

Global Anticipation Feedback as a User Modeling Technique for Dialog Systems*

Alassane Ndiaye

Anthony Jameson

Department of Computer Science, University of Saarbrücken
P.O. Box 151150, D-66041 Saarbrücken, Federal Republic of Germany
{ndiaye, jameson}@cs.uni-sb.de

Abstract

Human dialog participants regularly predict the responses of their dialog partners by hypothetically assuming the partner's role. This strategy of using *global anticipation feedback* has seldom been approximated in dialog systems. In the system PRACMA, the technical prerequisites for this strategy have been fulfilled, and the system is used as a testbed to explore the potential and limitations of the strategy. This paper first introduces a theoretical framework for analyzing possible realizations of global anticipation feedback. It then shows how the strategy can be realized in a dialog system that is capable of taking both roles within its dialog situation. An extension of these techniques is discussed that addresses the limited predictability of users' responses. The final sections discuss several approaches to minimizing the computational cost of using global anticipation feedback and address further uses of the anticipation strategy.

Résumé

Dans la communication humaine, les personnes impliquées dans une situation de dialogue prédisent les éventuelles réponses de leur interlocuteur. A cette fin, elles se mettent mentalement à la place de celui-ci, prennent son rôle dans la situation de dialogue pour anticiper ses réactions ou déterminer ses intentions. Cette stratégie dite d'anticipation de la rétroaction (*global anticipation feedback*) a été rarement utilisée dans les systèmes de dialogue homme-machine. Les outils et techniques préalables à l'implémentation de cette stratégie ont été développés dans le système de dialogue PRACMA qui sert aussi de plate-forme pour explorer les avantages et limites d'une telle stratégie. Dans cet article, nous présentons d'abord un cadre théorique pour analyser les possibilités de réalisation. Nous décrivons ensuite l'implémentation dans un système de dialogue capable de prendre les deux rôles dans une situation de dialogue. Enfin, nous décrivons une méthode pour améliorer les résultats de l'anticipation par la prise en compte de l'incertitude sur les mobiles et intentions de l'interlocuteur. Dans la dernière partie, nous discutons quelques approches permettant de minimiser le coût de l'utilisation de l'anticipation, et nous présentons d'autres possibilités d'utilisation.

Keywords: User Modeling, Anticipation Feedback, Dialog Systems, Transmutability

1 Local and Global Anticipation Feedback

As in everyday dialogs, intelligent dialog behavior of a natural language (NL) system involves the ability of the system to shift its cognitive perspective (cf. Flavell *et al.*, 1968; Higgins, 1981)

*This research was supported by the German Science Foundation (DFG) in its Special Collaborative Research Program on Artificial Intelligence and Knowledge Based Systems SFB 314, Project N1, PRACMA. This paper expands the shorter account in (Ndiaye & Jameson, 1996).

in order to take the role of the dialog partner and to simulate his or her dialog behavior.¹ One particular way in which a system can anticipate a user's responses is to make use of the system's own comprehension (and perhaps generation) capabilities, temporarily taking the role of the user and simulating his or her behavior. The term *anticipation feedback loop* (AFL) was introduced to characterize such cases (Jameson & Wahlster, 1982; Wahlster & Kobsa, 1989).

The types of user response that have been anticipated in AFLs include the following, among others: the correct interpretation of elliptical utterances (Jameson & Wahlster, 1982); the accurate visualization of scene descriptions (Novak, 1987; Schirra, 1995; Blocher & Schirra, 1995); the drawing of correct inferences (see, e.g., Joshi, Webber, & Weischedel, 1984; Zukerman, 1990); and the pragmatic interpretation of utterances (Jameson, 1989).

To date, almost all implemented systems that have employed anticipation feedback have used a limited part of the system to realize a *local* AFL. ANTLIMA (Schirra, 1995; Blocher & Schirra, 1995) is a rare example of a system that uses a *global* AFL (Wahlster & Kobsa, 1989, pp. 22–26): A large part of the system's own understanding capabilities is used to anticipate the user's responses—in ANTLIMA, the way the user will visualize verbal descriptions of events perceived by the system.

One way of viewing the role of an AFL in an interactive system is illustrated by the decision tree in the left-hand side of Figure 1. A system \mathcal{S} has to choose among several possible dialog moves $m_1 \dots m_n$ that will have some effect on the user \mathcal{U} . Each m_i has some immediate degree of appeal for \mathcal{S} , which can be conceived of as a utility $U_m(m_i)$. But instead of selecting the move with the highest $U_m(m_i)$, \mathcal{S} anticipates the response r_i that \mathcal{U} is likely to make to each m_i ; and each r_i is itself associated with a utility $U_r(r_i)$. \mathcal{S} chooses the move with the highest total utility $U_m(m_i) + U_r(r_i)$. An AFL can be invoked in the step where \mathcal{S} anticipates \mathcal{U} 's response r_i . The point of doing so is that the determination of $U_r(r_i)$ in addition to $U_m(m_i)$ may affect \mathcal{S} 's choice of a move.

Most systems that have used AFLs have not explicitly reasoned in terms of decision trees and utilities. Nonetheless, their approaches can mostly be viewed as variants on the scheme of the left-hand side of Figure 1. For example, \mathcal{S} may not explicitly compute utilities but may rather simply reject m_i if r_i is clearly undesirable (e.g., if it involves a misunderstanding by \mathcal{U}). And the candidate moves m_i can be generated one by one, instead of all at once; in this way, each m_i can represent an improvement on m_{i-1} , taking into account the results of the anticipations $r_1 \dots r_{i-1}$. Further variants will be mentioned below.

The present paper explores this relatively uncharted area of global anticipation feedback. The potential benefits of global AFLs, as well as the problems involved in realizing them, will be introduced in the remainder of this section with some examples from the dialog system PRACMA (Jameson *et al.*, 1994; Jameson *et al.*, 1995), which we use here as a testbed. The goals of this research are (a) to create a theoretical framework for analyzing global AFLs which is more precise and differentiated than the corresponding theories that can be found in the relevant psychological and user modeling literature; and (b) to lay a foundation for practically useful applications of global AFLs in dialog systems.

¹“Every communicator carries around with him an image of the receiver. He takes his receiver (as he pictures him to be) into account when he produces a message. He anticipates the possible responses of this receiver and tries to predict them ahead of time. These images affect his own message behaviors.” (Berlo, 1960, p. 117)

“Penser à la pensée d'autrui est une caractéristique essentielle de toute attitude sociale; chacun cherche à suivre et à devancer le progrès de la pensée de l'autre, l'avantage étant à celui qui devine une pensée de l'autre que celui-ci croit ignorée.” (Guillaume, 1954, p. 182)

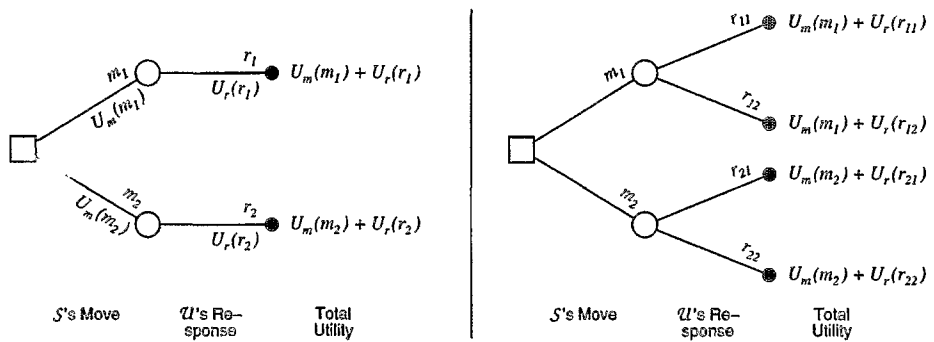


Figure 1: Decision trees that illustrate the strategy of anticipating the user's response before selecting a move.

The right-hand side represents a generalization of the left-hand side by taking into account uncertainty about the user's responses.

1.1 Dialog Situation of PRACMA

PRACMA models dialogs in which a person (to be called the *seller*) is trying to sell his or her used car to a potential *buyer* (cf. Figure 2). PRACMA is able to take the role of either the seller or the buyer within its dialog situation. This ability to switch roles can be seen as a particular variant of the property of dialog systems that Wahlster & Kobsa (1989, p. 30) define as *transmutability*: the property of being adaptable to applications that differ with respect to "dialog type, user type, and intended system behavior". This within-dialog transmutability in PRACMA enables the system to realize various types of AFLs—in particular truly global AFLs, in which the system consults a complete instantiation of itself in order to anticipate the dialog partner's responses.

The goals of the two dialog participants in PRACMA's example domain conflict to a certain degree: The buyer wants to get the best possible information on which to base a decision about the car, whereas the seller would like to sell the car, whether or not it is really suitable for the buyer. When S is the seller, this conflict increases the importance of anticipation feedback for S , because it increases the range of utility that U 's responses can have for S . For example, if U decides to ask about an attribute of the car, it makes a big difference to a noncooperative seller whether this attribute happens to be one on which the car rates highly or poorly; for a cooperative seller, this difference would not be so important.

1.2 A Simple and a Complex Local AFL

When PRACMA takes the role of the seller, one frequent task of the system is to decide what (if anything) to say about some attribute of the car (e.g., about the car's overall mileage). To do so, S anticipates the effect of various possible comments concerning that attribute on the buyer's evaluation process. Two types of local AFLs have been implemented for this purpose, a simple, fast one and a sophisticated, more time-consuming one.

When using the *simple* AFL, S simply invokes the same procedures that it would use in the role of the buyer to determine the evaluative implications of each comment it is considering, essentially

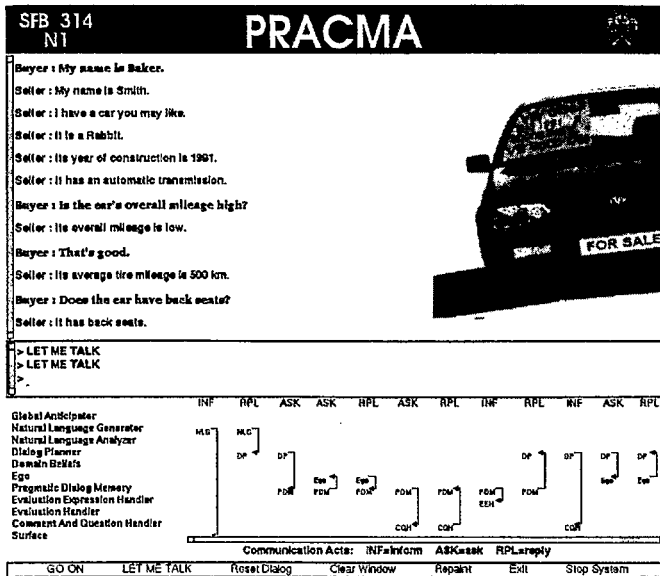


Figure 2: Beginning of an example dialog with PRACMA; the system can take the role of either the seller or the buyer.

asking “What effect would this comment have on *my* evaluation of a car?”.

When using the *complex* AFL, *S* takes into account a number of possible differences between its own knowledge and evaluation standards and those of the *U* the system is talking to. For example, *S* might estimate (a) how likely it seems to *U* that the car's overall mileage is high, (b) whether *U* knows that overall mileage has implications not only for safety but also for the likelihood of passing the next inspection, and (c) how much importance *U* assigns to the evaluation dimensions of safety and reliability (Jameson & Schäfer, 1994).

PRACMA uses the simple AFL when taking the role of a seller who is not able or willing to devote much attention to the dialog. Both AFLs are local in that the only part of the system's comprehension capabilities that *S* makes use of is *S*'s capability to derive the evaluative implications of a given comment about a car.

1.3 A Global AFL

The use of a global AFL in PRACMA becomes necessary when *S* tries to anticipate a more complex response by *U*: What *U*'s next dialog move will be if *S* makes a given comment. For example, if *S* states that the car's overall mileage is low, *U* might (a) say nothing and keep listening, (b) express some reaction like “That's good”, and/or (c) ask a further question, probably related in some way to the topic of overall mileage. It can be important for *S* to anticipate *U*'s next move. For example, if *U* asks a question, it might concern some topic that *S* would prefer to avoid (e.g., an attribute with respect to which the car is weak); if this seems likely, *S* should consider not making the comment it originally intended to make.

\mathcal{S} anticipates \mathcal{U} 's next dialog move by temporarily pretending that it is the buyer in the dialog. More concretely, \mathcal{S} consults a simultaneously active instantiation of PRACMA in which the system is taking the role of the buyer. \mathcal{S} basically asks "How would I respond if I were the buyer and if the seller told me (in this dialog context) that the car's overall mileage is low?". \mathcal{S} uses the results of this simulation as feedback to decide whether to make the comment about the car's overall mileage in the first place.

1.4 Potential and Limitations of Global AFLs

When Global AFLs are Especially Applicable

The reason why a global AFL is used in the second situation is that the response which has to be anticipated is determined by a number of different capabilities of \mathcal{U} : \mathcal{U} does not only have to interpret \mathcal{S} 's comment and determine its implications for the evaluation of the car (as in the first situation). Rather, \mathcal{U} also has to consider what kind of dialog move to produce next; this depends on the dialog strategy \mathcal{U} is pursuing. And if \mathcal{U} decides to ask a question about some attribute of the car, \mathcal{U} has to decide which of the many attributes to ask about and how to formulate the question.

It would in principle be possible for \mathcal{S} to use a combination of several local AFLs in order to anticipate \mathcal{U} 's next dialog move; but this approach would require the system designers to provide for some quite complex processing, which would be applicable only to this particular type of anticipation. For anticipating \mathcal{U} 's dialog moves in other dialog contexts, different solutions would have to be found.

By contrast, once the system has been given the capability to obtain global anticipation feedback, this single general technique can be used to anticipate many different types of (observable and unobservable) responses by the dialog partner.

Issues Raised by Global AFLs

This simplicity and generality is, however, associated with a number of limitations and challenges.

1. Within-dialog transmutability. A system that uses a global AFL must be able to take the role of the other participant in the type of dialog it conducts. By contrast, a local AFL presupposes only that the system be able to do some part of the processing required for the other role; and this common processing may involve a generic subtask, such as syntactic analysis, which is relatively independent of any particular dialog role. For human beings, transmutability is often given, because people learn to take many different roles in dialogs in the course of their everyday experience. (For example, even a professional salesperson often has the opportunity to act as a customer.) But systems that employ user modeling techniques are typically designed to play a particular role. It may therefore require a considerable additional investment to enable them to switch to the role of their dialog partner.²

2. Communication between system instantiations. It is not trivial for a dialog system to invoke itself in another role without interfering with its workings in its original role.

3. Uncertainty about factors that determine the responses of the dialog partner. The decision tree in the left-hand side of Figure 1 presupposes that \mathcal{S} can predict \mathcal{U} 's response to any given move m_i with certainty. The more general case is shown in the right-hand side of the figure: When considering a move m_i , \mathcal{S} can at best narrow the possible responses of \mathcal{U} down to a set $\{r_{ij}\}$. In the case of a global AFL, this uncertainty is due to the fact that \mathcal{U} 's response will be influenced by some factors that are not entirely known to \mathcal{S} . In other words, \mathcal{S} does not know

²An alternative way of realizing a global AFL will not be considered further here: The approach of linking a system to another completely different system that is capable of taking the other role in the dialog situation.

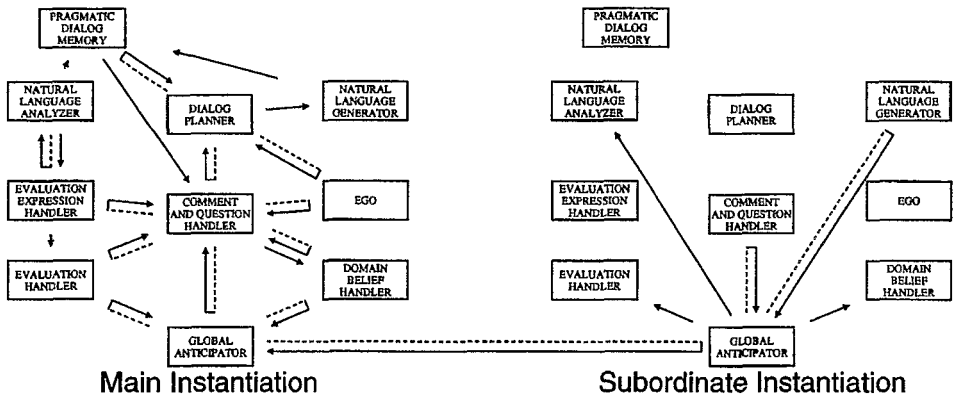


Figure 3: Overview of PRACMA's agents and their communication.

A single solid arrow represents an INFORM message, which passes unsolicited information from one agent to another. A dashed line followed by a solid arrow represents a sequence of two communication acts of the types ASK and REPLY, respectively. The subordinate instantiation of the system is used for global anticipation feedback. Its agents exchange the same types of message as the agents in the main instantiation, and they also exchange the types of message shown in the right-hand side of the figure.

exactly how to pretend to be the user. The question arises as to how \mathcal{S} can best deal with this uncertainty.

4. Efficiency. In general, it is relatively time-consuming for a dialog system to anticipate the complete processing by the user of a possible dialog move by the system. The question therefore arises of how the computational cost of using global AFLs can be minimized.

2 An Architecture for Anticipation Feedback

As background to the discussion of these issues, this section first introduces the overall architecture of PRACMA and shows how it supports both local and global AFLs.

2.1 Multi-Agent Architecture

PRACMA is implemented within the multi-agent architecture CHANNELS (Ndiaye & Jameson, 1994). CHANNELS uses techniques from distributed AI and from concurrent object-oriented programming. The system's modules, realized as agents, interact cooperatively through a communication-act-based protocol that governs the exchange of messages. The agents run concurrently as simulated processes over a local network using PVM (Geist *et al.*, 1994) and ICE (Amtrup, 1994). The left-hand side of Figure 3 shows PRACMA's agents and the types of messages that they exchange.

2.2 Realization of Within-Dialog Transmutability

The agents fall into three categories that illustrate different ways of realizing within-dialog transmutability:

1. The COMMENT AND QUESTION HANDLER (CQH for short) and the DIALOG PLANNER are responsible for high-level interpretation and generation of the system's utterances. Their basic workings are similar for both roles, but they use a good deal of role-specific declarative and procedural knowledge.

2. The EVALUATION HANDLER and the DOMAIN BELIEF HANDLER reason in both roles about the evaluations and beliefs, respectively, of the buyer. Each one uses the same basic formalism in both roles: Bayesian networks and a modal-logic-based knowledge representation system (Hustadt & Nonnengart, 1993), respectively. In the role of the buyer, this reasoning is quite simple; but when the system is the seller, it constitutes more complex meta-level reasoning.
3. The remaining agents function in essentially the same way in both roles.

2.3 Realization of Local Anticipation Feedback

The local AFL sketched in Section 1.2, invoked when S is taking the role of the seller, involves straightforward communication between CQH and the EVALUATION HANDLER of the instantiation of PRACMA in the role of the seller. Recall that the purpose is to anticipate U 's evaluative response (in the role of the buyer) to a comment C . In the *simple* variant, CQH sends a query to the EVALUATION HANDLER that was originally designed for the case where the system is taking the role of the buyer: The EVALUATION HANDLER is asked how much the system's *own* evaluation of a car would change on the basis of the comment C . In the *complex* variant, CQH sends a meta-level query that is only applicable when the system is the seller: It asks the EVALUATION HANDLER to reason on the meta-level about U 's evaluation process and predict U 's evaluative response. In either case, after obtaining this information from the EVALUATION HANDLER, CQH also considers other properties of C , such as the extent to which it is true and its degree of relatedness to the current dialog focus. CQH then returns to the DIALOG PLANNER a comment that rates well overall according to this set of criteria, or it reports to the DIALOG PLANNER that no such comment can be found.

3 Realization of Global Anticipation Feedback

The job of choosing the system's next dialog move is divided hierarchically between the DIALOG PLANNER and CQH. The DIALOG PLANNER, an incremental planner, decides what *type* of move to make. In doing so, it takes into account a variety of factors, including the dialog history (stored in the PRAGMATIC DIALOG MEMORY) and various motivational parameters (stored in EGO). Once it has decided on a particular type of move, it asks CQH to choose a specific move of that type. CQH does this by executing the algorithm BEST-MOVE (see Figures 4 and 5).³ This algorithm realizes the basic strategy sketched in the left-hand side of Figure 1. Its use is illustrated by the example in the left-hand side of Figure 6. The example presupposes that the system is the seller and that the dialog context is the one shown at the end of Figure 2—that is, S has just mentioned that the car has back seats and S now has a chance to say something else. The DIALOG PLANNER has decided that a move of the type "comment-on-attribute" should be made. There are a number of comments with some relevance to the topic of back seats that S could conceivably make, but only the ones shown on the left in Figure 6 have a sufficiently large UTILITY-OF-MOVE to be worth considering further. The procedure for assessing UTILITY-OF-MOVE is different for each type of move, as is shown in Figure 5. For the type "comment-on-attribute", the corresponding algorithm, PREDICTED-EVALUATION-SHIFT, is also shown in Figure 5.

If S did not use global anticipation feedback, in this example S would simply select the comment "The car has four doors", because of its high UTILITY-OF-MOVE.

To take into account UTILITY-OF-ANTICIPATED-RESPONSE as well, CQH uses global anticipation to predict how U would respond to each of these three comments. A global AFL cannot be realized as a query to one of the agents that make up the system, since it requires an invocation of the entire system. Therefore, the agent GLOBAL ANTICIPATOR maintains a *subordinate instantiation*

³The particular pseudocode notation used in this and later figures is that of (Russell & Norvig, 1995).

```

function BEST-MOVE(type, constraints) returns a dialog move of type type that fulfills constraints
/* Executed by CQH */
possible-moves ← ALLOWABLE-MOVES(type, constraints)
for m in possible-moves do
    UTILITY[m] ← UTILITY-OF-MOVE(m)
end
reasonable-moves ← subset of possible-moves with UTILITY >  $\delta$ 
if *global-anticipation?* = True then
    for m in reasonable-moves do
        UTILITY[m] ← UTILITY[m] + UTILITY-OF-ANTICIPATED-RESPONSE(m)
    end
return the m in reasonable-moves with the highest UTILITY[m]

```

Figure 4: CQH's algorithm for choosing a dialog move of a given type.

The square brackets in the term UTILITY[*m*] reflect the fact that UTILITY is not a function but rather an attribute of a dialog move. The *reasonable-moves* correspond to the $m_1 \dots m_n$ in the left-hand side Figure 1; UTILITY-OF-MOVE and UTILITY-OF-ANTICIPATED-RESPONSE correspond to U_m and U_r , respectively..

```

function UTILITY-OF-MOVE(move) returns a utility
/* Executed by CQH */
case TYPE[move]
    Moves possible in both roles:
    silence
        return 0
    ...
    Moves possible in role of seller:
    comment-on-attribute:
        return ASK(EVALUATION HANDLER, "PREDICTED-EVALUATION-SHIFT(move)")
    ...
    Moves possible in role of buyer:
    question-about-attribute:
        return ASK(EVALUATION HANDLER, "CURRENT-EVALUATIVE-UNCERTAINTY(TOPIC[move])")
    evaluative-reaction:
        return INFORMATIVENESS(move)
    ...

```

```

function PREDICTED-EVALUATION-SHIFT(comment) returns an estimate of the shift in  $\mathcal{U}$ 's evaluation that comment
would lead to
/* Executed by the EVALUATION HANDLER */
construct and evaluate a Bayesian network to predict  $\mathcal{U}$ 's evaluation shift, taking into account  $\mathcal{S}$ 's uncertain beliefs
about  $\mathcal{U}$ 's interests and knowledge
prediction ← the probability distribution representing the resulting belief concerning  $\mathcal{U}$ 's evaluation shift
return EXPECTED-VALUE(prediction)

```

Figure 5: Algorithms used by CQH to assess the immediate utility of a dialog move.

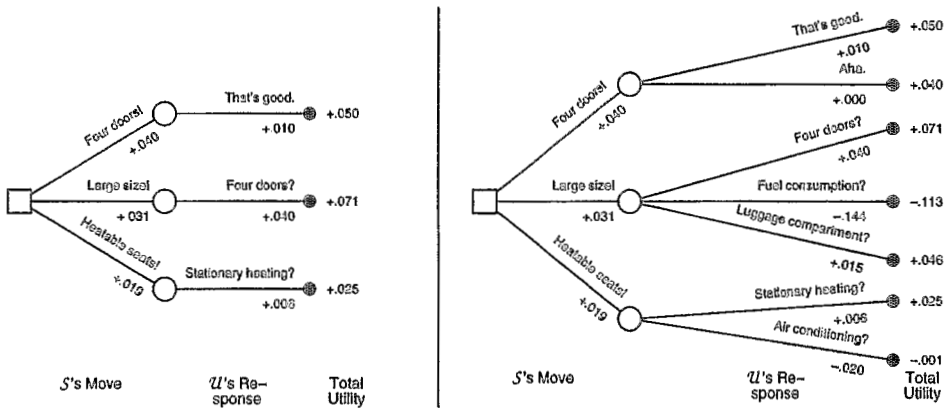


Figure 6: Example of use of global anticipation.

The left-hand side shows an example without taking uncertainty into account. The right-hand side shows the same example in which uncertainty is taken into account.

of PRACMA in a separate COMMONLISP image, which may be located on another computer.⁴ The GLOBAL ANTICIPATOR (nicknamed GAF for Global Anticipation Feedback) initializes the subordinate instantiation to take the other dialog role. The subordinate instantiation does not engage in any direct interaction with the user but rather responds to inputs from the GAF that controls it.

This subordinate image contains instantiations of all of PRACMA's agents, including GAF (see the right-hand side of Figure 3). One task of GAF is to ensure that the subordinate instantiation constitutes as realistic a model of the actual user as is possible given the information available to the main instantiation. To accomplish this, GAF regularly queries the agents EVALUATION HANDLER and DOMAIN BELIEF HANDLER in the main instantiation, asking them for their assessments of U's evaluation criteria and knowledge; the resulting estimates are used to initialize the corresponding agents in the subordinate instantiation.

In addition to maintaining the subordinate instantiation, GAF handles communication between the two instantiations, as shown in Figure 3.

The first two algorithms in Figure 7 summarize how the subordinate instantiation is used to help determine the UTILITY-OF-ANTICIPATED-RESPONSE corresponding to a possible move by S. Even after U's response has been anticipated, the task still remains of assessing how desirable that move is from the point of view of S. To do this really thoroughly, S would have to take into account the possible moves that S might subsequently make; that is, S would have to expand the decision tree in the left-hand side of Figure 6, adding nodes to the right. The feasibility of this approach will be discussed below; but in any case at some point S must stop looking ahead and must form at least a crude assessment of how desirable a particular outcome is (cf. the evaluation functions

⁴The CHANNELS architecture provides external agents that make possible communication between different instantiations. An external agent is a virtual agent within a PRACMA instantiation with which the agents within the system can communicate but which is actually located in another instantiation of the system.

```

function UTILITY-OF-ANTICIPATED-RESPONSE(move) returns a utility
/* Executed by CQH */
return UTILITY-OF-RESPONSE(ASK(GAF, "ANTICIPATED-RESPONSE(move)"))



---


function ANTICIPATED-RESPONSE(move) returns an anticipated response to move
/* Executed by GAF */
INFORM(NATURAL LANGUAGE ANALYZER, "UTTERANCE-INTERPRETED(move)")
return ASK(NATURAL LANGUAGE GENERATOR, "INTERNAL-REPRESENTATION(LATEST-UTTERANCE())")



---


function UTILITY-OF-RESPONSE(response) returns a utility
/* Executed by CQH */
case TYPE[response]
  Responses possible in both roles:
  silence:
    return 0
  ...
  Possible responses from buyer:
  question-about-attribute:
    return UTILITY-OF-MOVE(BEST-MOVE("comment-on-attribute", "(topic = TOPIC[response]"))))
  evaluative-reaction:
    return POSITIVENESS[response]
  ...
  Possible responses from seller:
  comment-on-attribute:
    return ASK(EVALUATION HANDLER, "RESULTING-UNCERTAINTY-REDUCTION(response)")
  ...

```

Figure 7: Algorithms for global anticipation feedback.

ANTICIPATED-RESPONSE causes the subordinate instantiation to process a possible move by \mathcal{S} ; it obtains from the NATURAL LANGUAGE GENERATOR an internal representation of the resulting response by the simulated \mathcal{U} . (The language-processing capabilities of the NATURAL LANGUAGE ANALYZER and the NATURAL LANGUAGE GENERATOR are not made use of within the subordinate instantiation.) UTILITY-OF-RESPONSE heuristically assesses the desirability of this response from the point of view of \mathcal{S} .

used to assess the desirability of positions in board games such as chess). Some of the heuristics that CQH uses for this purpose are shown in the definition of UTILITY-OF-RESPONSE in Figure 7. For example, the assessed utility to \mathcal{S} of having \mathcal{U} ask whether the car has four doors is equal to the utility of \mathcal{S} 's answering the question (e.g., by saying "It has four doors"). If \mathcal{U} does ever ask this particular question, the system may in fact find a better way to respond to it; but this rough assessment still has some value for planning purposes.

Similarly, explicitly positive or negative expressions of evaluation by \mathcal{U} are judged to be in themselves desirable or undesirable, respectively.⁵

The left-hand side of Figure 6 illustrates how, when the ANTICIPATED-UTILITY-OF-RESPONSE is taken into account, the relative utilities of the various possible moves by \mathcal{S} can be different than they would have been without global anticipation feedback. Specifically, \mathcal{S} chooses to comment on the car's size, even though the statement that it has four doors would make a better initial

⁵The reasoning here is, for example, that buyers who hear themselves express positive evaluations will come to perceive themselves as liking the car and will therefore be more inclined to decide in favor of it.

impression, because S anticipates that after the former comment U will proceed to ask a question about the number of doors anyway.

4 Taking Uncertainty into Account

The algorithms discussed so far have presupposed that, if S takes the trouble to use a global AFL, S will always anticipate U 's response correctly. As already mentioned, this assumption is less realistic than the conceptualization shown in the right-hand side of Figure 1. But the question remains: How can global anticipation, given a possible move m_i , return not just a single anticipated response r_{i1} , but rather a *set* of possible responses $\{r_{i1} \dots r_{in}\}$?

This question is difficult to answer in a general way. But within the framework presented here, the problem is manageable if S considers only other responses of the *same type* as the most likely response r_{i1} . The basic idea is to exploit the way in which CQH chooses moves of a given type, namely by evaluating all reasonable moves of that type (cf. Figure 4). Although this algorithm has been discussed so far only with respect to its use in the main instantiation, it is of course also used in the subordinate instantiation, when U is being simulated. For example, when the subordinate instantiation, in the role of the buyer, chooses a specific question to ask, the CQH of the subordinate instantiation first considers all questions that have some relevance to the current dialog focus and then chooses the one with the highest UTILITY-OF-MOVE. A consequence is that when the subordinate instantiation has produced a move r_{i1} for U as a response to the move m_i by S , GAF can ask CQH which moves it considered that had a UTILITY-OF-MOVE that was *almost as high* as that for r_{i1} . The assumption underlying this query is the following one: The moves that rated almost as high as r_{i1} for the simulated U represent the most likely alternative hypotheses about how U will respond to m_i .

The algorithms in Figure 8 realize this strategy, which will be called the *runner-up strategy*. They are generalizations of the corresponding algorithms in Figure 7. Note that anticipating a set of possible responses is in itself no more time-consuming than anticipating a single one.

The right-hand side of Figure 6 shows how the example in left-hand side turns out if S applies the runner-up strategy. Now S takes into account the possibility that U might respond to the comment on the car's size by asking about its fuel consumption. As this happens to be a major weak point of the car, the comment on size now appears to be the least desirable of the three possible comments.

The runner-up strategy, as realized here, suffers from a fundamental limitation: Even if ϵ is set high, the set returned by ANTICIPATED-RESPONSES will often not include U 's actual response—namely, in cases where that response is of a different type than the ones considered by the simulation of U . For example, S may take into account five possible questions by U but fail to take into account an evaluative reaction that is in fact more likely than most of the questions. There appears to be no straightforward way to anticipate the most likely responses of *all types*, given the hierarchical way in which the DIALOG PLANNER and CQH work together to choose moves.

This limitation does not appear to be specific to the particular implementation of global anticipation feedback described here. In general, it cannot be assumed that all of the responses that an agent *might have made* if its parameters had been slightly different are responses that the agent *considered making* during its selection of a single response.

Therefore, if the system requires a more thorough overview of possible responses by the user, it will have to invoke several different simulations of the user, initializing each one somewhat differently. This approach is currently being explored with PRACMA.

```

function UTILITY-OF-ANTICIPATED-RESPONSE (move) returns a utility
  /* Executed by CQH */
  possible-utilities ← set of utilities obtained by applying UTILITY-OF-RESPONSE to the results yielded by ASK(GAF,
  "ANTICIPATED-RESPONSE(move)")
  return AVERAGE(possible-utilities)

```

```

function ANTICIPATED-RESPONSES(move) returns a set of possible responses to move
  /* Executed by GAF */
  INFORM(NATURAL LANGUAGE ANALYZER, "UTTERANCE-INTERPRETED(move)")
  return ASK(CQH, "BEST-MOVES-CONSIDERED( $\epsilon$ )")

```

```

function BEST-MOVES-CONSIDERED( $\epsilon$ ) returns a set of moves
  /* Executed by CQH */
  reasonable-moves ← the set of reasonable moves considered during the most recent execution of BEST-MOVE
  best-move ← the m in reasonable-moves with the highest UTILITY[m]
  return the set of all m in reasonable-moves such that
    UTILITY[m] ≥ UTILITY[best-move] -  $\epsilon$ 

```

Figure 8: Generalizations of the global anticipation algorithms (Figure 7) that take into account uncertainty about \mathcal{U} 's response.

In ANTICIPATED-RESPONSES, instead of obtaining a single result from the NATURAL LANGUAGE GENERATOR, GAF obtains from CQH a by-product of its processing. If the relative likelihoods of the ANTICIPATED-RESPONSES(*move*) could be estimated, these estimates would enter into the calculation of UTILITY-OF-ANTICIPATED-RESPONSE; at present, the average of the *possible-utilities* is used as a rough approximation.

5 Efficiency Considerations

Because global anticipation feedback is computationally expensive, a system must be selective in applying it. This section discusses several possible types of selectivity.

Evaluating moves and responses selectively. The algorithms presented above are consistent with the idealized use of decision trees in that they presuppose that all branches are to be processed completely. But if it is acceptable to sacrifice some decision quality, computation can be done more selectively. For example, within a satisficing strategy the search for a move can be terminated as soon as one move with an acceptable overall utility has been found.

Minimizing look-ahead. As has already been noted, there is no reason in principle why \mathcal{S} cannot expand the decision trees shown in the above figures so as to look farther into the future. For example, game-playing programs often look at least several moves ahead.⁶ A different type of look-ahead can be achieved if global anticipation is allowed to occur within the subordinate instantiation: When anticipating \mathcal{U} 's next response, \mathcal{S} considers how \mathcal{U} will anticipate \mathcal{S} 's subsequent move, etc. An important limitation of both of these types of look-ahead is their relatively high computational cost. For example, to extend the decision tree in the right-hand side of Figure 1 beginning with one of the right-most nodes, \mathcal{S} has to go through the whole process of generating possible moves, a process which can involve all of the agents which make up the system. Note also that as the tree gets deeper, the additional expansions become less worthwhile,

⁶Note that the decision trees used here differ from the game trees used with techniques like minimax: The utility criteria of the user are not directly opposed to those of the system, and the moves of each participant are based in part on considerations that the other participant is not entirely aware of. Accordingly, the system views the user not as an adversary but simply as a cause of events that have a limited degree of predictability.

```

function REASONABLE-MOVE-POSSIBLE?(type, constraints) returns a truth value
/* Executed by CQH */
for move in ALLOWABLE-MOVES(type, constraints) do
  if UTILITY-OF-MOVE(move) >  $\delta$  then return True
end
return False

```

Figure 9: Algorithm used by CQH to answer relatively quickly a query by the DIALOG PLANNER as to whether an acceptable move of a given type is available in the current context.

The DIALOG PLANNER uses the results of such queries, along with other considerations, to decide what type of move actually to request. Note that this algorithm can yield *True* even in cases where BEST-MOVE (Figure 4) would not return a satisfactory move.

as they concern dialog moves which are increasingly unlikely ever to occur. One reasonable approach is to make the amount of look-ahead dependent on (a) the resources available to the system and (b) the assessed importance of correct anticipation.

Skipping global anticipation feedback during initial planning. It may be necessary to restrict the use of global anticipation feedback to a late stage in the utterance planning process. For example, when deciding what type of dialog move to make next, the DIALOG PLANNER often asks CQH whether it is possible to make a worthwhile move of a given type. Even if CQH responds positively, the DIALOG PLANNER may end up choosing a different type of dialog move, since other criteria are also relevant. Because this type of query by the DIALOG PLANNER comes frequently, it would be impractical for CQH to invoke GAF (perhaps repeatedly) every time it answers such a query. Instead, as is shown in Figure 9, CQH simply checks whether there is some move of the type in question that is acceptable with respect to the relatively simple criterion UTILITY-OF-MOVE. It is only when (and if) the DIALOG PLANNER subsequently asks CQH actually to select a move of this type that CQH takes the trouble to invoke GAF. When it does so, it may of course discover that all of the possible moves rate poorly with respect to UTILITY-OF-ANTICIPATED-RESPONSE. In such cases the system's behavior is similar to that of a person who begins to say something and then has second thoughts about the wisdom of doing so. The occasional appearance of this phenomenon seems to be a necessary consequence of the limited time that the system can spend anticipating the user's responses during the early planning of a dialog contribution.

Selective updating of the subordinate instantiation. One necessary aspect of a procedure for global anticipation feedback has only been mentioned briefly so far: the updating of the subordinate instantiation on the basis of estimates of the user's knowledge, evaluation criteria, and other characteristics. Performing this updating frequently can be not only time-consuming but also wasteful. For example, only a small part of the updates may actually have any effect on the anticipation of *U*'s next move. A simplified approach is to do the updating only occasionally—or even only once, at the beginning of a dialog, on the basis of the initially available information about the user.⁷

⁷The tendency that people sometimes show, especially at an early age, to ignore differences between themselves and their dialog partners (see, e.g., Astington, 1993; Flavell *et al.*, 1968; Higgins, 1981; Oléron *et al.*, 1981) may in some cases represent an application of this strategy of selective updating.

6 Further Uses of Global Anticipation Feedback

In addition to the uses of global anticipation feedback in PRACMA discussed in the previous sections, two further uses has been explored.

Global anticipation by PRACMA in the role of the buyer. Global AFLs are similarly applicable when PRACMA takes the role of the buyer in its dialog situation (Ndiaye, 1996a; Ndiaye, 1996b). For example, consider the buyer who is concerned that a car may have air conditioning, because air conditioning is somewhat harmful to the environment. If the buyer starts the dialog by asking whether the car has air conditioning, the seller is likely to infer that the buyer attaches high importance to comfort; the seller may therefore start volunteering information about other comfort-related attributes of the car. If the buyer can anticipate this response, he or she can postpone the question about air-conditioning; later, when the seller has had time to form a fairly accurate model of the buyer's evaluation standards, the buyer may ask the question, anticipating that this problem will no longer arise.

Anticipating internal responses rather than dialog moves. For concreteness, the discussion above has focused on the problem of anticipating what the user will do next in the dialog. But in many cases, what S needs to know about is some aspect of U 's internal processing. It is fairly straightforward to extend the methods proposed here to handle this sort of anticipation. In fact, Figure 8 already showed how GAF can return to the main instantiation information about internal states of the subordinate instantiation (here: concerning the responses that were considered by the CQH of the subordinate instantiation). But using a global AFL (as opposed to a local one) to anticipate internal responses will only be worthwhile if the responses are determined in a complex way; otherwise, a more local form of anticipation (such as the one sketched in the algorithm PREDICTED-EVALUATION-SHIFT in Figure 5) is likely to be feasible and preferable.

7 Conclusion

In sum, the presented framework has shown the potential benefits of the use of truly global anticipation feedback in a dialog system. The experience reported here has shown that global anticipation feedback is in fact a feasible technique with many potential uses in dialog systems; and that there exist enough degrees of freedom in realizing and applying the technique to enable designers to overcome some of the problems that may initially make the use of the technique appear to be impractical.

References

- Amtrup, J. 1994. ICE: INTARC Communication Environment — Design und Spezifikation. Technical Report VM-Memo 48, University of Hamburg, Hamburg.
- Astington, J. W. 1993. *The Child's Discovery of the Mind*. Cambridge, MA: Harvard University Press.
- Berlo, D. K. 1960. *The Process of Communication: An Introduction to Theory and Practice*. New York: Holt, Rinehart and Winston.
- Blocher, A., and Schirra, J. R. J. 1995. Optional deep case filling and focus control with mental images: ANTLIMA-KOREF. In Mellish, C. S., ed., *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*. San Mateo, CA: Morgan Kaufmann. 417–423.
- Flavell, J. H.; Botkin, P. T.; Fry Jr., C. L.; Wright, J. W.; and Jarvis, P. E. 1968. *The Development of Role-Taking and Communication Skills in Children*. New York: Wiley.
- Geist, A.; Beguelin, A.; Dongorra, J.; Jiang, W.; Manchek, R.; and Sunderman, V. 1994. *PVM: Parallel Virtual Machine. A User's Guide and Tutorial for Networked Parallel Computing*. Cambridge, MA: MIT Press.
- Guillaume, P. 1954. *Introduction à la psychologie*. Paris: J. Vrin.

- Higgins, E. T. 1981. Role taking and social judgment: Alternative developmental perspectives and processes. In Flavell, J. H., and Ross, L., eds., *Social Cognitive Development: Frontiers and Possible Futures*. Cambridge, England: Cambridge University Press. 119–153.
- Hustadt, U., and Nonnengart, A. 1993. Modalities in knowledge representation. In *Proceedings of the Sixth Australian Joint Conference on Artificial Intelligence*. Singapore: World Scientific Publishing. 249–254.
- Jameson, A., and Schäfer, R. 1994. Dynamically constructed Bayesian networks for modeling interests and knowledge. In *Proceedings of the Fourth International Conference on User Modeling*, 219.
- Jameson, A., and Wahlster, W. 1982. User modelling in anaphora generation: Ellipsis and definite description. In *Proceedings of the Fifth European Conference on Artificial Intelligence*, 222–227.
- Jameson, A.; Kipper, B.; Ndiaye, A.; Schäfer, R.; Simons, J.; Weis, T.; and Zimmermann, D. 1994. Cooperating to be noncooperative: The dialog system PRACMA. In Nebel, B., and Dreschler-Fischer, L., eds., *KI-94: Advances in artificial intelligence*. Berlin: Springer. 106–117. Available from <http://zaphod.cs.uni-sb.de/KI/>.
- Jameson, A.; Schäfer, R.; Simons, J.; and Weis, T. 1995. Adaptive provision of evaluation-oriented information: Tasks and techniques. In Mellish, C. S., ed., *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*. San Mateo, CA: Morgan Kaufmann. 1886–1893. Available from <http://zaphod.cs.uni-sb.de/KI/>.
- Jameson, A. 1989. But what will the listener think? Belief ascription and image maintenance in dialog. In Kobsa, A., and Wahlster, W., eds., *User Models in Dialog Systems*. Berlin: Springer. 255–312.
- Joshi, A.; Webber, B.; and Weischedel, R. M. 1984. Living up to expectations: Computing expert responses. In *Proceedings of the Fourth National Conference on Artificial Intelligence*, 169–175.
- Ndiaye, A., and Jameson, A. 1994. Supporting flexibility and transmutability: Multi-agent processing and role-switching in a pragmatically oriented dialog system. In Jorrand, P., and Sgurev, V., eds., *Proceedings of the Sixth International Conference on Artificial Intelligence: Methodology, Systems, Applications*. Singapore: World Scientific Publishing. 381–390. Available from <http://zaphod.cs.uni-sb.de/KI/>.
- Ndiaye, A., and Jameson, A. 1996. Predictive role taking in dialog: Global anticipation feedback based on transmutability. In Carberry, S., and Zukerman, I., eds., *Proceedings of the Fifth International Conference on User Modeling*. Boston, MA: User Modeling, Inc. 137–144. Available from <http://zaphod.cs.uni-sb.de/KI/>.
- Ndiaye, A. 1996a. *Globale Antizipation in einem transmutierbaren Dialogsystem [Global Anticipation Feedback within a transmutable Dialog System]*. Ph.D. Dissertation, University of Saarbrücken. Forthcoming.
- Ndiaye, A. 1996b. Rollenübernahme in Dialogsystemen: Globale Antizipation in einem transmutierbaren Dialogsystem [Role taking in dialog systems: global anticipation feedback within a transmutable dialog system]. Technical Report, University of Saarbrücken. Available from <http://zaphod.cs.uni-sb.de/KI/>.
- Novak, H.-J. 1987. Strategies for generating coherent descriptions of object movements in street scenes. In Kempen, G., ed., *Natural Language Generation: New Results in Artificial Intelligence, Psychology, and Linguistics*. Dordrecht, The Netherlands: Nijhoff. 117–132.
- Oléron, P.; Beaudichon, J.; Danset-Léger, J.; Melot, A.-M.; Nguyen-Xuan, A.; and Winnykamen, F. 1981. *Savoirs et savoir-faire psychologiques chez l'enfant*. Bruxelles: Pierre Mardaga.
- Russell, S. J., and Norvig, P. 1995. *Artificial Intelligence: A Modern Approach*. Englewood Cliffs, NJ: Prentice-Hall.
- Schirra, J. R. J. 1995. *Bildbeschreibung als Verbindung von visuellem und sprachlichem Raum — Eine interdisziplinäre Untersuchung von Bildvorstellungen in einem Hörermodell [Scene Description as a Linking of Visual and Verbal Space: An Interdisciplinary Study of Visual Imagery in a Listener Model]*. DISKI – AI Dissertations. Sankt Augustin, Germany: Infix Verlag.
- Wahlster, W., and Kobsa, A. 1989. User models in dialog systems. In Kobsa, A., and Wahlster, W., eds., *User Models in Dialog Systems*. Berlin: Springer. 4–34.
- Zukerman, I. 1990. A predictive approach for the generation of rhetorical devices. *Computational Intelligence* 6:25–40.