

Chapter 18

A PRIMER ON VERBAL PROTOCOL ANALYSIS

Susan Trickett and J. Gregory Trafton

We have been using verbal protocol analysis for over a decade in our own research, and this chapter is the guide we wish had been available to us when we first started using this methodology. The purpose of this chapter is simply to provide a very practical, “how-to” primer for the reader in the science—and art—of verbal protocol analysis.

It is beyond the scope of this chapter to discuss the theoretical grounding of verbal protocols; however, there are several resources that do so, as well as providing some additional practical guidelines for implementing the methodology: Ericsson and Simon (1993); van Someren, Barnard, and Sandberg (1994); Chi (1997); Austin and Delaney (1998); and Ericsson (2006). In addition, numerous studies—too many to list here—have effectively used verbal protocols, and these may provide additional useful information about the application of this technique.

Briefly put, verbal protocol analysis involves having participants perform a task or set of tasks and verbalizing their thoughts (“talking aloud”) while doing so. The basic assumption of verbal protocol analysis, to which we subscribe, is that when people talk aloud while performing a task, the verbal stream functions effectively as a “dump” of the contents of working memory (Ericsson & Simon, 1993). According to this view, the verbal stream can thus be taken as a reflection of the cognitive processes in use and, after analysis, provides the researcher with valuable information about not only those processes, but also the representations on which they operate. In addition, verbal protocols can reveal information about misconceptions and conceptual change, strategy acquisition, use, and mastery, task performance, affective response, and the like.

Verbal protocol analysis can be applied to virtual environments in several ways. For example, verbal protocols collected during performance by both experts and novices can inform the design of virtual environments to be used for training. By comparing expert and novice performance, a designer could identify specific areas where learners would benefit from support, remediation, or

instruction. Similarly, having experts talk aloud while performing a task could help a curriculum designer identify and develop the content to be delivered in the virtual environment. Verbal protocols are likely to be especially useful in evaluation—first, in determining how effective the virtual environment is as a training tool, and second, in assessing the student’s learning and performance, where verbal protocols can provide valuable information about the learner’s cognitive processes, beyond simple measures of accuracy and time on task.

Before explaining the “nuts and bolts” of verbal protocol analysis, we should note that there are several caveats to using this method. First, it is important to distinguish between *concurrent* and *retrospective* verbal protocols. Concurrent protocols are delivered at the same time as the participant performs the task and are ideally unprompted by the experimenter. Retrospective protocols are provided *after* the task has been completed in response to specific questions posed by the experimenter, such as “How did you solve this problem, and why did you choose that particular strategy?” The chief problem with retrospective protocols is that people do not necessarily have access either to what they did or why they did it (see Nisbett and Wilson, 1977, for a full discussion of this issue; Ericsson and Simon, 1993, for an overview of research on retrospective versus concurrent protocols; and Austin and Delaney, 1998, for a discussion of when retrospective protocols may be reliably used).

Second, it is important to differentiate between having a participant think aloud while problem solving and having him or her describe, explain, or rationalize what he or she is doing. The main problem with participants providing explanations is that this removes them from the *process* of problem solving, causing them to think about what they are doing rather than simply doing it, and as a result, such explanations may change their performances. Ericsson and Simon (1993) provide a complete discussion of the effects of such “socially motivated verbalizations” on problem solving behavior, and later in this chapter we suggest ways to reduce the likelihood that they will happen.

Third, it is important to be aware of the potential for subjective interpretation at all stages of collecting and analyzing verbal protocol data, from interpreting incomplete or muttered utterances to assigning those utterances a code. The danger is that a researcher may think he or she “knows” what a participant intended to say or meant by some utterance and thus inadvertently misrepresent the verbal protocol. We will describe ways to handle the data that reduce the likelihood of researcher bias affecting interpretation of the protocols, and we will suggest some methodological safeguards that can also help minimize this risk.

Another important consideration for the researcher is the level of time commitment involved in verbal protocol studies. Simply put, this type of research involves a serious investment of time and resources, not only in collecting data (resources in terms of experimenters, participants, and equipment), but also in transcribing, coding, and analyzing data. It is simply not possible for one researcher to manage all the tasks associated with verbal protocol analysis of a complex task/domain with many participants. Furthermore, both collecting

and coding the data require a high level of training, and thus it is important that members of the research team are willing to make a commitment to complete the project. A coder who quits midstream causes a major setback to the research, because a new recruit must be found and trained before the work can continue.

Verbal protocol data are “expensive data.” Depending on the task to be performed, experimental sessions generally last one or two hours, and participants must usually be “run” singly. In some studies, participants in an undergraduate psychology pool are not suitable candidates for the research; thus, time must be spent identifying and recruiting appropriate participants. In the case of expert studies or “real world” performance, data collection often involves travel to the participant’s work site, and may require coordination with supervisors and co-workers, or even special permission to enter a site. Those who agree to participate in the research usually must give up valuable work time, and schedule changes or other unexpected occurrences may prevent their participation at the last minute. Thus, perhaps even more than in laboratory research, it is crucial to plan as much as possible and to take into account many contingencies that can arise. Such contingencies range from the mundane (for example, triple-checking that the equipment is functioning properly, that backup batteries are available, and that all the relevant forms and materials are in place) to the unpredictable (for example, having backup plans for cancellation or illness).

Despite these caveats, we believe that collecting verbal protocol data is a flexible methodology that is appropriate for research in any research endeavor for which the goals are to understand the *processes* that underlie performance. We think it will be especially fruitful in the domain of virtual reality (VR)/virtual environments (VEs). Not much protocol analysis has been done in this area; however, the technique offers many insights into cognitive processing and is likely to be very useful in comparing VR/VEs to real world performance, for example. Researchers involved in education and training, including VE training, may have many goals—for example, understanding how effective a given training program is, what are its strengths and weaknesses, how it can be improved, how effectively students use the system, whether students actually learn what the instruction is designed to teach, what difficulties they have in mastering the material to be learned, what difficulties they have using the delivery system, and so on. Outcome measures, such as time on task or accuracy on a knowledge assessment test, provide only partial answers and only to some of these questions. Process measures, on the other hand, provide a much richer, more detailed picture of what occurs during a learning or problem solving session and allow a much finer-grained analysis of both learning and performance. Verbal protocols are the raw data from which extremely useful process data can be extracted. Although our emphasis in this chapter is only on verbal protocol data, combining different types of process data (such as verbal protocols, eye-track data, and mouse-click data) can create an extremely powerful method to understand not only *what* is learned, but also *how* learning occurs.

DATA COLLECTION

A number of issues related to data collection should be addressed before the first participant even arrives! Several confidentiality issues surround the act of collecting data, which stem from the potentially sensitive nature of video and/or audio recording participants as they perform tasks. First, in obtaining informed consent, it is crucial that participants fully understand the recording process and have a specific option to agree or disagree to have their words, actions, and bodies recorded. Thus, consent forms need to contain two lines, as follows:

____ I agree to have my voice, actions, and upper body videotaped.

____ I do not agree to have my voice, actions, and upper body videotaped.

Of course, the actual wording will depend on the scope of recording planned—it might additionally include gestures or walking patterns, for example.

Second, agreement to have one's words and actions recorded should not be taken as agreement to have one's words and actions shared with others. Although participants might be quite willing to participate in a study and to be recorded while doing so, they might be uncomfortable at the thought of their data being shared, for example, at conferences and other research presentations. If the researcher is planning to use video clips as illustrations during research talks, it is important to have an additional section of the consent form that specifically asks participants whether or not they agree to have their data shared. They have the right to refuse, and this right must be honored. As mentioned earlier, this does not necessarily disqualify participants from participating in the study; it just means that a participant's video cannot be shared with others. The third issue relating to confidentiality is that of storing the data. It is especially important to reassure participants that the data will be stored confidentially—that data will be identified by a number rather than a participant's name, that data will be locked away or otherwise secured, and that access to the data will be strictly limited to members of the research team.

One of the advantages of verbal protocol data is its richness; the downside of this richness is that the data quickly become voluminous. Even a relatively short session can result in pages of transcribed protocol that must be coded and analyzed. In addition, participants, such as experts, may be scarce and difficult to recruit. For these reasons, verbal protocol studies generally have fewer subjects than other kinds of psychological studies. How many participants are needed? The answer depends on the nature of the study. We (and others) have had as few as one person participate in a particular study. Such single-subject case studies are extremely useful for generating hypotheses that can then be tested experimentally on a larger sample. In general, however, we have found that the more expert the participants, the less variance there is in relevant aspects of performance. With less expert participants, we generally find slightly higher numbers of participants work well. Ideally, we try to include 5 to 10 novices, although this is not always possible. Studies with larger numbers of participants generally involve less exploratory research, for which a coding scheme is already well established and for which only a subset of the data needs to be coded and

analyzed. As with all research, the precise number of participants will rest on the goals of the study (exploratory, confirmatory, case study, and so forth), the nature of the participants (experts, novices, or naive subject pool participants), and practical considerations concerning recruitment and accessibility of participants.

COLLECTING DATA

Adequate recording equipment is a must. Although we have occasionally used audio recording alone, we have found that the additional information contained in video recording can be very valuable when it comes to analyzing the data. Furthermore, we have found that having more than one camera is often worthwhile, for several reasons. First, an additional camera (or cameras) provides an automatic backup in the case of battery or other equipment failure. In addition, a second camera can be aimed at a different area of the experimental setup, thus reducing the need for the experimenter to attempt to follow the participant around while he or she performs the task. In some cases, we have even used three cameras, as additional cameras recording from different angles allow the experimenter to reconstruct more precisely what was happening. We are also somewhat extravagant in our use of videotapes. Even when a session is quite short and leaves a lot of unused space on a tape, we use a new videotape for each participant or new session. Especially with naturalistic tasks, it is difficult to anticipate how long a participant will take. Some people work quite slowly, and underestimating the time they will take can lead to having to change a tape mid-session. Our view is that videotapes and batteries are relatively cheap, whereas missing or lost data are extremely expensive. In short, we have found that thinking ahead to what data we *might* need has helped us to have fewer regrets once data collection is completed.

The most critical aspect of verbal protocol data is the sound, and a high quality lapel microphone is well worth the cost. Many participants are inveterate mutterers, and most participants become inaudible during parts of the task. These periods of incoherence often occur at crucial moments during the problem solving process, and so it is important to capture as much of what may be transcribed as possible. The built-in microphone on most cameras is simply not powerful enough to do so. Using a lapel microphone also allows the researcher to place the camera at a convenient distance behind the participant, capturing both the sound and the participant's interactions with experimental materials but not his or her face. Attaching an external microphone (such as a zoom or omnidirectional microphone) to any additional cameras that are used will allow the second, or backup, camera to capture sound sufficiently well in most cases, although muttering or very softly spoken utterances may be lost. Having such a reliable backup sound recording system is always a good idea.

It is crucial to begin the session with good batteries and to check, before the participant begins the task, that sound is actually being recorded. The experimenter should also listen in during the session to make sure that nothing has gone

amiss. Although this point may seem obvious, most researchers who collect verbal protocols have had the unhappy experience of reviewing a videotape only to find there is no sound. We hope to spare you that pain!

As with any psychological research, participants should be oriented to the task before beginning the actual study. In the case of verbal protocols, this orientation includes training in actually giving verbal protocols. We have a standard training procedure that we use with all participants.¹ During the task, if a participant is silent for more than three or four seconds, it is important to jump in with a prompt. The researcher should simply say, "Please keep talking" or even "Keep talking." More "polite" prompts, such as "What are you thinking?" are actually detrimental, because they invite a more social response from the participant (such as, "Well, I was just wondering whether . . .") and thus remove the participant from the process of problem solving into the arena of social interaction.

Occasionally, a participant will be unable to provide a verbal protocol, and this difficulty generally manifests itself in one of two ways: either the participant remains essentially silent throughout the task, talking aloud only when prompted and then lapsing back into silence, or the participant persists in explaining what he or she is doing or about to do. In the first case, there is very little the experimenter can do, and depending on the situation, it may be better to bail out rather than subject the participant to this uncomfortable situation, since the data will be unusable in any case. One way to help prevent this from happening is to ensure beforehand that participants are fully at ease in whatever language the protocol is to be given. In the second case, the experimenter must make a judgment call as to whether to pause the session and provide some retraining. Retraining will be disruptive and may make the participant even more nervous; however, it is unlikely that a participant will change strategy midstream. In either case, protocols that consist mostly of explanations cannot be used. These problems should not be confused with protocols that contain a great deal of muttering, broken grammatical structures, or incomplete or incoherent thoughts. Although a challenge to the transcriber and the coder, such protocols are quite acceptable and, in fact, represent "good" protocols, since people rarely think in complete sentences.

The objective during verbal protocol collection is to keep the participant focused on the task, without distraction. One obvious way to do this is to make sure that the room in which the session is taking place is not in a public place and that a sign on the door alerts others not to interrupt. Less obviously, it means that the experimenter should resist the urge to ask questions during the session. If domain related questions do arise, the experimenter should make a note of them and ask after the session is completed. Querying participants' explicit knowledge after task completion does not jeopardize the integrity of the verbal protocols, because it taps into stable domain knowledge that is not altered by the retrospective process. Asking questions during problem solving, however, is disruptive, and may alter the participants' cognitive processes.

¹A copy of this procedure is available at <http://www.nrl.navy.mil/aic/iss/aas/cog.complex.vis.php>

PROCESSING THE DATA

Once data collection is completed, verbal protocols must be transcribed and coded before data analysis can begin. Unfortunately, there are no consistently reliable automated tools to perform transcription. Fortunately, transcription can be carried out by less-skilled research assistants, with a couple of caveats. First, transcribers should be provided with a glossary of domain relevant terms; otherwise, when they encounter unfamiliar words, they are likely to misinterpret what is said. Second, transcribers should be trained in segmenting protocols. Both these steps will reduce the need for further processing of the transcriptions before coding can begin.

How protocols are segmented depends on the grain size of the analyses to be performed and will therefore vary from project to project. Careful thought should be given to the appropriate segmentation, because this process will affect the results of the analyses. It is beyond the scope of this chapter to illustrate in detail the many options for segmenting protocols; however, Chi (1997) provides a thorough discussion of different methods and levels of segmentation and their relationship with the type of analyses to be performed.

In general, we have found that segmenting according to complete thought provides us with the flexibility to code data on a number of different dimensions. By “complete thought” we basically mean a clause (a subject and a verb), although given the incomplete nature of human speech, utterances do not always fall neatly into a clausal structure. However, we have had excellent agreement among transcribers by using the complete thought rubric. Table 18.1 shows an example

Table 18.1. Example of Segmenting by Complete Thought

Utterance
OK, 0Z surface map
shows a low pressure in the Midwest moving up
there’s a little bubble high ahead of it
that’s pretty much stacked well with 500 millibars with another low
so that’s looking at severe weather in the Midwest, easily
it’s at 700 millibars
just shows increasing precipitation
clouds are going to be increasing for the next 24 hours in western Ohio
Oh, what level are you?
850
they’re kind of in a little bit of a bubble high
so 0Z looking at low cloud cover
no precipitation at all

of this segmentation scheme from the meteorology domain. As the illustration shows, the length of each utterance can vary from a single word(s) (for example, “850”) to fairly lengthy (for example, “clouds are going to be increasing for the next 24 hours in western Ohio”). In this domain, we are generally interested in the types of mental operation participants perform on the visualizations, such as reading off information, transforming information (spatially or otherwise), making comparisons, and the like, and these map well to our segmentation scheme. The utterances “it’s at 700 millibars” and “clouds are going to be increasing for the next 24 hours in western Ohio” would each be coded as one “read-off information” event, even though they differ quite a bit in length. Of course, it would be possible to subdivide longer utterances further (for example, “clouds are going to be increasing // for the next 24 hours // in western Ohio”); however, according to our coding scheme this is still one read-off information event, which would now span three utterances. Having a single event span more than one utterance makes data analysis harder. The two most important guidelines in dividing protocols are (1) consistency and (2) making sure that the segments map readily to the codes to be applied.

Although transcription cannot be automated, a number of video analysis software programs are available to aid in data analysis, and some of these programs allow protocols to be transcribed directly into the software. We have not found the perfect protocol analysis program yet, though we have used MacShapa, Transana, and Noldus Observer. Looking ahead and deciding what kind of analysis software will be used will also save a great deal of time by reducing the need for further processing of transcripts in order to successfully import them to the video analysis software. Another option that we have successfully used is to transcribe protocols in a spreadsheet program, such as Excel, with each segment on a different row. We then set up columns for our different coding categories and can easily perform frequency counts, create pivot tables, and import the data into a statistical analysis program. This method works very well if the video portion of the protocol is irrelevant; however, in most cases we are interested not only in what people are saying, but what they are doing (and looking at) while they are saying it. In this case, using the spreadsheet method is disadvantageous, because the transcription cannot easily be aligned with the video. Video analysis software has the advantage that video timestamps are automatically entered at each line of transcription, thus facilitating synchronizing text and video.

Once data have been transcribed and segmented, they are ready for coding. Unless the research question can be answered by some kind of linguistic coding (for example, counting the number of times a certain word or family of words is used), a coding scheme must be developed that maps to the cognitive processes of interest. Establishing and implementing an effective and reliable coding scheme lies at the heart of verbal protocol analysis, and this is usually the most difficult and time consuming part of the whole process. In some cases, researchers will have strong a priori notions about what to look for. However, verbal protocol analysis is a method that is particularly useful in exploratory research, and in these cases, researchers may approach the protocol with only a general idea.

Regardless of the stage of research, coding schemes must frequently be elaborated or refined to match the circumstances of a given study.

To illustrate this point, consider a study we performed to determine how meteorologists handle uncertainty in the forecasting task and in the visualizations they use. Our hypothesis was that when uncertain, they would use more spatial transformations than when certain. Our coding scheme for spatial transformations was already in place. However, we had several possible ways that we could code uncertainty. One option was to code it linguistically. We developed a range of linguistic pointers to uncertainty, such as markers of disfluency (um, er, and so forth), hedging words (for example, “sort of,” “maybe,” and “somewhere around”), explicit statements of uncertainty (for example, “I have no idea” and “what’s that?”), and the like. Although this scheme proved highly reliable, with very little disagreement between coders, it turned out to be less useful for the research question we were trying to answer. A forecaster could be highly *uncertain* overall, but because he or she had no doubt about that uncertainty, the utterances contained few, if any, linguistic markers of uncertainty. Table 18.2 shows an example of this mismatch. The participant was trying to forecast precipitation for a particular period, and the two models (ETA and Global Forecast System [GFS]) she was using disagreed. Model disagreement is a major source of uncertainty for forecasters, so we were confident that at a more “global” level, the participant was highly uncertain at this juncture. Her uncertainty was also supported by her final utterance (that the forecast was going to be hard). However, the

Table 18.2. Mismatch between Linguistic and Global Coding Schemes: Uncertain Episode but Certain Utterances

Utterance	Linguistic Coding of Uncertainty	Global Coding of Uncertainty
OK, now let’s compare to the ETA . . . and they differ Interesting Well, they don’t differ too much they differ on when the precipitation’s coming Let’s go back to 36 hours OK, 42, 48, 54 Yeah, OK, so they have precip coming in 48 hours from now Let me try to go back to GFS and go to 48 hours and see what they have Well, OK, well they don’t differ Well, yeah, cause they’re calling for more so maybe this is gonna be a hard precip forecast	No linguistic uncertainty expressed in any utterance	Model disagreement indicates overall uncertainty (explicitly articulated in last utterance at left)

coding scheme to a small portion of the data. At this stage, we usually find problems coding specific utterances. It is important to keep very careful, detailed notes during this phase as to the nature of the coding problem. We can then look for patterns in the problematic segments, which often leads to subdividing or otherwise revising a particular code.

There are two equally important goals to keep in mind in developing a coding scheme. First, the scheme must be *reliable*; that is, it should be sufficiently precisely delineated that different coders will agree on the codes applied to any given utterance. Second, the scheme must be *useful*; that is, it must be at the right “level” to answer the question of interest. The linguistic scheme for coding uncertainty described above was highly reliable, but did not answer our need. The episodic coding scheme captured the relevant level of uncertainty, but required several refinements before we were confident that it was reliable. First, we had to obtain agreement on our division of the protocols into episodes, and then we had to obtain agreement on our coding of each episode as certain, uncertain, or mixed.

Obtaining agreement, or establishing *inter-rater reliability*, is essential in order to establish the validity of a coding scheme (see Cohen, 1960, for a discussion of establishing agreement). Once a coding scheme appears to be viable, the next step is to have two independent coders code the data and then compare their results. Obviously, coders must be well trained prior to this exercise. We establish a set of training materials, based on either a subset of the data, or, preferably, from a different dataset to which the same coding scheme can be applied. Using a different dataset reduces bias when it comes to the final coding of the data, because the coders are not prejudiced by previously having seen and worked with—and having discussed—the transcripts. The training materials consist of a set of written instructions describing each of the codes and a set of examples and nonexamples of the code. We find examples that are as clear and obvious as possible, but we also use “rogue” examples that meet some aspect of the coding criteria, but are nonetheless *not* illustrations of the code. These “near-misses” help define the boundaries of what does and does not count as a particular example of a code. Beside each example or nonexample, we provide an explanation of the coding decision. Once the coders believe they understand the coding scheme, they code a small subset of the data, and we meet to review the results. Once again, it is important for the coders to keep detailed notes about issues they encounter during this process. The more specific the coder can be, the more likely the fundamental cause of the coding problem can be addressed and resolved.

Initially, we look simply at points of agreement and disagreement in order to get a general sense of how consistently the codes are being applied. Although it is tempting to ignore points of agreement, we find it useful to look at individual coders’ reasoning in these instances in order to make sure that agreement is as watertight as possible. Points of disagreement, of course, provide a wealth of useful information in refining and redefining the coding scheme. Often, resolution is simply a matter of providing a crisper definition of the code. Sometimes, however, resolution involves subdividing one coding category into two or collapsing

two separate codes into one when the distinction proves blurry or otherwise not useful. The purpose of this stage is to identify and resolve the coders' uncertainties about the coding scheme in order that inter-rater reliability (IRR) might be established.

In order to obtain IRR, one coder codes the entire dataset, and the second coder codes a subset of the data. How much double-coding must be done depends in some measure on the nature of the transcripts and the coding scheme. The more frequently the codes occur, the less double-coding is necessary—usually 20–25 percent of the data is sufficient. For codes that occur rarely, a larger portion of the data must be double-coded in order to have sufficient instances of the code. If a code is infrequent, it is less likely to occur in any given subset of the data, and coders are more likely to be in agreement. However, they are agreeing that the code did *not* occur, which does not speak to the consistency with which they identify the code. After both coders have completed this exercise, we investigate their level of agreement.

There are two approaches to establishing agreement. The first is simply to count the number of instances in which coders agreed on a given code and calculate the percent agreement. This approach, however, generally results in an inflated measure of agreement, because it does not take into account the likelihood of agreement by chance. A better measure of agreement is Cohen's kappa (Cohen, 1960). This involves constructing a contingency table and marking the number of times the coders agreed on the code's occurrence or nonoccurrence² and the number of instances in which each coder said the code occurred but the other coder did not. Even if codes are scarce and the majority of instances fall into the "no-occurrence" cell, Cohen's kappa makes the appropriate adjustment for agreement by chance. Cohen's kappa is easily calculated as follows:

$$\kappa = (Po - Pc)/(1 - Pc)$$

Po is the proportion of observed agreement, and Pc is the proportion of agreement predicted by chance. Some researchers define poor reliability as a kappa of less than 0.4, fair reliability as 0.4 to 0.6, good reliability as 0.6 to 0.8, and excellent as greater than 0.8. We believe that kappa values lower than 0.6 indicate unacceptable agreement and that in such cases the coding scheme must undergo further revision.

A low value for kappa need not be considered a deathblow to the coding scheme. Most coding schemes go through a cycle of construction, application, evaluation, and revision several times before they can be considered valid. Although the process can be tedious, it is a vital part of using verbal protocols, because of the dangers of subjectivity and researcher bias raised earlier in this chapter. It should also be noted that coding schemes will most likely never achieve perfect agreement, nor is perfect agreement necessary. When coders disagree, there are two options. First, all disagreements can be excluded from

²Nonoccurrence is coded in addition to occurrence in order for coding to be symmetrical. Overlooking agreement about nonoccurrence can lead to a lower estimate of agreement than is warranted. This is especially an issue if instances of a code are relatively rare.

analysis. The advantage of this approach is that the researcher can be confident that the included data are reliably coded. The disadvantage is obviously that some valuable data are lost. A second option is to attempt to resolve agreements by discussion. Often, for example, one coder may simply have misunderstood some aspect of the transcript; in such cases, the coder may decide to revise his or her original coding. If, after discussion, coders still disagree, these instances will have to be excluded from analysis. However, the overall loss of data using this method will be less than if disagreements are automatically discarded. Either method is acceptable, although it is important to specify how disagreements were handled when presenting the results of one's research. After IRR is established and all the data are coded, the analysis proceeds as with any other research, using quantitative statistical analyses or qualitative analytical methods, depending on the research questions and goals.

PRESENTING RESEARCH RESULTS

There are some specific issues in presenting research based on verbal protocol analysis. Although the use of verbal protocols has been established as a sound method and has therefore gained increased acceptance in the last few years, such research may nevertheless be difficult to publish. Some reviewers balk at the generally small sample sizes. One approach to the sample size issue is to conduct a generalizability analysis (Brennan, 1992), which can lay to rest concerns that variability in the data is due to idiosyncratic individual differences rather than to other, more systematic factors.

Another concern is the possibility of subjectivity in coding the data and the consequential danger that the results will be biased by the researcher's predilection for drawing particular conclusions. There are two ways to address this concern. The first is to take extreme care not only in establishing IRR, but also in describing how the process was conducted, making sure, for example, that a sufficient proportion of the data was double-coded and that this is reported. Specifying how coders were trained, how coders worked independently, without knowledge of each other's codings, and how disagreements were resolved can also allay this concern. The second action is to provide very precise descriptions of the coding scheme, with examples that are crystal clear, so that the reader feels that he or she fully understands the codes and would be able to apply them. Although this level of detail consumes a great deal of space, many journals now have Web sites where supplementary material can be posted. When describing a coding scheme, good examples may be worth many thousands of words!

Even with these safeguards, research using verbal protocols may still be difficult to publish, especially if it is combined with an *in vivo* methodology. *In vivo* research is naturalistic in that it involves going into the environment in which participants normally work and observing them as they perform their work. The stumbling block to publication is that the research lacks experimental control. However, the strength of the method is precisely that it is less likely to change people's behavior, as imposing experimental controls may do. We have found

that an excellent solution to this vicious cycle is to follow our *in vivo* studies with controlled laboratory experiments that test the conclusions we draw from the *in vivo* work. In one project, we observed several expert scientists in a variety of domains as they analyzed their scientific data. Our observations suggested that they used a type of mental model we called “conceptual simulation” and that they did so especially when the hypothesis involved a great deal of uncertainty. However, our conclusions were correlational; in order to test whether the uncertainty was causally related to conceptual simulation, we conducted an experimental laboratory study in which we manipulated the level of participants’ uncertainty. In this type of follow-up study, it is especially important to design a task and materials that accurately capture the nature of the original task in order to elicit behavior that is as close as possible to that observed in the natural setting.

CONCLUSION

We have found verbal protocol analysis an invaluable tool in our studies of cognitive processes in several domains, such as scientific reasoning, analogical reasoning, graph comprehension, and uncertainty. Our participants have ranged from undergraduates in a university subject pool, to journeymen meteorology students, to experts with advanced degrees and decades of experience. Although we have primarily used verbal protocol analysis as a research tool, we believe that it can be used very effectively in applied settings, such as the design, development, and evaluation of virtual environments for training, where it is important to discern participants’ cognitive processes.

There are other methods of data collection that also require participants’ verbal input, such as structured and unstructured interviews, knowledge elicitation, and retrospective protocols. These methods assume that people have accurate access to their own cognitive processes; however, this is not necessarily the case. Consider, for example, the difficulty of explaining how to drive a stick-shift vehicle, compared with actually driving it. People frequently have faulty recall about their actions and motives or may describe what they *ought*—or were taught—to do rather than what they *actually* do. In addition, responses may be at a much higher level of detail than meets the researcher’s interest, or they may skip important steps in the telling that they would automatically perform in the doing. Although in some interview settings, such as semistructured interviews, the interviewer has the opportunity to probe the participant’s answers in more depth, unless the researcher already has some rather sophisticated knowledge of the task, the questions posed may not elicit the desired information.

The strength of concurrent verbal protocols is that they can be considered a reflection of the actual processes people use *as they perform a task*. Properly done, verbal protocol analysis can provide insights into aspects of performance that might otherwise remain inside a “black box” accessible only by speculation. Although they are time consuming to collect, process, and analyze, we believe that the richness of the data provided by verbal protocols far outweighs the costs. Ultimately, as with all research, the purpose of the research and the resources

available to the researcher will determine the most appropriate method; however, the addition of verbal protocol analysis to a researcher's repertoire will open the door to a potentially very productive source of data that invariably yield interesting and often surprising results.

REFERENCES

- Austin, J., & Delaney, P. F. (1998). Protocol analysis as a tool for behavior analysis. *Analysis of Verbal Behavior, 15*, 41–56.
- Brennan, R. L. (1992). *Elements of generalizability theory* (2nd ed.). Iowa City, IA: ACT Publications.
- Chi, M. T. H. (1997). Quantifying qualitative analyses of verbal data: A practical guide. *The Journal of the Learning Sciences, 6*(3), 271–315.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement, 20*, 37–46.
- Ericsson, K. A. (2006). Protocol analysis and expert thought: Concurrent verbalizations of thinking during experts' performance on representative tasks. In K. A. Ericsson, N. Charness, P. J. Feltovich, & R. R. Hoffman (Eds.), *The Cambridge handbook of expertise and expert performance* (pp. 223–241). New York: Cambridge University Press.
- Ericsson, K. A., & Simon, H. A. (1993). *Protocol analysis: Verbal reports as data* (2nd ed.). Cambridge, MA: MIT Press.
- Nisbett, R. E., & Wilson, T. D. (1977). Telling more than we can know: Verbal reports on mental processes. *Psychological Review, 84*, 231–259.
- van Someren, M. W., Barnard, Y. F., & Sandberg, J. A. C. (1994). *The think aloud method: A practical guide to modeling cognitive processes*. London: Academic Press.