# Combining Deictic Gestures and Natural Language
## for Referent Identification

Alfred Kobsa, Jürgen Allgayer, Carola Reddig, Norbert Reithinger
Dagmar Schmauks, Karin Harbusch, Wolfgang Wahlster
SFB 314: AI - Knowledge-Based Systems
University of Saarbrücken
D-6600 Saarbrücken 11
West Germany

## Abstract

In virtually all current natural-language dialog systems, users can only refer to objects by using linguistic descriptions. However, in human face-to-face conversation, participants frequently use various sorts of deictic gestures as well. In this paper, we will present the referent identification component of XTRA, a system for a natural-language access to expert systems. XTRA allows the user to combine NL input together with pointing gestures on the terminal screen in order to refer to objects on the display. Information about the location and type of this deictic gesture, as well as about the linguistic description of the referred object, the case frame, and the dialog memory are utilized for identifying the object. The system is tolerant in respect to impreciseness of both the deictic and the natural language input. The user can thereby refer to objects more easily, avoid referential failures, and employ vague everyday terms instead of precise technical notions.

*Keywords:* Deixis, referent identification, NP analysis, parsing

## 1. Introduction

Various aspects of *referent identification* by hearers have been investigated in the last few years: It has been studied as a process of noun phrase resolution and attribute comparison (Lipkis 1982), as a planned action (Cohen 1981, 84), as a process which depends on focus (Grosz 1981), context (Reichman 1981), the mutual beliefs shared between speaker and hearer (Clark & Marshall 1981) and the modality of linguistic communication (telephone vs. teletype, cf. Cohen 1984), and as a process which is prone to various sorts of conversational failure (Goodman 1985). In all of these studies, natural language is the only conversational medium. For identifying objects under discussion, the hearer can therefore only utilize the NL descriptions provided by the speaker, and information about the previous dialog and the task domain at hand.

In face-to-face conversation, however, participants also frequently use extralinguistic means for referent identification, in particular, various sorts of deictic gestures (such as pointing at something by ones hand, finger, pencil, head or eyes). One

---

may assume that this is done for simplifying and speeding up the identification process for both the hearer and the speaker, as well as avoiding referential failures. Certain technical innovations in the last few years (e.g., high-resolution graphic displays, window systems, touch-sensitive screens, input via a pointing device such as the mouse or the light-pen) have made it possible for computational linguistics to also experiment with and study a certain class of these deictic gestures, namely, tactile gestures for identifying objects on a terminal screen.

From an application-oriented perspective as well, such an ability is certainly a desirable characteristic for natural language dialog systems. In current systems, referring to visual objects involves the user either to employ unambiguous labels displayed together with the objects (cf. Phillips 1985), or purely linguistic descriptions which sometimes become rather complex (e.g. the "bright pink flat piece of hippopotamus face shape piece of plastic" in Goodman 1985). In Woods et al. (1979), a combination of deictic and natural language input has already been envisaged, but solely with restricted flexibility. Since an analyzer for pointing gestures is independent of a particular language, one might also consider transferring it to other NL dialog systems.

In this paper, we will present the referent identification component of XTRA, a system for a natural-language access to expert systems currently under development at the University of Saarbrücken. In its present application domain, XTRA is intended to assist a user in filling out his/her annual withholding tax adjustment form. The system will respond to terminological questions of the user, extract from the user's natural-language input the relevant data that is to be entered in the application form, and verbalize the inferences of the tax expert system. During the dialog, the relevant page of the application form is displayed on one window of the screen (for a simplified example, see Fig. 1; only the tax form is visible to the user).

For referring to single regions in the form, to the entities stored therein, or to larger regions which contain embedded regions, the user can employ linguistic descriptions (which we will call *descriptors*), pointing gestures with a pointing device (mouse), or both. From now on, the noun 'deictic' will refer to the use of a pointing device, and the term 'deictic expression' to the use of a descriptor plus a deictic (such as 'these deductibles' + deictic), or of a deictic alone.

In Bühler's (1982) terminology, the kind of deixis used in our situation is a *demonstratio ad oculos*. The objects on the display are visually observable, upon which the user and the system share a *common* visual field. In Clark & Marshall's (1981) terms, they are in a situation of *physical copresence*. Therefore, objects on the display need not be introduced by the user, but can immediately be referred to by a descriptor, a deictic, or both.

In many cases, however, neither kind of reference will be precise. Referential expressions, on the one hand, will often apply to more than one region in our form (as is the case when the user employs the term 'the deductibles' in order to refer to specific deductible sums such as dues for the membership in a professional organization). Deictic gestures, on the other hand, are also often imprecise in that they are not aimed at the region in which the user actually wants to refer to. Reasons for this might be inattentiveness, an oversized pointing device, or the user's intention not to hide the data entered in the respective field. Another factor of uncertainty is the *pars-pro-toto deictic*. In this case, the user points at an embedded region when actually intending to refer to a superordinated region. This is particularly the case when a form region is completely partitioned into a number of embedded sub-regions.

Therefore, in our model, we utilize several sources of information for identifying the region the user probably wants to refer to: the descriptor s/he uses, the location and the type of his/her pointing gesture, intrasentential context (case frames), and the dialog context. The information from each of these sources alone may be ambiguous or imprecise. Combined, however, they almost always allow for a precise identification of a referent.

## 2. Knowledge sources of the system

### 2.1. The tax form and the form hierarchy

During the dialog with the user, the system displays the relevant page of the income tax form on the terminal screen. As is illustrated in Fig. 1, such a form consists of a number of rectangular regions, which may themselves contain embedded regions, etc. We will abbreviate these regions by R1, R2, etc. The user can apply deictic operations to all regions.

For representing hierarchical relationships between regions, the system maintains an internal *form hierarchy*. Every region in the form has a corresponding element in the form hierarchy. Hierarchical relationships between form elements can then be expressed by father-son relationships within the form hierarchy. There are two reasons for introducing such a hierarchical order:

- *Geometrical reasons*: If region Rj is geometrically embedded in region Ri, then the element in the form hierarchy corresponding to Rj becomes a son of the element corresponding to Ri. An example is given in Fig. 1 where regions R2 and R3 are geometrically embedded in R1. Hence, their corresponding elements in the form hierarchy are subordinated to the element corresponding to R1.
- *Semantic reasons:* In many cases, there is a semantic coherence between regions in the form not directly expressed by the geometrical hierarchy. For example, see regions R15 and R16, and regions R33 and R34 in Fig. 1, which intuitively form units within the form for which no direct geometrical equivalents exist. Therefore, so-called *abstract regions* are introduced in the form hierarchy to which conceptually coherent regions can be connected. These regions even need not be geometrically adjacent and can be subordinated to

more than one abstract region. In Fig. 1, abstract regions are denoted by the symbol 'AR' (as e.g. AR48, the father of R15 and R16). It is not surprising that abstract units in the form hierarchy are often directly related to higher-level representational elements in the conceptual knowledge base of the system (cf. section 2.3.).

Moreover, we discern two types of bottom regions: *Label regions* contain the official inscriptions on the form (e.g. LR9 for 'Professional Expenses'), *value regions* contain the space for the user's data (e.g. VR28 for educational expenses). From now on, we will no longer distinguish between the form as displayed on the screen and the form hierarchy stored in the system. Since a close relationship between both structures exist, no problems will arise thereby.

### 2.2. The pointing gestures

Following Clark et al. (1983), we will call the region(s) at which the user pointed to the *demonstratum*, and the region which s/he intended to refer to the *referent*. Three cases can then be discerned:

a) The demonstratum is identical to the referent.
b) The demonstratum is a descendant of the referent (pars-pro-toto deixis). In this case, the referent may be a geometrical or an abstract region.
c) The demonstratum is geometrically adjacent to the referent. This occurs when the user points below the referent, to its right, etc. (e.g., by inattentiveness or because of not wanting to hide the data entered in the respective region).

In most cases, obviously, the location of a deictic does not identify its referent, but only restrains the set of possible referential candidates. Therefore, information about the pointing gesture usually has to be combined with information from other knowledge sources.

Another observation was that most subjects use several types of pointing gestures differing in exactness. Their choice seems to depend on the size of the target region. The larger the referent and the more sub-regions it contains, the vaguer is the pointing gesture. Therefore, our system allows the user to choose among several degrees of accuracy in his/her deictic. The user's decision, in turn, is taken into account when the system has to choose between referential candidates differing in size or to the degree of embedment (cf. section 3.1.2.).

### 2.3. The conceptual knowledge base

In our system, conceptual knowledge is represented by a frame-based language that shows a strong resemblance to Brachman's (1978) KL-ONE. The general part of the representation contains concepts and attribute descriptions of concepts. Attribute descriptions mainly consist of roles and value restrictions for possible role fillers. In Fig. 1, concepts are depicted by ovals and roles by small circles (the figure has been somewhat simplified). For object concepts (as e.g. 'MEMBERSHIP FEE' and 'ORGANIZATION'), attribute descriptions specify the properties of the objects described by the concept. For action concepts (as e.g. 'PHYSICAL TRANSFER', 'ADD' etc.), they specify the case frame.

General concepts can be ordered in a concept hierarchy, allowing the attribute descriptions of concepts to be inherited from the superordinated concepts. In Fig. 1, the bold arrows denote such superconcept relations. More specific concepts can be defined by introducing additional attribute descriptions or by further restraining the value restrictions of role fillers. It is possible for

**Conceptual Knowledge Base**

THING

COLLECTION
member
PERSON
agent
TO ADD
AMOUNT OF MONEY
result

= source
value
NUMBER
(R 14)
(R 16)
(R 18)
(R 34)
(R 36)

ORGANIZATION
recipient
MEMBER - SHIP FEE

name
STRING

CHARITABLE ORGANIZATION
PROFESSIONAL ORGANIZATION
PROF. ORG. MEMB. FEE
CHAR. ORG. MEMB. FEE

(AR 48)
(AR 51)

(AR 48)   80   40   (AR 51)

(R 16)   (R 34)

**Lexicon**

ADD
MEMBERSHIP FEE

ADD - VP
indir - obj
MEMB. FEE PP
MEMB. FEE NP
det          head
DEF - PL
MEMB. FEE NOUN

**Functional - Semantic Structures**

**Form Hierarchy**

LR9
R13 - - - - VR25
AR47
R14 - - - - VR28
AR60
R15 - - - - VR26
AR48
R16 - - - - VR29
R2
R17 - - - - VR27
AR49
R18 - - - - VR30
AR10
LR4
R1
AR12
R3
AR54
R31 — AR53
AR51
R33 - - - - VR43
R34 - - - - VR44
LR11
AR52
R35 - - - - VR45
R36 - - - - VR46
LR37

| R1 | Deductibles | | |
|---|---|---|---|
| R2 | Professional Expenses | | |
| R13 | Educational expenses | R14 | 250.00 |
| R15 | Professional organ. membership fee | R16 | 80.00 |
| R17 | Business trips | | |
| R3 | Other Deductibles | | |
| R31 | Charitable organizations | | |
| R33 | Membership fees | R34 | 40.00 |
| R35 | Donations | R36 | 20.00 |

"Can I add my annual $15.00 ACL dues to these membership fees?"

*Fig. 1: The knowledge sources of the system*

a concept to be subordinated to more than one superconcept, thus inheriting the properties of several superconcepts.

Natural-language input of the user containing new facts relevant for tax adjustment, as well as data entered directly into the form, causes structures of the general part to be *individualized*. Individualized concepts (depicted by ovals with lateral strokes in Fig. 1) and individualized attribute descriptions are thereby created. In Fig. 1, the individualized structures express the facts that the user spent $80 and $40 as professional organization and charitable organization membership fees, respectively.

Concepts and roles can be linked to elements in the form hierarchy if they conceptually correspond to a region in the form. In Fig. 1, for instance, the concept 'NUMBER' is associated with regions R16 and R34, amongst others, and the concept 'PROF.ORGAN.MEMB.FEE' with region AR48.

## 2.4 The functional-semantic structure

Before individualizations of the conceptual knowledge base are created, the natural-language input of the user is first mapped onto individualizations of the so-called *functional-semantic structure (FSS)*. The task of the FSS (cf. Allgayer & Reddig 1986) is to express the syntactic and semantic relationships between the constituents of the input sentence. It is also represented in a KL-ONE-like scheme. Amongst other things, the word stem entries in the lexicon determine which parts of the FSS are to be individualized. During this process, information about the location and the type of the occuring pointing gestures is assigned to the noun phrases to which they belong. Fig. 1 shows part of the individualized FSS generated by the input sentence.

The FSS forms the starting point for the referential analysis of the natural-language input, i.e. the mapping onto individualized structures of the conceptual knowledge base. This task is performed by an interpreter using appropriate mapping rules.

## 2.5. The dialog memory

Our current provisional approach is to regard the dialog memory as a structured list containing individualizations of the concepts in the conceptual knowledge base. When a referent is recognized as not having been mentioned before (because it is not contained in the dialog memory), the respective concept is individualized, linked to the referent, and entered as the most relevant element of the dialog memory. In Fig. 1 we assume that regions R16, R34, AR48 and AR51, amongst others, have been addressed before. Thus the concepts PROF.ORG.MEMB.FEE, CHAR.ORG.MEMB.FEE and NUMBER have been individualized and linked to these regions.

## 3. Referent identification processes

In a user's NL input, a deictic can be used at any position where a noun phrase or a (locative) adverbial phrase is to be expected. From a syntactic point of view, a deictic can serve two functions:

- it supplements a syntactically saturated description, i.e. takes the form of an additional attribute.

- it replaces a syntactically obligatory constituent (e.g. the head of a noun phrase).

The position of a deictic may be before, within, or after a noun phrase. Syntactic vicinity is taken into account if an ambiguity occurs in embedded noun phrases.

In the XTRA system, four sources of information are utilized in order to identify the referent of a deictic expression: The location of the user's pointing gesture, the descriptor s/he uses, case frame restrictions, and the contents of the dialog memory. The three former sources can be found in the functional-semantic structure, the latter source in the individualized part of the conceptual knowledge base. Referent identification, then, is performed in the following order:

a) Generation of potential referents by the most appropriate knowledge source. Source-specific partial plausibility values are thereby assigned to each generated candidate. Only deictic, descriptor and case frame are considered in this step, the dialog memory is only used in step (b).

b) Re-evaluation of each candidate by consecutively considering the information from all other knowledge sources.

c) Overall evaluation by considering all partial plausibility assignments; selection of the candidate with the highest plausibility factor.

In the following section we will describe how the most appropriate knowledge source for referent generation is selected and how referential candidates are generated. Since we are particularly concerned with referent identification through pointing gestures, we will only describe the referent generation strategy of the deixis analyzer (also cf. Allgayer 1986). For generating candidates through descriptors and case frames, we use the "classical" way leading from the lexicon via the FSS over to individualized concepts in the conceptual knowledge base and to the form hierarchy. In section 3.2., we then describe how the deixis analyzer re-evaluates candidates supplied by descriptor and case frame analysis, and how candidates generated by the deixis analyzer are re-evaluated by considering the information of all other knowledge sources. The example depicted in Fig. 1, to which we constantly refer to in the upcoming section, was chosen to demonstrate that, in many cases, all, or nearly all of these knowledge sources are necessary to correctly identify a referent.

### 3.1. Generating potential referents

#### 3.1.1. Deciding for the most appropriate knowledge source

In order to restrain the computational complexity of the identification process, it must be decided first whether referential candidates should be generated by analyzing the pointing gesture, the descriptor, or the case frame of the user's input. To assure that only a small number of candidates must be re-evaluated in the subsequent steps, it is certainly advisable to choose the knowledge source which yields the smallest set of plausible candidates that still contains the referent. The evaluation of each knowledge source is performed according to the following criteria:

- *Deixis:* The quality of a user's deictic for candidate generation is inversely proportional to the number of regions contained in the demonstratum and the number of ancestors of the demonstratum. A deictic to R3 in Fig. 1, for instance, will yield less candidates than a deictic to R34.

- *Descriptor:* If a descriptor does not contain a head, it cannot be used for candidate generation. Otherwise, its quality is inversely proportional to the number of subconcepts of its conceptual representation and the number of regions linked to these concepts. E.g., for the representation in Fig. 1, the descriptor 'number' will yield by far more candidates than the descriptor 'membership fee'.

- *Case frame:* The quality of a case restriction for referent generation depends on the quality of the selection restriction concept of the corresponding role in the conceptual knowledge base. This quality can be computed in the previous manner mentioned. In Fig. 1, the selection restrictions for the ADD concept do not seem to be profitable for candidate generation.

### 3.1.2. Generating candidates by analyzing the user's pointing gesture

As was mentioned above, our system allows for the use of several types of deictic gestures differing in precision. A so-called *deictic field* is associated with each type of pointing gesture, its size corresponding to the degree of exactness of the deictic. An example for three different types of pointing gestures is given in Fig. 2.
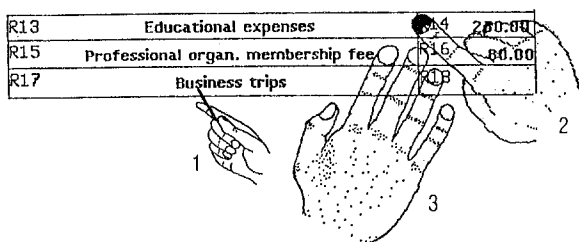


*Fig.2:* Three types of pointing gestures

A deictic field may either be completely contained in a basic region (as is the case for deictic 1 in Fig. 2) or overlap two or more basic regions (deictics 2 and 3, respectively). All basic regions that are overlapped by a deictic field serve as first referential candidates in our system. The ratio of that part of a region covered by a deictic field in relation to the size of the total region yields the plausibility value for the region. Deictic 3, for instance, generates R18, R16, R17 and R15 as first candidates, in order of descending plausibility (cf. Allgayer 1986).

In a second step, the system accounts for the possibility of pars-pro-toto deixis. All regions semantically or geometrically superordinated to any of the current candidates are also considered as candidates. The plausibility assignment of a superordinated region depends on its type, the plausibility of its candidate subregions, and the type of pointing gesture employed by the user (the vaguer the pointing gesture, the higher is the plausibility of the superordinated regions). In Fig. 2, regions AR49 and AR48 would be added in the case of deictic 3, both with higher plausibility than any of the first candidates. This upward *propagation* through the hierarchy can be applied iteratively, yielding even more candidates (the valuation function smoothly declines thereby to render high-level regions less plausible). The resulting set of candidates has to be re-evaluated by the processes described below.

### 3.2. Re-evaluating the set of candidates

### 3.2.1. Re-evaluation by analysis of the pointing gesture

If the optimization process of section 3.1.1. decided that descriptor or case frame analysis were the most appropriate knowledge sources for candidate generation, analysis of the deictic is employed in our system for re-evaluating the candidates supplied by these components. This evaluation process is rather similar to candidate generation described above. For example,

see Fig. 1 (we assume that the deictic in this example is the same as deictic 3 in Fig. 2): If the desciptor analyzer generated AR48, AR51, R16 and R34 as potential referents (since the descriptor was 'membership fee', see below), the deixis component would assign high plausibility values to the former, and very low ones to the latter.

### 3.2.2. Re-evaluation by descriptor analysis

This process determines to what extent the conceptual representation of the descriptor applies to the current candidates. Each candidate is tested as to whether the representation of the descriptor, a subconcept of this representation, or (if existent) the restriction concept of the value slot of one of these concepts is linked to the candidate. The more concepts in between the representation of the descriptor and the linked subconcept, the lower the new partial plausibility assignment. Let us assume for our example in Fig. 1 that the deixis analyzer, in order of decreasing plausibility, has generated regions AR49, AR48, R18, R16, R17 and R15 as potential referents. If the descriptor is 'these membership fees', the descriptor analysis will prefer AR48 and R16, since a subconcept of the representation of this descriptor is linked to AR48, and the restriction concept of its value slot is linked to R16.

### 3.2.3. Re-evaluation by case frame analysis

This process determines to what extent the selection restriction concept of the respective slot in the conceptual representation of the verb applies to the referential candidates under investigation. This evaluation process is performed almost identically to that of the descriptor. In our example, high plausibility would be attributed to regions R16 and R18, since the concept NUMBER (the restriction concept of the relevant slot of the concept ADD) is linked to these regions.

### 3.2.4. Restriction by dialog memory

This process determines whether a referent has recently been mentioned by checking whether or not an individualized concept connected with it is contained in the dialog memory. The better the position of such an individualized concept in the list, the better the plausibility of the candidate. In Fig. 1, we assume that both the professional and the charitable society memberships and their values have been addressed just recently. Therefore, in our example, high plausibility values are assigned to regions R16 and AR48. The overall evaluation will then select R16, it having obtained the highest total plausibility.

### 4. Discussion

Our system demonstrates that spatial deixis is a valuable source of information for identifying referents which also can be investigated and utilized in natural language dialog systems with pictoral display. Three reasons sum up the advantages of using pointing gestures: They save the speaker the generation, and the hearer the analysis of complex referential descriptions and thus simplify the natural-language dialog; they often allow for reference in situations in which linguistic reference is simply not possible (think of referring to one out of a dozen similar objects); and they permit the speaker to be vague, imprecise, or ambiguous, and to use everyday terms instead of precise technical terms unknown to him/her.

In natural-language dialog systems, deixis analysis can be combined well with standard methods for referent identification.

Some of the identification processes (e.g., tests with case frame, descriptor and dialog memory) are rather similar to the classical methods used for anaphora and ellipsis resolution. Others, such as the generation and evaluation of candidates by the deixis analyzer, are typical with respect to this particular kind of conversational medium.

It should be pointed out, however, that our environment for spatial deixis is, in several ways, somewhat simpler than those occurring in person-to-person dialogs (cf. Schmauks 1986). The deictic field is only two-dimensional, and the objects that can be pointed at are clearly separated from each other. Compared to real-life situations, the number of possible referents is relatively small. "Left" and "right" mean the same thing for the user and the system (which is not the case, e.g., in face-to-face conversation). However, this relative simplicity need not be a drawback. Instead, one might regard our environment as a study *in vitro*, eliminating a number of uncertainty factors so that the essential characteristics of spatial deixis become more salient.

Another question is whether the deictic behavior of subjects who use a pointing device is the same as that of subjects who touch the display with their fingers (and thus, whether deixis via a pointing device is a valid simulation of tactile deixis). One might argue, e.g., that people point more precisely with a mouse than with their fingers, or vice versa. We are currently conducting an informal experiment to answer these questions. In any case, only the propagation functions are perhaps affected by a change of the deictic medium, whereas the referent identification processes will remain the same.

Attempts are currently being made to also integrate visual and conceptual salience in our model (cf. Clark et al. 1983). When a pointing gesture is ambiguous, it appears that regions set off by bold frame or coloring, as well as regions containing important data for the task domain are preferred. We expect this preference to be taken into account in the evaluation processes of the deixis analyzer. Another possible extension which we would like to investigate is in replacing the strategy described in section 3.1.1. by a certain form of incremental referent identification. There is strong empirical evidence (e.g. Goodman 1985) that people begin with referent identification immediately after receiving initial information about it, instead of waiting until the speaker's referential act is terminated. Since all components described above are strictly separated, it appears basically possible to also use them in an incremental identification process. In one-processor systems, however, great care must be taken that the knowledge source first adressed does not block the system by generating too many candidates. Therefore, some process controlling will be necessary, either by ressource limitation or by taking into account the heuristics listed in section 3.1.1.

## References

Allgayer, J. (1986): Eine Graphikkomponente zur Integration von Zeigehandlungen in natürlichsprachliche KI-Systeme. 16. GI-Jahrestagung, Berlin, FRG (in print).

Allgayer, J. und C. Reddig (1986): Systemkonzeption zur Verarbeitung kombinierter sprachlicher und gestischer Referentenbeschreibungen. SFB 314, Dept. of Computer Science, University of Saarbrücken, FR Germany.

Brachman, R. J. (1978): A Structural Paradigm for Representing Knowledge. Report No. 3605, Bolt, Beranek and Newman Inc., Cambridge, MA.

Bühler, K. (1982): The Deictic Field of Language and Deictic Words. Abridged translation of K. Bühler (1934): Sprachtheorie, part 2, chapters 7 and 8. In: R. J. Jarvella and W. Klein, eds.: Speech, Place, and Action. Chichester etc.: Wiley.

Clark, H. H. and C. R. Marshall (1981): Definite Reference and Mutual Knowledge. In: A. K. Joshi, B. L. Webber and I. A. Sag, eds.: Elements of Discourse Understanding. Cambridge: Cambridge Univ. Press.

Clark, H. H., R. Schreuder and S. Buttrick (1983): Common Ground and the Understanding of Demonstrative Reference. Journal of Verbal Learning and Verbal Behavior 22, 245-258.

Cohen, P. R. (1981): The Need for Referent Identification as a Planned Action. Proceedings of the 7th International Joint Conference on Artificial Intelligence, Vancouver, Cda., 31-36.

Cohen, P. R. (1984): The Pragmatics of Referring and the Modality of Communication. Computational Linguistics 10, 97-146.

Goodman, B. A. (1985): Repairing Reference Identification Failures by Relaxation. Proceedings of the 23rd ACL Meeting, Chicago, IL, 204-217.

Grosz, B. J. (1981): Focusing and Description in Natural Language Dialogues. In: A. K. Joshi, B. L. Webber and I. A. Sag, eds.: Elements of Discourse Understanding. Cambridge: Cambridge Univ. Press.

Phillips, B., M. J. Freiling, J. H. Alexander, S. L. Messick, S. Rehfuss and S. Nicholl (1985): An Eclectic Approach to Building Natural Language Interfaces. Proceedings of the 23rd ACL Meeting, Chicago, IL, 254-261.

Lipkis, Thomas (1982): A KL-ONE Classifier. Proceedings of the 1981 KL-ONE Workshop. Report No. 4842, Bolt, Beranek and Newman Inc., Cambridge, MA, 128-145.

Reichman, R. (1981): Plain Speaking: A Theory and Grammar of Spontaneous Discourse. Report No. 4681, Bolt, Beranek and Newman Inc., Cambridge, MA.

Schmauks, D. (1986) : Formulardeixis und ihre Simulation auf dem Bildschirm. Ein Überblick aus linguistischer Sicht. Memo No. 4, Sonderforschungsbereich 314, Dept. of Computer Science, University of Saarbrücken, FRG.

Woods, W. A., R. J. Brachman, R. J. Bobrow, R. R. Cohen and J. W. Klovstad (1979): Research in Natural Language Understanding: Annual Report. TR 4274, Bolt, Beranek & Newman, Cambridge, MA.