

LIVING UP TO EXPECTATIONS: COMPUTING EXPERT RESPONSES¹

Aravind Joshi and Bonnie Webber
Department of Computer and Information Science
Moore School/D2
University of Pennsylvania
Philadelphia PA 19104

Ralph M. Weischedel²
Department of Computer & Information Sciences
University of Delaware
Newark DE 19716

ABSTRACT

In cooperative man-machine interaction, it is necessary *but not sufficient* for a system to respond truthfully and informatively to a user's question. In particular, if the system has reason to believe that its planned response might mislead the user, then it must block that conclusion by modifying its response. This paper focusses on identifying and avoiding potentially misleading responses by acknowledging types of "informing behavior" usually expected of an expert. We attempt to give a formal account of several types of assertions that should be included in response to questions concerning the achievement of some goal (in addition to the simple answer), lest the questioner otherwise be misled.

1. Introduction]

In cooperative man-machine interaction, it is necessary *but not sufficient* for a system to respond truthfully and informatively to a user's question. In particular, if the system has reason to believe that its planned response might mislead the user to draw a false conclusion, then it must block that conclusion by modifying or adding to its response.

Such cooperative behavior was investigated in [5], in which a modification of Grice's *Maxim of Quality* - "Be truthful" - is proposed:

If you, the speaker, plan to say anything which may imply for the hearer something that you believe to be false, then provide further information to block it.

This behavior was studied in the context of interpreting certain definite noun phrases. In this paper, we investigate this revised principle as applied to responding to users' plan-related questions. Our overall aim is to:

1. characterize tractable cases in which the system as respondent (R) can anticipate the possibility of the user/questioner (Q) drawing false conclusions from its response and hence alter it so as to prevent this happening;
2. develop a formal method for computing the projected inferences that Q may draw from a

¹This work is partially supported by NSF Grants MCS 81-07200, MCS 83-05221, and IST 83-11400.

²At present visiting the Department of Computer and Information Science, University of Pennsylvania PA 19104.

particular response, identifying those factors whose presence or absence catalyzes the inferences;

3. enable the system to generate modifications of its response that can defuse possible false inferences and that may provide additional useful information as well.

In responding to any question, including those related to plans, a respondent (R) must conform to Grice's first *Maxim of Quantity* as well as the revised *Maxim of Quality* stated above:

Make your contribution as informative as is required (for the current purposes of the exchange).

At best, if R's response is not so informative, it may be seen as uncooperative. At worst, it may end up violating the revised *Maxim of Quality*, causing Q to conclude something R either believes to be false or does not know to be true: the consequences could be dreadful. Our task is to characterize more precisely what this expected informativeness consists of. In question answering, there seem to be several quite different types of information, over and beyond the simple answer to a question, that are nevertheless expected. For example,

1. When a task-related question is posed to an expert (R), R is expected to provide additional information that he recognizes as necessary to the performance of the task, of which the questioner (Q) may be unaware. Such response behavior was discussed and implemented by Allen [1] in a system to simulate a train information booth attendant responding to requests for schedule and track information. In this case, not providing the expected additional information is simply uncooperative: Q won't conclude the train doesn't depart at any time if R fails to volunteer one.
2. With respect to discussions and/or arguments, a speaker contradicting another is expected to support his contrary contention. Again, failing to provide support would simply be viewed as uncooperative [2, 3].
3. With respect to an expert's responses to questions, if Q expects that R would inform him of P if P were true, then Q may interpret R's silence regarding P as implying P is not true.³ Thus if R knows P to be true, his silence may lead to Q's being misled. This third type of expected informativeness is the basis for the potentially misleading responses that we are trying to avoid and that constitute the subject of this paper.

What is of interest to us is *characterizing* the Ps that Q would expect an expert R to inform him of, if they hold. Notice that these Ps differ from script-based expectations [8], which are based on what is taken to be the ordinary course of events in a situation. In describing such a situation, if the speaker doesn't explicitly reference some element P of the script, the listener simply assumes it is true. On the other hand, the Ps of interest here are based on normal cooperative discourse behavior, as set out in Grice's maxims. If the speaker doesn't make explicit some information P that the listener believes he would possess and inform the listener of, the listener assumes it is false.

In this paper, we attempt to give a formal account of a subclass of Ps that should be included (in addition to the simple answer) in response to questions involving Q's achieving some goal⁴ - e.g., "Can I

³This is an interactional version of what Reiter [13] has called the "Closed World Assumption" and what McCarthy [9] has discussed in the context of "Circumscription".

⁴A companion paper [6] discusses responses which may mislead Q into assuming some default which R knows not to hold. Related work [4] discusses providing indirect or modified responses to yes/no questions where a direct response, while truthful, might mislead Q.

drop CIS577?", "I want to enrol in CIS577?", "How do I get to Marsh Creek on the Expressway?", etc., lest that response otherwise mislead Q. In this endeavor, our first step is to specify that knowledge that an expert R must have in order to identify the Ps that Q would expect to be informed of, in response to his question. Our second step is to formalize that knowledge and show how the system can use it. Our third step is to show how the system can modify its planned response so as to convey those Ps. In this paper, Section 2 addresses the first step of this process and Sections 3 and 4 address the second. The third step we mention here only in passing.

2. Factors in Computing Likely Informing Behavior]

Before discussing the factors involved in computing this desired system behavior, we want to call attention to the distinction we are drawing between *actions* and *events*, and between the *stated goal* of a question and its *intended goal*. We limit the term *action* to things that Q has some control over. Things beyond Q's control we will call *events*, even if performed by other agents. While events may be *likely* or even *necessary*, Q and R nevertheless can do nothing more than wait for them to happen. This distinction between actions and events shows up in R's response behavior: if an action is needed, R can suggest that Q perform it. If an event is, R can do no more than inform Q.

Our second distinction is between the *stated goal* or "S-goal" of a request and its *intended goal* or "I-goal". The former is the goal most directly associated with Q's request, beyond that Q know the information. That is, we take the S-goal of a request to be the goal directly achieved by using the information.

Underlying the stated goal of a request though may be another goal that the speaker wants to achieve. This *intended goal* or "I-goal" may be related to the S-goal of the request in any of a number of ways:

- The I-goal may be the same as the S-goal.
- The I-goal may be more abstract than the S-goal, which addresses only part of the I-goal. (This is the standard goal/sub-goal relation found in hierarchical planning [14].) For example, Q's S-goal may be to delete some files (e.g., "How can I delete all but the last version of FOO.MSS?"), while his I-goal may be to bring his file usage under quota. This more abstract goal may also involve archiving some other files, moving some into another person's directory, etc.
- The S-goal may be an enabling condition for the I-goal. For example, Q's S-goal may be to get read/write access to a file, while his I-goal may be to alter it.
- The I-goal may be more general than the S-goal. For example, Q's S-goal may be to know how to repeat a control-N, while his I-goal may be to know how to effect multiple sequential instances of a control character.
- Conversely, the I-goal may be more specific than the S-goal - for example, Q's S-goal may be to know how to send files to someone on another machine, while his I-goal is just to send a particular file to a local network user, which may allow for a specialized procedure.

Inferring the I-goal corresponding to an S-goal is an active area of research [1, Carberry83, 10, 11]. We assume for the purposes of this paper that R can successfully do so. One problem is that the relationship that Q believes to hold between his S-goal and his I-goal may not actually hold: for example, the S-goal

may not fulfill part of the I-goal, or it may not instantiate it, or it may not be a pre-condition for it. In fact, the S-goal may not even be possible to effect! This failure, under the rubric "relaxing the appropriate-query assumption", is discussed in more detail in [10, 11]. It is also reason for augmenting R's response with appropriate Ps, as we note informally in this section and more formally in the next.

Having drawn these distinctions, we now claim that in order for the system to compute both a direct answer to Q's request and such Ps as he would expect to be informed of, were they true, the system must be able to draw upon knowledge/beliefs about

- the events or actions, if any, that can bring about a goal
- their enabling conditions
- the likelihood of an event occurring or the enabling conditions for an action holding, with respect to a state
- ways of evaluating methods of achieving goals - for example, with respect to simplicity, other consequences (side effects), likelihood of success, etc.
- general characteristics of cooperative expert behavior

The roles played by these different types of knowledge (as well as specific examples of them) are well illustrated in the next section.

3. Formalizing Knowledge for Expert Response

In this section we give examples of how a formal model of user beliefs about cooperative expert behavior can be used to avoid misleading responses to task-related questions - in particular, what is a very representative set of questions, those of the form "How do I do X?". Although we use logic for the model because it is clear and precise, we are not proposing theorem proving as the means of computing cooperative behavior. In Section 4 we suggest a computational mechanism. The examples are from a domain of advising students and involve responding to the request "I want to drop CIS577". The set of individuals includes not only students, instructors, courses, etc. but also states. Since events and actions change states, we represent them as (possibly parameterized) functions from states to states. All terms corresponding to events or actions will be underlined. For these examples, the following notation is convenient:

Q	the user
R	the expert
Sc	the current state of the student
RB(P)	R believes proposition P
RBQB(P)	R believes that Q believes P
<u>admissible(e(S))</u>	<u>event/action e can apply in state S</u>
<u>likely(a,S)</u>	<u>a is a likely event/action in state S</u>
holds(P,S)	P, a proposition, is true in S
want(x,P)	x wants P to be true

To encode the preconditions and consequences of performing an action, we adopt an axiomatization of STRIPS operators due to [Chester83, 7, 15]. The preconditions on an action being applicable are encoded using "holds" and "admissible" (essentially defining "admissible"). Namely, if c1, ..., cn are preconditions on an action a,

$\text{holds}(c1,s) \ \&\dots\ \&\text{holds}(cn,s) \Rightarrow \text{admissible}(a(s))$

a 's immediate consequences $p1, \dots, pm$ can be stated as

$\text{admissible}(a(s)) \Rightarrow \text{holds}(p1, a(s)) \ \&\dots\ \&\text{holds}(pm, a(s))$

A frame axiom states that only $p1, \dots, pm$ have changed.

$\neg(p=p1) \ \&\dots\ \&\neg(p=pm) \ \&\text{holds}(p,s) \ \&\text{admissible}(a(s)) \Rightarrow \text{holds}(a(s))$

In particular, we can state the preconditions and consequences of dropping CIS577. (h and n are variables, while C stands for CIS577.)

$\text{RB}(\text{holds}(\text{enrolled}(h, C, \text{fall}), n) \ \&\text{holds}(\text{date}(n) < \text{Nov16}, n) \Rightarrow \text{admissible}(\text{drop}(h, C)(n)))$

$\text{RB}(\text{admissible}(\text{drop}(h, C)(n)) \Rightarrow \text{holds}(\neg\text{enrolled}(h, C, \text{fall}), \text{drop}(h, C)(n)))$

$\text{RB}(\neg(p=\text{enrolled}(h, C, \text{fall})) \ \&\text{admissible}(\text{drop}(h, C)(n)) \ \&\text{holds}(p, n) \Rightarrow \text{holds}(p, \text{drop}(h, C)(n)))$

Of course, this only partially solves the frame problem, since there will be implications of $p1, \dots, pm$ in general. For instance, it is likely that one might have an axiom stating that one receives a grade in a course only if the individual is enrolled in the course.

Q 's S -goal in dropping CIS577 is not being in the course. By a process of reasoning discussed in [10, 11], R may conclude that Q 's likely intended goal (I -goal) is not failing it. That is, R may believe:

$\text{RBQB}(\text{holds}(\neg\text{fail}(Q, C), \text{drop}(Q, C)(Sc)))^5$

$\text{RB}(\text{want}(Q, \neg\text{fail}(Q, C)))$

What we claim is: (1) R must give a truthful response addressing at least Q 's S -goal; (2) in addition, R may have to provide information in order not to mislead Q ; and (3) R may give additional information to be cooperative in other ways. In the subsections below, we enumerate the cases that R must check in effecting (2). In each case, we give both a formal representation of the additional information to be conveyed and a possible English gloss. In that gloss, the part addressing Q 's S -goal will appear in normal type, while the additional information will be underlined>.

For each case, we give two formulae: a statement of R 's beliefs about the current situation and an axiom stating R 's beliefs about Q 's expectations. Formulae of the first type have the form $\text{RB}(P)$. Formulae of the second type relate such beliefs to performing an informing action. They involve a statement of the form

$\text{RB}[P] \Rightarrow \text{likely}(i, Sc)$,

where i is an informing act. For example, if R believes there is a better way to achieve Q 's goal, R is likely to inform Q of that better way. Since it is assumed that Q has this belief, we have

$\text{QB}(\text{RB}[P] \Rightarrow \text{likely}(i, Sc))$.

⁵It will also be the case that $\text{RBQB}(\text{admissible}(\text{drop}(Q, C)(Sc)))$ if Q 's asks "How can I drop CIS577?", but not if he asks "Can I drop CIS577?". In the latter case, Q must of course believe that it may be admissible, or why ask the question. In either case, R 's subsequent behavior doesn't seem contingent on his beliefs about Q 's beliefs about admissibility.

where we can equate "Q believes *i* is likely" with "Q expects *i*." Since R has no direct access to Q's beliefs, this must be embedded in R's model of Q's belief space. Therefore, the axioms have the form (modulo quantifier placement)

$$RBQB(RB[P] \Rightarrow \text{likely}(i, Sc)).$$

An informing act is meant to serve as a command to a natural language generator which selects appropriate lexical items, phrasing, etc. for a natural language utterance. Such an act has the form *inform-that*(*R, Q, P*) R informs Q that P is true.

3.1. Failure of enabling conditions

Suppose that it is past the November 15th deadline or that the official records don't show Q enrolled in CIS577. Then the enabling conditions for dropping it are not met. That is, R believes Q's S-goal cannot be achieved from Sc.

$$[1] RB(\text{want}(Q, \neg \text{fail}(Q, C)) \ \& \ \neg \text{admissible}(\text{drop}(Q, C)(Sc)))$$

Thus R initially plans to answer "You can't drop CIS577". Beyond this, there are two possibilities.

3.1.1. A way

If R knows another action *b* that would achieve Q's goals (cf. formula [2]), Q would expect to be informed about it. If not so informed, Q may mistakenly conclude that there is no other way. Formula [3] states this belief that R has about Q's expectations.

$$[2] RB((\exists b) [\text{admissible}(b(Sc)) \ \& \ \text{holds}(\neg \text{fail}(Q, C), b(Sc))])$$

$$[3] RBQB(RB[\text{want}(Q, \neg \text{fail}(Q, C)) \ \& \ \neg \text{admissible}(\text{drop}(Q, C)(Sc))] \ \& \ RB[(\exists b) [\text{admissible}(b(Sc)) \ \& \ \text{holds}(\neg \text{fail}(Q, C), b(Sc))]]) \Rightarrow \text{likely}(\text{inform-that}(R, Q, (\exists b) [\text{admissible}(b(Sc)) \ \& \ \text{holds}(\neg \text{fail}(Q, C), b(Sc)) \ \& \ \text{can}(Q, b), Sc)]))$$

R's full response is therefore "You can't drop 577; *you can b*." For instance, *b* could be changing status to auditor, which may be performed until December 1.

3.1.2. No way

If R doesn't know of any action or event that could achieve Q's goal (cf. [4]), Q would expect to be so informed. Formula [5] states this belief about Q's expectations.

$$[4] RB(\neg(\exists a) [\text{admissible}(a(Sc)) \ \& \ \text{holds}(\neg \text{fail}(Q, C), a(Sc))])$$

$$[5] RBQB(RB(\text{want}(Q, \neg \text{fail}(Q, C)) \ \& \ \neg(\exists a) [\text{admissible}(a(Sc)) \ \& \ \text{holds}(\neg \text{fail}(Q, C), a(Sc))]) \Rightarrow \text{likely}(\text{inform-that}(R, Q, \neg(\exists a) [\text{admissible}(a(Sc)) \ \& \ \text{holds}(\neg \text{fail}(Q, C), a(Sc))]), Sc))$$

To say only that Q cannot drop the course does not exhibit expert cooperative behavior, since Q would be uncertain as to whether R had considered other alternatives. Therefore, R's full response is "You can't drop 577; *there isn't anything you can do to prevent failing*."

Notice that R's analysis of the situation may turn up additional information which a cooperative expert

could provide that does not involve avoiding misleading Q. For instance, R could indicate enabling conditions that prevent there being a solution: suppose the request to drop the course is made after the November 15th deadline. Then R would believe the following, in addition to [1]

RB(holds(enrolled(Q,C,fall),Sc) & holds(date(Sc)>Nov15,Sc))

More generally, we need a schema such as the following about Q's beliefs:

RBQB(RB[want(Q,¬fail(Q,C))
& (holds(P1, S) &...& holds(Pn, S) ⇒ admissible(a(S))
& (¬holds(Pi, S), for some Pi above))
⇒ likely(inform-that(R,Q,¬holds(Pi,S)),S))

In this case the response should be "You can't drop 577; *Pi isn't true.*" Alternatively, the language generator might paraphrase the whole response as, "if *Pi* were true, you could drop."

Of course there are potentially many ways to try to achieve a goal: by a single action, by a single event, or by an event and an action, ... In fact, the search for a sequence of events or actions that would achieve the goal may consider many alternatives. If all fail, it is far from obvious which blocked condition to notify Q of, and knowledge is needed to guide the choice. Some heuristics for dealing with that problem are given in [12].

3.2. An nonproductive act

Suppose the proposed action does not achieve Q's I-goal, cf. [6]. For example, dropping the course may still mean that failing status would be recorded as a WF (withdrawal while failing). R may initially plan to answer "You can drop 577 by ...". However, Q would expect to be told that his proposed action does not achieve his I-goal. Formula [7] states R's belief about this expectation.

[6] RB(¬holds(¬fail(Q,C), drop(Q,C)(Sc)) & admissible(drop(Q,C)(Sc)))

[7] RBQB(RB[want(Q,¬fail(Q,C)) & ¬holds(¬fail(Q,C),drop(Q,C)(Sc))
& admissible(drop(Q,C)(Sc))]
⇒ likely(inform-that(R,Q,
¬holds(¬fail(Q,C),drop(Q,C)(Sc)),Sc))

R's full response is, "You can drop 577 by *However, you will still fail.*" Furthermore, given the reasoning in section 3.1.1 above, R's full response would also inform Q if there is an action *b* that the user can take instead.

3.3. A better way

Suppose R believes that there is a better way to achieve Q's I-goal, cf. [8] - for example, taking an incomplete to have additional time to perform the work, and thereby not losing all the effort Q has already expended. Q would expect that R, as a cooperative expert, would inform him of such a better way, cf. [9]. If R doesn't, R risks misleading Q that there isn't one.

[8] RB((∃b)[holds(¬fail(Q,C), b(Sc)) &
admissible(b(Sc)) & better(b,drop(Q,C)(Sc))])

[9] RBQB(RB[want(Q,¬fail(Q,C)) &
RB[(∃b)[holds(¬fail(Q,C), b(Sc)) & admissible(b(Sc)) &
better(b,drop(Q,C)(Sc))]
⇒ likely(inform-that(R,Q,

$$(\exists b)[\text{holds}(\neg \text{fail}(Q,C), b(\text{Sc})) \ \& \ \text{admissible}(b(\text{Sc})) \ \& \ \text{better}(b, \text{drop}(Q,C)(\text{Sc}))], \text{Sc}]]$$

R's direct response is to indicate how f can be done. R's full response includes, in addition, " b is a better way."

Notice that if R doesn't explicitly tell Q that he is presenting a better way (i.e., he just presents the method), Q may be misled that the response addresses his S-goal: i.e., he may falsely conclude that he is being told how to drop the course. (The possibility shows up clearer in other examples - e.g., if R omits the first sentence of the response below

Q: How do I get to Marsh Creek on the Expressway?
 R: It's faster and shorter to take Route 30. Go out Lancaster Ave until....

Thus even when adhering to expert response behavior in terms of addressing an I-goal, we must keep the system aware of potentially misleading aspects of its modified response as well.

Note that R may believe that Q expects to be told the best way. This would change the second axiom to include within the scope of the existential quantifier

$$(\forall a)\{\neg(a=b) \Rightarrow [\text{holds}(\neg \text{fail}(Q,C), a(\text{Sc})) \ \& \ \text{admissible}(a(\text{Sc})) \ \& \ \text{better}(b,a)]\}$$

3.4. The only way

Suppose there is nothing inconsistent about what the user has proposed - i.e., all preconditions are met and it will achieve the user's goal. R's direct response would simply be to tell Q how. However, if R notices that that is the only way to achieve the goal (cf. [10]), it could optionally notify Q of that, cf. [11].

$$[10] \text{RB}((\exists! a)[\text{holds}(\neg \text{fail}(Q,C), a(\text{Sc})) \ \& \ \text{admissible}(a(\text{Sc})) \ \& \ a = \text{drop}(Q,C)(\text{Sc})])$$

$$[11] \text{RBQB}(\text{RB}(\text{want}(Q, \neg \text{fail}(Q,C))) \ \& \ \text{RB}((\exists! a)[\text{holds}(\neg \text{fail}(Q,C), a(\text{Sc})) \ \& \ \text{admissible}(a(\text{Sc})) \ \& \ a = \text{drop}(Q,C)(\text{Sc})]) \Rightarrow \text{likely}(\text{inform-that}(R, Q, (\exists! a)[\text{holds}(\neg \text{fail}(Q,C), a(\text{Sc})) \ \& \ \text{admissible}(a(\text{Sc})) \ \& \ a = \text{drop}(Q,C)(\text{Sc})], \text{Sc})))$$

R's full response is "You can drop 577 by That is the only way to prevent failing."

3.5. Something Turning Up

Suppose there is no appropriate action that Q can take to achieve his I-goal. That is,

$$\text{RB}(\neg(\exists a)[\text{admissible}(a(\text{Sc})) \ \& \ \text{holds}(g, a(\text{Sc}))])$$

There may still be some event e out of Q's control that could bring about the intended goal. This gives several more cases of R's modifying his response.

3.5.1. Unlikely event

If e is unlikely to occur (cf. [12]), Q would expect R to inform him of e , while noting its implausibility, cf. [13]

$$[12] \text{RB}((\exists e)[\text{admissible}(e(\text{Sc})) \ \& \ \text{holds}(\neg \text{fail}(Q,C), e(\text{Sc})) \ \& \ \neg \text{likely}(e, \text{Sc})])$$

[13] RBQB(RB(want(Q, ¬fail(Q,C)) &
 RB(¬(∃a)[admissible(a(Sc)) & holds(¬fail(Q,C),a(Sc))] &
 (∃e)[admissible(e(Sc)) & holds(¬fail(Q,C),e(Sc))
 & ¬likely(e,Sc)])
 ⇒ likely(inform-that(R, Q,
 (∃ e)[admissible(e,Sc) & holds(¬fail(Q,C), e(Sc))
 & ¬likely(e, Sc)]), Sc))

Thus R's full response is, "You can't drop 577. *If e occurs, you will not fail 577, but e is unlikely.*"

3.5.2. Likely event

If the event *e* is likely (cf. [14]), it does not seem necessary to state it, but it is certainly safe to do so. A formula representing this case follows.

[14] RB((∃e)[admissible(e(Sc)) &
 holds(¬fail(Q,C),e(Sc))] & likely(e,Sc))

R's beliefs about Q's expectations are the same as the previous case except that likely(*e*, Sc) replaces ¬likely(*e*, Sc). Thus R's full response may be "You can't drop 577. *However, e is likely to occur, in which case you will not fail 577.*"

3.5.3. Event followed by action

If event *e* brings about a state in which the enabling conditions of an effective action *a* are true, cf. [15]

[15] RB((∃e)(∃a)[admissible(e(Sc)) & admissible(a(e(Sc))) &
 holds(¬fail(Q,C), a(e(Sc)))]))

[16] RBQB(RB((∃e)(∃a)[want(Q, ¬fail(Q,C)) & admissible(e(Sc))
 & admissible(a(e(Sc))) & holds(¬fail(Q,C),a(e(Sc)))]))
 ⇒ likely(inform-that(R,Q,
 (∃e)(∃a) [holds(¬fail(Q,C),a(e(Sc)))] &
 admissible(a(e(Sc)))]),Sc))

then the same principles about informing Q of the likelihood or unlikelihood of *e* apply as they did before. In addition, R must inform Q of *a*, cf. [16]. Thus R's full response would be "You can't drop 577. *If e were to occur, which is (un)likely, you could a and thus not fail 577.*"

4. Reasoning

Our intent in using logic has been to have a precise representation language whose syntax informs R's reasoning about Q's beliefs. Having computed a full response that conforms to all these expectations, R may go on to 'trim' it according to principles of brevity that we do not discuss here.

Our proposal is that the informing behavior is "pre-compiled". That is, R does not reason explicitly about Q's expectations, but rather has compiled the conditions into a case analysis similar to a discrimination net. For instance, we can represent informally several of the cases in section 3.

```

If admissible(drop(Q,C)(Sc))
then if ¬holds(¬fail(Q,C),drop(Q,C)(Sc))
then begin nonproductive act
  if (∃b)[admissible(b(Sc)) & holds(¬fail(Q,C),b(Sc))]
  then a way
  else no way
end
else if (∃b)[admissible(b(Sc)) &

```

```

                holds( $\neg$ fail(Q,C),b(Sc)) & better(b,f)]
            then a better way
    else if ( $\exists$  b){admissible(b(Sc)) & holds( $\neg$ fail(Q,C), b(Sc))}
            then a way
            else no way

```

...

Note that we are assuming that R assumes the most demanding expectations by Q. Therefore, R can reason solely within its own space without missing things.

5. Conclusion

Since the behavior of expert systems will be interpreted in terms of the behavior users expect of cooperative human experts, we (as system designers) must understand such behavior patterns so as to implement them in our systems. If such systems are to be truly cooperative, it is not sufficient for them to be simply truthful. Additionally, they must be able to predict limited classes of false inferences that users might draw from dialogue with them and also to respond in a way to prevent those false inferences. The current enterprise is a small but non-trivial step in this direction. In addition to questions about achieving goals, we are investigating other cases where a cooperative expert should prevent false inferences by another agent, including preventing inappropriate default reasoning [6, JWW84nonmon].

Future work should include

- identification of additional cases where an expert must prevent false inferences by another agent,
- formal statement of a general principle for containing the search for possible false inferences, and
- design of a natural language planning component to carry out the informing acts assumed in this paper.

ACKNOWLEDGEMENTS

We would like to thank Martha Pollack, Deborah Dahl, Julia Hirschberg, Kathy McCoy and the AAAI program committee reviewers for their comments on this paper.

References

1. Allen, J. Recognizing Intentions from Natural Language Utterances. In *Computational Models of Discourse*, M. Brady, Ed., MIT Press, Cambridge MA, 1982.
2. Birnbaum, L., Flowers, M. & McQuire, R. Towards an AI Model of Argumentation. Proceedings of 1980 Conference, American Assoc. for Artificial Intelligence, Stanford CA, August, 1980.
3. Cohen, R. A Theory of Discourse Coherence for Argument Understanding. Proceedings of the 1984 Conference, Canadian Society for Computational Studies of Intelligence, University of Western Ontario, London Ontario, May, 1984, pp. 6-10.
4. Hirschberg, J. Scalar Implicature and Indirect Responses in Question-Answering. Proc. CSCSI-84, London, Ontario, May, 1984.
5. Joshi, A.K. Mutual Beliefs in Question Answering Systems. In *Mutual Belief*, N. Smith, Ed., Academic Press, New York, 1982.
6. Joshi, A., Webber, B. & Weischedel, R. Preventing False Inferences. Proceedings of COLING-84, Stanford CA, July, 1984.
7. Kowalski, Robert. *Logic for Problem Solving*. North Holland, New York, 1979.
8. Lehnert, W. A Computational Theory of Human Question Answering. In *Elements of Discourse Understanding*, A. Joshi, B. Webber & I. Sag, Ed., Cambridge University Press, 1981.
9. McCarthy, John. "Circumscription -- A Form of Non-Monotonic Reasoning". *Artificial Intelligence* 18 (1980), 27-39.
10. Pollack, Martha E. Goal Inference in Expert System. MS-CIS-84-07, University of Pennsylvania, 1984. Doctoral dissertation proposal.
11. Pollack, M. Good Answers to Bad Questions. Proc. Canadian Society for Computational Studies of Intelligence (CSCSI), Univ. of Western Ontario, Waterloo, Canada, May, 1984.
12. Ramshaw, Lance and Ralph M. Weischedel. Problem Localization Strategies for Pragmatics Processing in Natural Language Front Ends. Proceedings of COLING-84, July, 1984.
13. Reiter, R. Closed World Databases. In *Logic and Databases*, H. Gallaire & J. Minker, Ed., Plenum Press, 1978, pp. 149-177.
14. Sacerdoti, Earl D.. *A Structure for Plans and Behavior*. American Elsevier, New York, 1977.
15. Warren, D.H.D. WARPLAN: A System for Generating Plans. Proceedings of IJCAI-75, August, 1975.