official definition of “household furniture.”

For each user, half the scenarios described typical situations and half atypical situations.

Study 1: Text interface

In this study, we varied the way the survey system provided clarification. When clarification was user-initiated, users could request the official definition for a survey concept by clicking the mouse on highlighted text

in the question. When clarification was system-initiated, the system would also offer a definition when users took “too long” to respond. This was defined as taking longer than the median response time for atypical scenarios when no clarification was available (see experimental design below). Clarification was offered as a Windows dialog box; users could reject the offer by clicking “no” if they didn’t think clarification was needed.

We also varied instructions to the users about how precisely they would need to understand the system’s questions—that is, we varied the grounding criterion. Some users were told that clarification was essential; they were encouraged to obtain definitions from the computer because their everyday definitions might differ from the survey’s. Other users were told merely that clarification was available, that definitions would be available if users wanted them.

So there were five experimental conditions:

	<u>Type of clarification</u>	<u>User instructed that...</u>
1	no clarification	
2	at user’s request	Clarification essential
3	at user’s request	Clarification available
4	when user takes too long <u>or</u> at user’s request	Clarification essential
5	when user takes too long <u>or</u> at user’s request	Clarification available

The 54 users, recruited from an advertisement in the *Washington Post*, were paid to participate. Most (44) reported using a computer every day.

Results

Users’ responses were virtually perfectly accurate (their responses fit the official definitions) when they answered about typical scenarios. For atypical scenarios, users were more accurate when they could get clarification than when they couldn’t (see Figure 1). Response accuracy mainly depended on the instructions to the user about the grounding criterion. When users had been told that definitions were merely available, their accuracy was as poor as when they couldn’t get clarification. When they had been told that definitions were essential, response accuracy was much better.

Response accuracy was strongly related to how often users received clarification. As Figure 2 shows, when users had been told that definitions were essential, they received help most of the time; in fact, they frequently asked for help for typical scenarios, when clarification presumably wasn’t necessary. They asked for clarification so often that the system rarely initiated clarification—the users requested help before the system-initiated clarification was triggered. In contrast, users who had been told that clarification was merely available rarely asked for clarification, and they responded so quickly that

system-initiated clarification was rarely triggered. It seems that it didn’t occur to these users that their interpretation of ordinary terms like “bedroom” and “job” might be different from the system’s, and so they answered confidently, quickly, and inaccurately.

As Figure 3 shows, clarification took time. Response times were much longer in cases where users received clarification. As we anticipated, improved accuracy from clarification can be costly.

Users’ ratings of their satisfaction with the system indicated two things. First, users who could not get clarification reported that they would have asked for clarification if they could. This suggests that interacting with dialog survey systems that don’t allow clarification may be relatively unsatisfying. Second, users’ grounding criteria affected their perceptions of the system. System-initiated clarification was rated useful and not annoying by “clarification essential” users, and less useful and more annoying by “clarification available” users. Presumably users who had been told that clarification was available found it jarring for the system to offer unsolicited help for seemingly straightforward questions.

Overall, these results suggest that the success of human-machine collaboration may depend both on users’ grounding criteria—how important they believe it is to understand accurately—and also on whether users recognize that system concepts may differ from theirs.

Study 2: Speech interface

This study used a Wizard-of-Oz technique to simulate a speech interface. Users believed they were interacting with a computer, when actually a hidden experimenter presented the questions and scripted clarification. To enhance believability, we used an artificial-sounding text-to-speech computer voice (Apple’s “Agnes” voice).

This study used exactly the same questions and scenarios as Study 1. Users participated in one of four experimental conditions. In the first condition, the system never provided clarification. In the second condition, clarification was user-initiated—the system would provide clarification if users asked for it explicitly. In the third condition, clarification was not only user-initiated but also system-initiated—the system would “automatically” provide full definitions when users displayed specific uncertainty markers that had been shown to be more prevalent in atypical situations in human-human interviews collected with these materials (Bloom and Schober 2000). These included *ums*, *uhs*, pauses, repairs, and talk other than an answer. In the fourth condition, the system always provided clarification; no matter what the user did, the system would present the full official definition for every question.

All users were told that the survey concepts might differ from theirs. This instruction is less extreme than the “clarification essential” instruction in Study 1, in which users were told that they would likely need definitions to

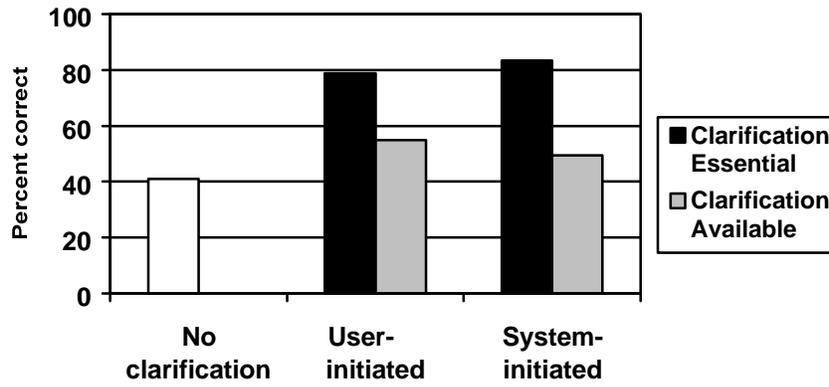


Figure 1: Response accuracy for atypical scenarios, Study 1

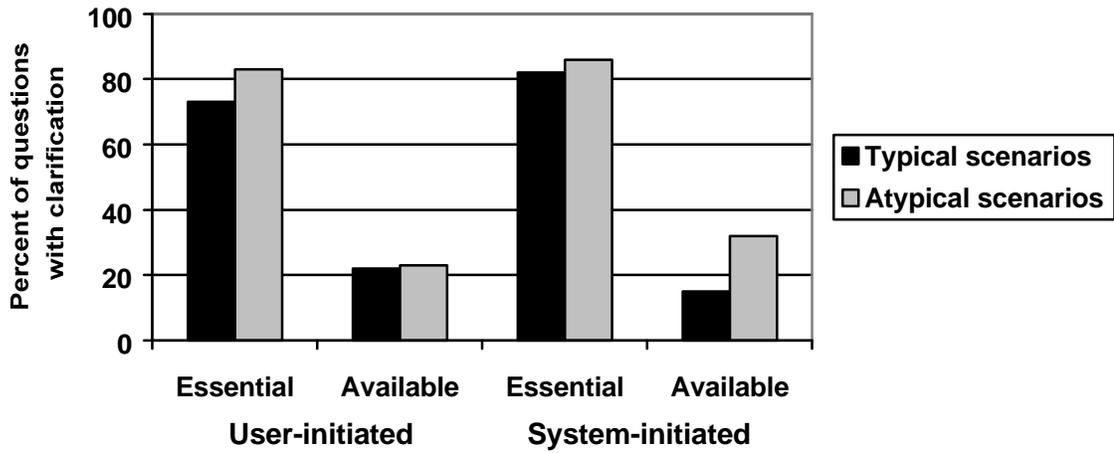


Figure 2: How often users received clarification, Study 1

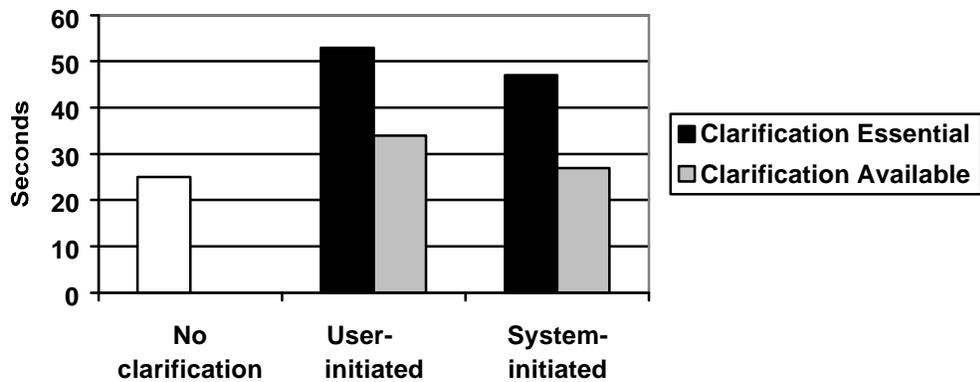


Figure 3: Response time per question, Study 1

get the answers right, but it goes beyond the “clarification available” instruction from Study 1, where users weren’t told that the system’s definitions might differ from theirs. The users in the user- and system- initiated clarification conditions were told that they should ask for clarification if they felt they needed it. The users in the system-initiated clarification condition were also told that the computer would sometimes provide unsolicited assistance. The users in the clarification-always condition were told that for each question they would hear definitions of key concepts.

40 users, recruited from an advertisement in the *Village Voice* or members of the New School community, were paid to participate.

Results

As in Study 1, users’ responses were almost perfectly accurate when they answered about typical scenarios. For atypical scenarios, users were substantially more accurate when they were always given clarification than when they were never given clarification (see Figure 4). When clarification was user-initiated, response accuracy was no better than when users received no clarification, because users almost never asked for clarification (one user asked one time). As in Study 1, it seems likely that it didn’t occur to users that clarification was necessary.

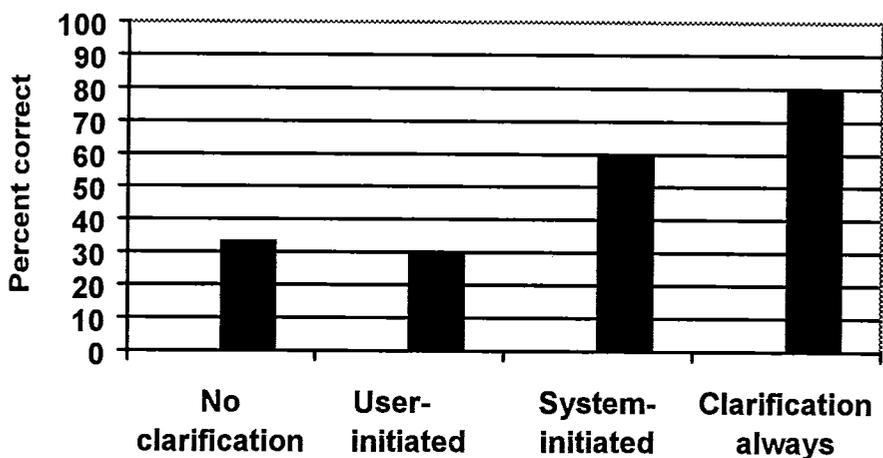


Figure 4: Response accuracy for atypical scenarios, Study 2

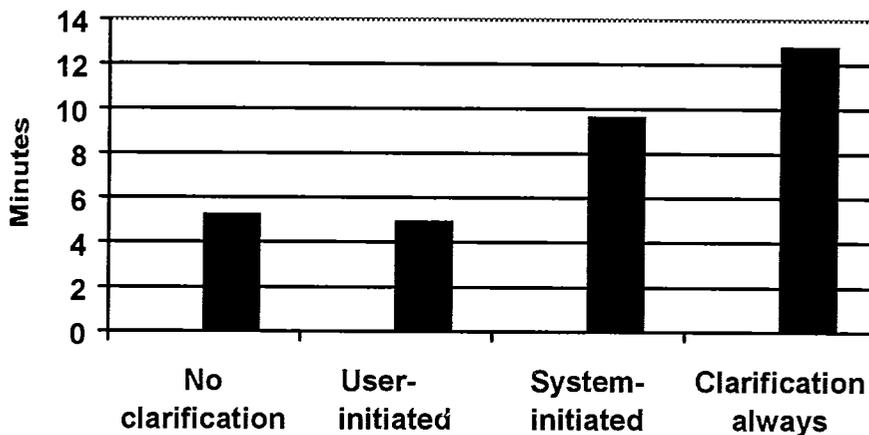


Figure 5: Interview duration, Study 2

Response accuracy was better when clarification was system-initiated, although it was not as good as when clarification was given always. Users in the system-initiated clarification condition were given clarification 52% of the time (46% of the time for typical scenarios, 59% of the time for atypical scenarios); their accuracy for atypical scenarios corresponded perfectly with how often they were given clarification.

System-initiated clarification increased the amount of user-initiated clarification: users were more likely to ask questions in the system-initiated condition, presumably because they were more likely to recognize that clarification might be useful. These users also spoke less fluently, producing more *ums* and *uhs*. The system's unsolicited clarification might either have made the users less certain that their concepts matched the system's, leading them to unintentionally display their uncertainty by being less fluent, or it might have led users to signal that they wanted clarification by being less fluent.

Overall, the users in this study requested clarification (user-initiated) far less often than the users in Study 1. This might result from any or all of the differences between our text and speech interfaces. In the speech interface, clarification was harder to request; requests had to be formulated into explicit questions rather than being triggered by simple mouse clicks. Also, in the speech interface the definition unfolded over time (sometimes a substantial amount of time, up to 108 seconds), rather than appearing all at once, and in our application it was impossible to shut off; in the text interface, the definition appeared all at once and could be closed with a simple mouse click. Also, unlike in the text study, users couldn't reject system-initiated offers of clarification; here the system immediately provided clarification when triggered, without giving the option of rejecting the help.

As in Study 1, clarification took time. The more clarification a user received, the more time the interviews took. Interviews where clarification was always provided took twice as long as interviews with no clarification; system-initiated clarification took an intermediate amount of time (see Figure 5).

Also as in Study 1, users rated the system more positively when it was responsive (user- or system-initiated conditions). When the system was not responsive (no clarification or clarification always), users wanted more control and felt that interacting with the system was unnatural. Users didn't report finding system-initiated clarification particularly more annoying than user-initiated clarification—which they almost never used.

Overall, these results suggests that enhancing the collaborative repertoire of a speech system can improve comprehension accuracy without harming user satisfaction, as long as the system provides help only when it is necessary. But these improvements come at the cost of increased task duration, which could make such systems less practical in real-world survey situations.

Conclusions

Our findings demonstrate that a collaborative view can indeed transfer to interaction with non-human agents. Increased system clarification abilities can improve users' comprehension (and thus their response accuracy), while increasing (or not reducing) user satisfaction. But this comes at the cost of increased task duration, which could lower survey completion rates in the real world.

Our findings also demonstrate that extended clarification sequences are likely to be rare or unnecessary when users' conceptions are likely to be the same as the system's, as in our typical scenarios. The need for building survey systems with enhanced collaborative abilities may depend on the likelihood of potential misunderstandings; if this likelihood is high or unknown, enhanced collaborative abilities may be worth implementing.

The benefits for collaboratively enhanced survey systems come even with our rudimentary implementations, which are based on the most generic of user models (see Kay 1995 for discussion of different types of user models). A stronger test of collaborative approaches requires more customized interfaces, in which, for example, the system would reason about which parts of definitions would be appropriate to present at any given moment, what particular users are likely to misunderstand, etc. (see Moore 1995 for discussion of better explanation-giving systems).

Our findings demonstrate that computer implementations of surveys seem to run into exactly the same problems as human-human survey and instructional situations, where people don't always recognize they need help or aren't willing or able to ask for help (e.g., Graesser and McMahan 1993; Graesser et al. 1996; Schober and Conrad 1997).

But our findings also show that in some situations (our text interface, when users were told that clarification was essential), users are indeed willing to ask for clarification more often than they are with human interviewers (Schober and Conrad 1997). This is consistent with findings in other domains that interaction with a computer can lead to better task outcomes than interaction with a person. For example, people may feel better about working with an intelligent computer tutor than doing the same work in the classroom (Schofield 1995), and people are more willing to admit to sensitive behaviors when asked about them on self-administered computer surveys than in human-administered surveys (Tourangeau and Smith 1996).

We propose that some of these improvements from interacting with computers don't arise simply from the fact that the computer isn't a person. They arise in part from the fact that the costs and constraints of grounding vary in different media, as Clark and Brennan (1991) argued. Most tutoring and survey systems to date have been direct manipulation or simple (textual) character

entry systems like our text interface; in such interfaces the user's costs of presenting information to the system can be low. The human interactions to which such systems are often compared are speech interactions, where people have to formulate clarification requests explicitly and clarification takes significant amounts of time. Any differences in task performance may just as likely result from the differences between direct manipulation and speech as from the differences between computers and humans.

We believe our findings also require us to refine a theory of human-human collaboration by explicitly introducing the human-computer notion of initiative. Our findings that comprehension success can vary depending on whether the user or system takes the initiative should be extended to the human realm; a collaborative theory should include who takes the responsibility for clarifying meaning. In many cases speakers are responsible for what they mean, and listeners assume that what speakers say is readily interpretable to them in the current context (the "interpretability presumption," in Clark and Schober's 1991 terms). But in situations where the speaker is less competent or knowledgeable than the addressee, the addressee may take responsibility for the meaning, and may initiate clarification (Brennan, 1990; Schober, 1998). Who should be responsible under what circumstances, and what determines how speakers decide whose effort should be minimized, are important questions for a theory of collaboration.

Altogether, our results suggest that user-initiated clarification will work only if users recognize that clarification will help, recognize that the system's concepts may differ from theirs, are motivated to understand precisely, and are willing to take the extra turns to ground understanding. Explicit instructions to users can help make this happen—help set a high grounding criterion—but it's unclear whether such instruction is feasible in real-world situations. Our results suggest that system-initiated clarification will work only if users give reliable evidence of misunderstanding and if they are willing to accept offers of clarification. It won't work if users are confident in their misinterpretations.

In sum, the opportunity for dialog won't help if users don't recognize it's needed.

Acknowledgments

This material is based upon work supported by the National Science Foundation under Grant No SBR9730140 and by the Bureau of Labor Statistics. The opinions expressed here are the opinions of the authors and not those of the Bureau of Labor Statistics. We thank Susan Brennan and Albert Corbett for help on an earlier version of this paper.

References

- Akmajian, A., Demers, R.A., Farmer, A.K., and Harnish, R.M. 1990. *Linguistics: An Introduction to Language and Communication, 3rd Edition*. Cambridge, MA: The MIT Press.
- Bloom, J.E. 1999. Linguistic Markers of Respondent Uncertainty during Computer-Administered Survey Interviews. Ph.D. diss., Dept. of Psychology, New School for Social Research.
- Bloom, J.E., and Schober, M.F. 1999. A Speech Interface for Collaborating with a Computer-Administered Survey System. Manuscript in preparation.
- Bloom, J.E., and Schober, M.F. 2000. Respondent Cues that Survey Questions Are in Danger of Being Misunderstood. In *Proceedings of the American Statistical Association, Section on Survey Research Methods, 1999*. Alexandria, VA: American Statistical Association. Forthcoming.
- Brennan, S.E. 1990. Seeking and Providing Evidence for Mutual Understanding. Ph.D. diss., Dept. of Psychology, Stanford University.
- Clark H.H. 1992. *Arenas of Language Use*. Chicago: University of Chicago Press.
- Clark, H.H. 1996. *Using Language*. Cambridge, UK: Cambridge University Press.
- Clark, H.H., and Brennan, S.E. 1991. Grounding in Communication. In Resnick, L.B., Levine, J.M., and Teasley, S.D. eds. *Perspectives on Socially Shared Cognition*, 127-149. Washington, DC: APA.
- Clark, H.H., and Schaefer, E.F. 1987. Collaborating on Contributions to Conversations. *Language and Cognitive Processes* 2:1-23.
- Clark, H.H., and Schaefer, E.F. 1989. Contributing to Discourse. *Cognitive Science* 13:259-294.
- Clark, H.H., and Schober, M.F. 1991. Asking Questions and Influencing Answers. In Tanur, J.M. ed., *Questions About Questions: Inquiries into the Cognitive Bases of Surveys*, 15-48. New York: Russell Sage Foundation.
- Clark, H.H., Wilkes-Gibbs, D. 1986. Referring as a Collaborative Process. *Cognition* 22:1-39.
- Conrad, F.G., and Schober, M.F. 1999a. A Conversational Approach to Computer-Administered Questionnaires. In *Proceedings of the American Statistical Association, Section on Survey Research Methods, 1998*, 962-967. Alexandria, VA: American Statistical Association.

- Conrad, F.G., and Schober, M.F. 1999b. Conversational Interviewing and Measurement Error in a Household Telephone Survey. *Public Opinion Quarterly*, under revision.
- Couper, M.P., Baker, R.P., Bethlehem, J., Clark, C.Z.F., Martin, J., Nicholls II, W.L., and O'Reilly, J.M. eds. 1998. *Computer Assisted Survey Information Collection*. New York: John Wiley & Sons, Inc.
- Fowler, F.J., and Mangione, T.W. 1990. *Standardized Survey Interviewing: Minimizing Interviewer-Related Error*. Newbury Park, CA: SAGE Publications, Inc.
- Graesser, A.C., and McMahan, C.L. 1993. Anomalous Information Triggers Questions When Adults Solve Quantitative Problems and Comprehend Stories. *Journal of Educational Psychology*, 85:136-151.
- Graesser, A.C., Swamer, S.S., Baggett, W.B., and Sell, M.A. 1996. New Models of Deep Comprehension. In Britton, B.K. and Graesser, A.C. eds., *Models of Understanding Text*, 1-32. Mahwah, NJ: Lawrence Erlbaum Associates.
- Kay, J. 1995. Vive la difference! Individualized interaction with users. In Mellish, C.S. ed., *Proceedings of the 14th International Joint Conference on Artificial Intelligence*, 978-984. San Mateo, CA: Morgan Kaufmann Publishers.
- Moore, J.D. 1995. *Participating in Explanatory Dialogues: Interpreting and Responding to Questions in Context*. Cambridge, MA: MIT Press.
- Schober, M.F. 1998. Different Kinds of Conversational Perspective-Taking. In Fussell, S.R., and Kreuz, R.J. eds., *Social and Cognitive Psychological Approaches to Interpersonal Communication*, 145-174. Mahwah, NJ: Erlbaum.
- Schober, M.F. 1999a. Conversational Evidence for Rethinking Meaning. *Social Research* 65:511-534.
- Schober, M.F. 1999b. Making Sense of Questions: An Interactional Approach. In Sirken, M.G., Herrmann, D.J., Schechter, S., Schwarz, N., Tanur, J.M., and Tourangeau, R. eds., *Cognition and Survey Research*, 77-93. New York: John Wiley & Sons.
- Schober, M.F., and Clark, H.H. 1989. Understanding by Addressees and Overhearers. *Cognitive Psychology* 21:211-232.
- Schober, M.F., and Conrad, F.G. 1997. Does Conversational Interviewing Reduce Survey Measurement Error? *Public Opinion Quarterly* 61:576-602.
- Schober, M.F., and Conrad, F.G. 1999a. A Collaborative View of Standardized Surveys. In Maynard, D., Houtkoop, H., Schaeffer, N.C., and van der Zouwen, Z. eds., *Standardization and Tacit Knowledge: Interaction and Practice in the Survey Interview*. New York: John Wiley & Sons. Forthcoming.
- Schober, M.F., and Conrad, F.G. 1999b. Response Accuracy When Interviewers Stray from Standardization. In *Proceedings of the American Statistical Association, Section on Survey Research Methods, 1998*, 940-945. Alexandria, VA: American Statistical Association.
- Schober, M.F., Conrad, F.G., and Fricker, S.S. 2000. Further Explorations of Conversational Interviewing: How Gradations of Flexibility Affect Costs and Accuracy. In *Proceedings of the American Statistical Association, Section on Survey Research Methods, 1999*. Alexandria, VA: American Statistical Association. Forthcoming.
- Schofield, J.W. 1995. *Computer and Classroom Culture*. Cambridge, UK: Cambridge University Press.
- Tourangeau, R., and Smith, T. 1996. Asking Sensitive Questions: The Impact of Data Collection Mode, Question Format, and Question Context. *Public Opinion Quarterly* 60:275-304.