

Pragmatic effects on reference resolution in a collaborative task: evidence from eye movements

Joy E. Hanna^{a,*}, Michael K. Tanenhaus^b

^a*Department of Psychology, State University of New York at Stony Brook,
SUNY Stony Brook, Stony Brook, NY 11794, USA*

^b*University of Rochester, Rochester, NY, USA*

Received 5 November 2002; received in revised form 10 September 2003; accepted 7 October 2003

Abstract

In order to investigate whether addressees can make immediate use of speaker-based constraints during reference resolution, participant addressees' eye movements were monitored as they helped a confederate cook follow a recipe. Objects were located in the helper's area, which the cook could not reach, and the cook's area, which both could reach. Critical referring expressions matched one object (helper's area) or two objects (helper's and cook's areas), and were produced when the cook's hands were empty or full, which defined the cook's reaching ability constraints. Helper's first and total fixations showed that they restricted their domain of interpretation to their own objects when the cook's hands were empty, and widened it to include the cook's objects only when the cook's hands were full. These results demonstrate that addressees can quickly take into account task-relevant constraints to restrict their referential domain to referents that are plausible given the speaker's goals and constraints.

© 2003 Cognitive Science Society, Inc. All rights reserved.

Keywords: Psychology; Communication; Language understanding; Human experimentation; Reference resolution; Perspective; Common ground; Conversation

1. Introduction

In conversation, the most basic arena of language use (Clark, 1992), partners must coordinate their individual knowledge and actions in order to communicate successfully. However, we know little about the temporal grain with which this joint activity affects moment-by-moment

* Corresponding author. Tel.: +1-631-632-4173; fax: +1-631-632-7876.

E-mail address: jhanna@sunysb.edu (J.E. Hanna).

language comprehension and production. This issue is important for both psycholinguistic models of real-time language processing and for computational models of conversational dialog (for review see Allen et al., 2001).

Empirical studies of on-line processing have tended to ignore the effects of being engaged in conversation (cf. Clark, 1997), focusing instead on core processes that are often viewed as rapid, encapsulated, and autonomous. In addition, only with the relatively recent advent of head-mounted eye-trackers have researchers been able to explore the fine-grained time course of language processing in controlled, yet relatively natural conversational settings (Henderson & Ferreira, 2004; Tanenhaus, Spivey-Knowlton, Eberhard, & Sedivy, 1995). The current report presents an important step in combining research in on-line language comprehension with research in face-to-face interactive conversation. We monitored eye movements to examine how a speaker's goals and constraints in a collaborative task affect an addressee's initial interpretation of utterances containing definite noun phrases.

We investigated definite reference for two reasons. First, the resolution of referential expressions, especially those involving definite noun phrases, has long been viewed as a central component of language processing, and it interacts with other basic processes such as syntactic ambiguity resolution. Second, definite noun phrases can be used to refer to or introduce uniquely identifiable entities, but only with respect to a particular context or *referential domain of interpretation* (e.g., Barwise, 1989; Searle, 1969). Without a restricted domain (e.g., the set of objects in a kitchen workspace), the set of potential referents for a definite noun phrase (e.g., *the bowl*) would be unlimited.

We know that what counts as a referential domain for an addressee is dynamically updated taking into account both the unfolding utterance and addressee-based task-specific pragmatic constraints (Chambers, Tanenhaus, Eberhard, Filip, & Carlson, 2002). However, a speaker's choice of referential expression is likely to be influenced by her own task-specific constraints (Horton & Keysar, 1995), which may shift as a function of task-relevant goals. For an addressee to take into account these constraints, he would have to dynamically modify his referential domain by taking into account the goals of the speaker. To date, there is no direct evidence about whether or not addressees can do this rapidly enough to influence the earliest moments of reference resolution. Moreover, it has been argued by Keysar and colleagues that initial reference resolution is primarily egocentric, possibly to minimize the resource demands associated with constructing representations of mutual belief, including another person's pragmatic constraints (Keysar, Barr, Balin, & Brauner, 1996, 2000; but cf. Hanna & Tanenhaus, *in press*; Hanna, Tanenhaus, & Trueswell, 2003; Nadig & Sedivy, 2002). However, these studies explored whether addressees can ignore potential referents that the speaker cannot see or does not know about in a visual display, a methodology which could make the addressee's own egocentric constraints easier to use or more salient. Addressees might be able to take a speaker's pragmatic constraints into account more quickly if they arise from task-specific constraints, evidence for the constraints is perceptually based, and monitoring and using the constraints is integral to the goals of the task.

We investigated whether an addressee could take into account the speaker's pragmatic constraints quickly enough to influence the initial domain of interpretation for a definite noun phrase by having participants follow a recipe in a cooking simulation. A (female) confederate speaker (C) played the role of the cook who read a recipe, and participant addressees played

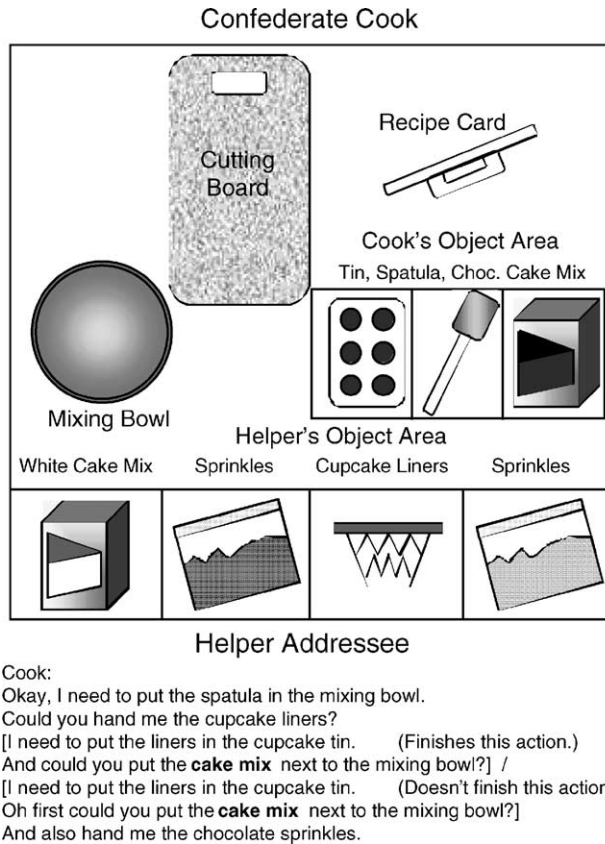


Fig. 1. Example display and critical instruction sequences for the two objects conditions. In the one object conditions the chocolate cake mix was replaced with frosting. This recipe was for cupcakes.

the role of the helper (**H**). Reading a recipe aloud is quite common during cooking, thus **C**'s utterances could be scripted without appearing unnatural, and instructions could be directed to herself or to **H** to move and manipulate objects in a particular order.

Fig. 1 shows a schematic of the set-up and a sample instruction sequence. **C** and **H** were seated across a table, and had spatially separate sets of objects for which they were responsible. **H** was told that most of the time he would be asked to move objects from his own area because **C** clearly couldn't reach them, but that sometimes when **C** was in the middle of something and had her hands full, he might be asked to move an object from **C**'s area, which **H** could reach.

The critical instructions to **H** manipulated two factors. The first was whether only one object on the table matched the definite referring expression, in **H**'s object area, or whether two objects matched, one in **H**'s area and one in **C**'s area. In order to create definite references that could apply equally well to either of two objects, modified names were initially used but the critical referring expression never included the modifying adjective. For example, in Fig. 1 there is a white cake mix in **H**'s object area and a chocolate cake mix in **C**'s object area, and the critical instruction always referred to *the cake mix*. The second factor was whether or not **C** had her hands full when she produced the instruction. On those trials, **C** began but did not finish the

preceding instruction, which required her to actively manipulate two to three other objects, such as to put liners in the cupcake tin. (**C** wore mirrored sunglasses and remained relatively still to avoid providing non-linguistic cues for reference resolution.)

The ability of the cook to reach the various objects provides a potential domain restriction, one that requires the helper to take into account information about her pragmatic constraints. When **C**'s hands are empty, **H** should only consider objects in his own area as possible referents. If **H** can take **C**'s empty hands and reaching ability into account, this restriction should persist even when there is another object that matches the referring expression in **C**'s area, since **C** would reach for that object herself if that was the intended referent. Examining Fig. 1, if **Hs** are able to restrict the domain of interpretation to their own area when **C**'s hands are empty, there should be no interference from the chocolate cake mix when **C** refers to *the cake mix*, and **Hs** should not ask which one was intended. In contrast, when **C**'s hands are full, the domain should widen to include the objects in **C**'s area, since **C** can no longer reach those objects on her own. **Hs** should consider the chocolate cake mix at least as often as the white cake mix, perhaps asking which one was meant. Importantly, this change in the domain of interpretation requires the helper to both attend to and understand, from the cook's perspective, which if any of the objects the cook can reach at the moment she utters the referring expression.

2. Method

2.1. Participants

Twelve native English speakers were paid for their participation. The confederate was a trained undergraduate research assistant.

2.2. Materials and design

The experiment took place on a 40 in. square table. **H**'s object area contained four spaces. **C**'s object area contained three spaces and was located equally distant from **C** and **H**. A cutting board was always present in front of **C**. On most trials, either a mixing bowl or a burner was located next to the cutting board. **C** had a vertical mount for the recipe cards, placed so that she could easily see them but **H** could not. Recipes were groups of instructions for doing such things as making cupcakes, tea, and bread. The objects were common kitchenware and food items, such as bowls, measuring cups, spices, and pasta. Critical objects were modified along either a kind or size dimension, for example either skim or 2% milk, or a big or small spatula. The set of objects present, and the subset and sequence of object movements and manipulation, were arbitrary so that references were not predictable.

Scripts contained five instructions: two initial instructions, one to **C** and one to **H**; a potential occupying action directed at **C**, which could either be finished or not finished before **C** continued with the critical instruction; the critical instruction to **H** to move an object to a particular location; and, a final instruction to either **C** or **H**. For the hands empty conditions, **Hs** were asked *And could you put the* [unmodified object name] *next to/on/in the cutting board/burner/mixing bowl*. For the hands full conditions, the same preamble was used, except we added *Oh, and*

first could you to make the instruction sound more natural and to draw attention to the fact that **C** was in the middle of something. The onset of the unmodified critical object name (e.g., *cake mix*) was always different than the other object names in the display.

Twelve object groups were combined with the recipes to create 12 critical items. By the beginning of the potential occupying action there was always one object in **C**'s object area and three objects in **H**'s object area. For each item there were two displays: a display with only one object that matched the critical referring expression and a display with two objects that matched. In the one object conditions, the single critical object was located in **H**'s object area, and the object in **C**'s area was unrelated; in the two object conditions, one of the critical objects was in **H**'s object area and one was the object in **C**'s area.

Within each experimental item, the four non-critical instructions served as fillers, counterbalancing object locations and the person to whom the instructions were directed. In addition to the (one or two) critical objects which varied along a kind or size dimension, items contained another set of two objects that varied along one of these dimensions, and references to these objects included the modifying adjective.

In order to reinforce **C**'s preference to move an object from her own object area when she could, two filler items matched the hands empty/two objects condition but included an instruction that was "mistakenly" directed at **H** and then repaired by **C**. For example, in the pumpkin pie filler, there was a large pie tin in **C**'s area and a small pie tin in **H**'s area. **C**, with her hands empty, began the third instruction *And could you put the pie tin*, then noticed her mistake and corrected herself, continuing with *actually I can put the large pie tin on the cutting board*. To counteract the underspecified references to a critical object in **C**'s area when her hands were full, two fillers matched the hands full/two objects condition and included fully specified references.

2.3. Procedure

Participants were told they would follow recipes that simulated real ones, and that their role was to be the cook's helper. The cook was introduced as a lab assistant who was uninformed about the experiment. It was explained that on each trial both they and **C** would get some objects in their areas, and that in order to follow the recipe these objects needed to be moved to different locations in a particular order. **Hs** were told that they would have to help **C** out either by moving the objects that **C** couldn't reach or by moving one of **C**'s own objects if she couldn't do it herself.

H's eye movements were monitored using an Applied Sciences Laboratories E5000 eye-tracker. On each trial, the experimenter placed objects on the table and handed **C** a recipe card. **C** named the recipe, asked the experimenter for either the mixing bowl or the burner if necessary, and listed the objects in her own area and in **H**'s area, looking between the recipe card and the objects to make it clear that she was becoming familiar with them. **C** then read aloud the recipe instructions. **C** always completed her object manipulation before moving on to the next instruction, except in the hands full conditions where she glanced at the recipe card and gave the critical instruction during the action.

In the two object conditions, the target was the one in **H**'s area when **C**'s hands were empty, and was designated as the one in **C**'s area when **C**'s hands were full. If **H** chose the wrong object

in these conditions, or asked which object was meant, **C** responded with the fully specified object name.

In order to appear natural, **C** varied non-critical instructions slightly. Critical instructions were always read exactly as written. To avoid providing referential cues with her eye gaze or body movements, **C** wore mirrored sunglasses, aimed her gaze and body towards the recipe card during all critical instructions, and kept her hands visible and still (either resting on the table or paused during object manipulation). **Hs** were told that the sunglasses were necessary because the experiment was concerned with verbal behavior.

Two practice trials familiarized **H** with the procedure, and with **C**'s ability to move an object from her own area when her hands were empty, and inability to do so when her hands were full. **Hs** were told that it was okay to ask the cook questions if they were confused at any point.

3. Results

A debriefing revealed that participants thought the confederate was practiced at the task, but none thought she was reading something other than a normal recipe. Data were analyzed from the videotape records using an editing VCR with frame-by-frame control and synchronized video and audio channels. Fixations were scored by noting which location **Hs** were fixating, including an object or space in his own object area, an object or space in **C**'s object area, one of the locations on the table, or **C**'s face or hands. Fixations were scored beginning from the onset of the critical instruction and ending with the fixation prior to **H** initiating a final reaching movement to pick up an object, noting the onset of the critical object name. Two of the 144 trials were discarded because of track loss.

3.1. *Proportion of fixations over time*

Figs. 2 and 3 present the proportion of fixations to the objects and locations in the display in each condition in 33 ms intervals, beginning with the onset of the critical instruction. The proportions don't always sum to 1.0 because fixations to objects other than those in the object areas of the display or to the cook were not included. At the beginning of the critical instruction participants were most often (96% of the time) fixating on the cook's hands or face. The average onset of the critical instruction and the object name are indicated on each graph with the light vertical lines. The average duration of the instruction before the object name was 647 ms for the hands empty conditions, and a longer 1337 ms for the hands full conditions due to the attention-drawing preamble. The average duration of the noun was 415 ms in the hands empty conditions and 431 ms in the hands full conditions.

Examination of the figures provides clear evidence that **Hs** largely restricted their domain of interpretation to the objects in their own area when **C**'s hands were empty, and expanded their domain to include **C**'s objects only when her hands were full. In both of the one object conditions (Fig. 2), fixations to the target diverged quickly from other referents, within 100 ms after the offset of the noun. Since it takes about 150–200 ms for an eye movement to be programmed and launched (Matin, Shao, & Boff, 1993), these fixations were programmed during the object name. Note that in the hands full/one object condition, there were frequent

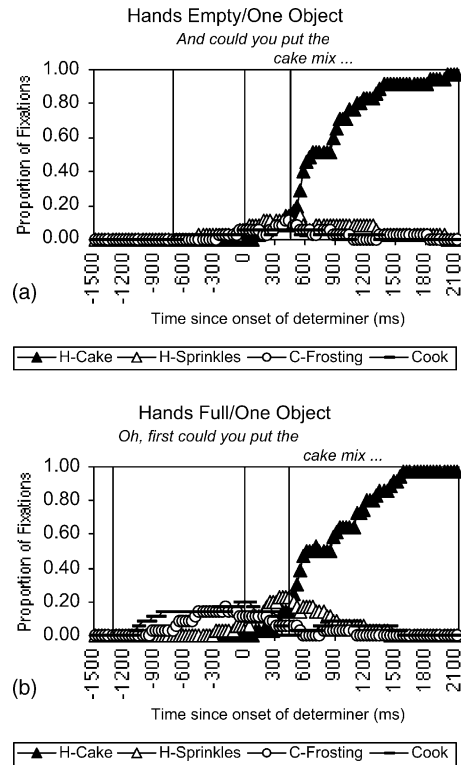


Fig. 2. Proportion of fixations to each object type over time in the hands empty/one object condition (a) and the hands full/one object condition (b), including the single target object in the helper's area (H-cake), any other object in the helper's area (H-sprinkles), the unrelated object in the cook's area (C-frosting), or the cook.

looks to **C** and **C**'s area during the preamble, indicating that **Hs** were anticipating that **C** might need help with something in her area. In the hands full/two objects condition (Fig. 3), **Hs** were equally likely to look at both object areas, with an initial bias towards **C**'s area, presumably because attention was focused on **C**. Crucially, the object in **C**'s area did not compete to this degree when **C**'s hands were empty; while there is a later rise in looks to **C**'s object, **Hs** show an immediate preference for their own object, just as they did in the one object conditions.

3.2. Total and first looks analyses

We evaluated this data pattern by conducting an ANOVA on the proportion of total and initial looks to critical objects in the four conditions. Because domain restriction in this task involves interpreting the referring expression with respect to either the set of objects in **H**'s area or the one object in **C**'s area, these analyses compared the percentage of looks to any of the objects in **H**'s area, looks to the one object in **C**'s area, and looks to any other object in the display (including **C**, the cutting board, the mixing bowl, or the burner), although overall the latter was a small percentage (3.5%) and had no significant effects.

For total looks to **H**'s area, there were main effects of hands and the number of objects, but these factors interacted significantly ($F(1, 8) = 21.18, p < .01, \text{MSE} = 0.01$). **Hs**

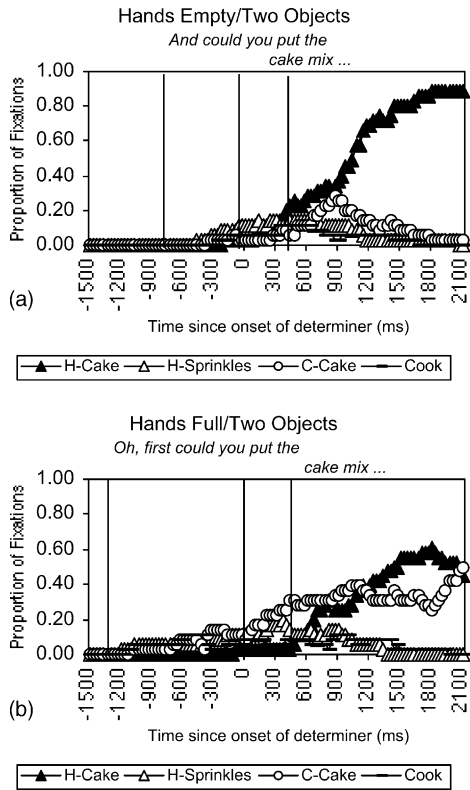


Fig. 3. Proportion of fixations to each object type over time in the hands empty/two objects condition (a) and the hands full/two objects condition (b), including the target object in the helper's area (**H**-cake), any other object in the helper's area (**H**-sprinkles), the critical object in the cook's area (**C**-cake), or the cook.

were equally likely to look at an object in their own area in the hands empty/one object (mean = 79%), hands full/one object (mean = 75%), and, crucially, the hands empty/two objects (mean = 71%) conditions, all P values $> .1$. It was only in the hands full/two objects condition that **H**s looked significantly less often to an object in their own area (mean = 45%) in comparison to when **C**'s hands were empty ($F(1, 8) = 25.60, p < .01, \text{MSE} = 0.02$) or to when there was only one object in the display ($F(1, 8) = 49.45, p < .001, \text{MSE} = 0.01$). The same pattern emerges for looks to the object in **C**'s area, with a significant interaction ($F(1, 8) = 28.37, p < .01, \text{MSE} = 0.01$). **H**s looked at the unrelated object in **C**'s area an equally small percentage of time when **C**'s hands were empty (mean = 10%) or full (mean = 9%) ($F(1, 8) = 0.04, p > .5, \text{MSE} = 0.01$), and in the two objects conditions looked more often at the critical object in **C**'s area when **C**'s hands were full (mean = 41%) in comparison to when they were empty (mean = 17%) ($F(1, 8) = 20.62, p < .01, \text{MSE} = 0.02$).

A similar pattern was found analyzing only the first look after the onset of the object name, indicating that the referential domain restriction was immediate. The interaction of hands and the number of objects was significant ($F(1, 8) = 8.76, p < .05, \text{MSE} = 0.02$). Simple effects revealed that the percentage of first looks to an object in **H**'s area in the hands empty/one object

condition (mean = 78%), the hands full/one object condition (mean = 89%), and the hands empty/two objects condition (mean = 68%) did not differ significantly, all P values $> .1$. **Hs** looked first at an object in their own area numerically less often in the hands full/two objects condition (mean = 56%) than in the hands empty/two objects condition, although this difference only approached significance ($F(1, 8) = 3.86$, $p < .08$, $MSE = 0.02$). First looks to the object in **C's** area showed a stronger pattern of effects. There was a significant main effect of objects, but this effect interacted with whether **C's** hands were empty or full ($F(1, 8) = 8.26$, $p < .05$, $MSE = 0.02$). There were no significant differences between the low percentage of first looks to **C's** area in the hands empty/one object (mean = 8%), hands full/one object (mean = 3%), and importantly hands empty/two object (mean = 18%) conditions, all P values $> .2$. There was, however, a significantly higher percentage of first looks to the object in **C's** area in the hands full/two object condition (mean = 36%) in comparison to both the hands full/one object condition ($F(1, 8) = 10.29$, $p < .05$, $MSE = 0.06$), and, crucially, the hands empty/two object condition ($F(1, 8) = 5.83$, $p < .05$, $MSE = 0.03$).

Finally, when **C's** hand were empty and there were two objects in the display, **Hs** never chose **C's** object, asking which one was meant 8% of the time, and choosing their own object 92% of the time. Planned comparisons showed that this percentage was no different from their 100% choice of the single target object in the hands empty/one object condition ($F(1, 8) = 1.80$, $p = .2$, $MSE = 0.02$), and was significantly greater than the 39% choice of their own object in the hands full/two objects condition ($F(1, 8) = 15.04$, $p < .01$, $MSE = 0.11$).

4. Discussion

The results present clear evidence that the referential domain of interpretation for the helper was rapidly modified by consideration of the cook's pragmatic constraints. When the cook's hands were empty, helpers immediately interpreted a request containing a definite noun phrase as referring to the referent in their own area, even when the environment included a second potential referent in the cook's area. Their own object was the single plausible referent given the cook's reaching ability and goals. In contrast, when the cook's hands were full, helpers considered the object in the cook's area as a potential referent, showing a preference for that interpretation and only excluding it after the onset of the object name in the one object conditions. Thus, the referential domain within which the addressee's uniqueness constraints were interpreted was immediately modulated by an awareness of the speaker's goals, as defined by the task-based constraints and signaled by a non-verbal cue.¹

These results are in line with increasing evidence for constraint-based models (e.g., MacDonald, 1994; Tanenhaus & Trueswell, 1995). From a constraint-based perspective, taking into account information that might shape the intentions of an interlocutor provides constraints that are integrated simultaneously and continuously with other lexical, structural, and discourse-based constraints to provide probabilistic evidence for alternative interpretations. We suggest that a constraint-based perspective can reconcile results like those presented here with demonstrations by Keysar and coworkers (e.g., Keysar et al., 2000) that addressees sometimes consider perceptually salient information even though that information is not available to the speaker. We propose that the degree to which an addressee will monitor the goals, intentions,

and likely knowledge of an interlocutor will vary with the salience and accessibility of that information and its importance for the goals of the addressee. In the current experiment, it was important for the addressee to monitor the speaker. Moreover, the speaker's abilities were embedded in the task constraints, reflective of her goals, and perceptually signaled by a simple cue. Taken together, these characteristics create the necessary conditions for strong effects of the speaker's perspective on the addressee's referential domain of interpretation.

While the particular cue available in this experiment was clearly made salient by the task, our finding is nonetheless of central theoretical importance for two reasons. First, it complements recent demonstrations that addressee-based pragmatic constraints can determine whether or not perceptually salient entities are included in the domain used to compute the uniqueness conditions of definite reference (e.g., Chambers et al., 2002). Second, and most importantly, goal-based domain restrictions, especially those signaled by relatively simple non-verbal cues, are likely to provide an important source of information that allows participants to coordinate their language without requiring resource demanding models of each other's mutual beliefs. Further exploring the effects of this type of information on reference resolution should contribute fundamentally to our understanding of language processing during conversation.

Note

1. The hands full conditions had an additional cue—the presence of the preamble in the critical instruction. The preamble was included to help the interruption of the object manipulation seem natural, and this linguistic cue appeared to primarily reinforce the non-verbal cue that the cook's hands were still occupied. The preamble did not slow down reference resolution in the hands full/one object condition, and participants were able to notice the emptiness of the cook's hands without an additional cue in the hands empty conditions. Therefore, we consider the primary pragmatic cue in this experiment to be the cook's reaching ability.

Acknowledgments

This material is based upon work supported by the National Institutes of Health under National Research Service Award 1F32MH1263901A1 and Grant No. HD-27206, and by the National Science Foundation under SBR Grant No. 95-11340, and Grants No. 0082602 and 9980013. We are very grateful to Katie Schack for being the confederate and for assistance in data coding. Some of the research in this paper was reported at the Fourteenth Annual CUNY Conference on Human Sentence Processing (March, 2001) in Philadelphia, PA.

References

- Allen, J. F., Byron, D. K., Dzikovska, M., Ferguson, G., Galescu, L., & Stent, A. (2001). Towards conversational human-computer interaction. *AI Magazine*, 22, 27–35.

- Barwise, J. (1989). *The situation in logic*. Stanford, CA: Center for the Study of Language and Information.
- Chambers, C. G., Tanenhaus, M. K., Eberhard, K. M., Filip, H., & Carlson, G. N. (2002). Circumscribing referential domains in real-time sentence comprehension. *Journal of Memory and Language*, 47, 30–49.
- Clark, H. H. (1992). *Arenas of language use*. Chicago: University of Chicago Press.
- Clark, H. H. (1997). Dogmas of understanding. *Discourse Processes*, 23, 567–598.
- Hanna, J. E., & Tanenhaus, M. K. (in press). The use of perspective during referential interpretation. In J. C. Trueswell & M. K. Tanenhaus (Eds.), *World situated language use: Psycholinguistic, linguistic, and computational perspectives on bridging the product and action traditions*. Cambridge, MA: MIT Press.
- Hanna, J. E., Tanenhaus, M. K., & Trueswell, J. C. (2003). The effects of common ground and perspective on domains of referential interpretation. *Journal of Memory and Language*, 49, 43–61.
- Henderson, J. M., & Ferreira, F. (Eds.). (2004). *The interface of language, vision, and action: Eye movements and the visual world*. New York: Psychology Press.
- Horton, W. S., & Keysar, B. (1995). When do speakers take into account common ground? *Cognition*, 59, 91–117.
- Keysar, B., Barr, D. J., Balin, J. A., & Brauner, J. S. (1996, November). *Common ground: An error correction mechanism in comprehension*. Paper presented at the 37th Annual Meeting of the Psychonomic Society, Chicago, IL.
- Keysar, B., Barr, D. J., Balin, J. A., & Brauner, J. S. (2000). Taking perspective in conversation: The role of mutual knowledge in comprehension. *Psychological Science*, 11(1), 32–37.
- MacDonald, M. C. (1994). Probabilistic constraints and syntactic ambiguity resolution. *Language and Cognitive Processes*, 9(2), 157–201.
- Matin, E., Shao, K. C., & Boff, K. R. (1993). Saccadic overhead: Information processing time with and without saccades. *Perception and Psychophysics*, 53(4), 372–380.
- Nadig, A. S., & Sedivy, J. C. (2002). Evidence of perspective-taking constraints in children's on-line reference resolution. *Psychological Science*, 13(4), 329–336.
- Searle, J. R. (1969). *Speech acts. An essay in the philosophy of language*. Cambridge: Cambridge University Press.
- Tanenhaus, M. K., Spivey-Knowlton, M. J., Eberhard, K. M., & Sedivy, J. E. (1995). Integration of visual and linguistic information in spoken language comprehension. *Science*, 268, 632–634.
- Tanenhaus, M. K., & Trueswell, J. C. (1995). Sentence comprehension. In J. Miller & P. Eimas (Eds.), *Handbook of cognition and perception*. San Diego, CA: Academic Press.