AD-A175 282

# INTEGRATING SYNTAX, SEMANTICS, AND DISCOURSE DARPA NATURAL LANGUAGE UNDERSTANDING PROGRAM

## R&D STATUS REPORT
## SDC -- A BURROUGHS COMPANY

The views and conclusions contained in this document are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the U.S. Government.

DTIC FILE COPY

DTIC
ELECTE
DEC 2 2 1986
S      D
E

CONTRACT DEPARTMENT

NOV 21 1986

SDC - A BURROUGHS CO.
CUSTOM PRODUCTS GROUP

86  12  04  069

## 1. Description of Progress

### 1.1. Lexical Lookup

The lexical reader has been modified to allow the PUNDIT Natural Language Processing system to recognize abbreviations ending in periods. PUNDIT previously had not been able to distinguish between periods ending abbreviations and periods ending sentences; consequently, it was not able to parse sentences with internal punctuation.

### 1.2. Grammar

#### 1.2.1. Re-test of TESTA Corpus

Effort over the last quarter has focused on obtaining results from an interactive procedure for "learning" selectional information. We have chosen a sub-set of the CASREP corpus for our initial experiments. The TESTA corpus contains 31 CASREP sentences which generate a well-formed regularized intermediate syntactic representation (ISR) (e.g., one containing no unresolved lambda terms). We have written a definite-clause grammar (DCG) to translate ISRs into sublanguage patterns, for use as data on domain-specific co-occurrence patterns.

The purpose of re-running the entire TESTA corpus was twofold: first, testing the DCG on as many ISRs as possible was an excellent method of debugging the grammar; and second, parsing a large number of CASREP sentences with the selection pattern generator switched on yielded valuable information concerning the reduction in the search space traversed by the parser. If an incorrect parse can be ruled out early in the parsing, say, by failing an incorrect noun pattern, the parser will not consume unnecessary time and space by chasing down what would only later prove to be a blind alley.

#### 1.2.2. Search Focus

This work has produced some striking results: using the TESTA corpus described above, 30 out of 31 reports received a correct first parse when selection was guided interactively by the linguist, as opposed to only 18 out of 31 correct first parses without selection. Two sentences which received NO correct parse without selection both receive correct first parses using interactive selection. Other interesting results are a decrease from an average of 15 parses/sentence without selection to approximately 1.4 parses/sentence with selection. There is also some "focusing effect" which reduces the search required to obtain the correct parse by about 10% - 15% using selection. Of course, the search of the entire search space was reduced much more (indicated by the decrease in the average number of parses per sentence), but this was not explicitly measured. (Search focus is measured by the ratio of nodes attached divided by nodes in the correct parse tree.) These results are the basis for a technical report that will be prepared before the end of the year.

The results of running the TESTA corpus of 31 sentences are as follows:

-- average number of parses found WITH DCG    =  1.45161
   average number of parses found WITHOUT DCG = 14.9677

-- search focus ratio with dcg/without DCG    =   0.892455

-- 2 sentences receiving NO correct parse WITHOUT DCG
        received CORRECT FIRST parse WITH DCG.

-- 13 sentences receiving INCORRECT first parse WITHOUT DCG
        received  CORRECT   first parse WITH DCG

### 1.3. Syntax/Semantics Interaction

Several meetings have been held to identify a discrete yet useful arena for developing a more flexible control strategy in syntactic and semantic processing of text. Currently, Pundit's control mechanism flows from the syntactic module to the semantic module. (Semantic data is used as a filter on the syntax, via the recently implemented selectional component, but this filtering does not change the *control* structure). The meetings have focused

on the distribution of lexical information, which is dispersed both in the lexicon used by the syntactic parser and in the semantic lexical entries with their associated syntactic mapping rules. One proposal that has emerged would exploit the semantic component's potential to predict verb transitivity on the basis of selectional information and thereby handle the phenomenon of transitivity alternation, i.e., verbs which have transitive and intransitive uses. This would be useful because many verbs in English exhibit transitivity alternations. For example, the verb *melt* is transitive in the sentence *The candle melts the ice cube*, and intransitive in the sentence *The ice cube melts*. In some instances, it is possible to predict a verb's transitivity or intransitivity after the subject and the verb have been parsed in order to prune the search space for the verb's object. For example, since *ice* cannot fill the agent of *melt*, it is possible for semantics to fill the theme role with the subject immediately at this point in processing the sentence; since the agent, if present, must occur as the subject of the sentence, it is now possible to predict that there will be no second argument, and hence to instruct the semantic module not to look for a direct object. Within a particular subdomain, there appear to be many constraints on the instantiation of thematic roles (beyond the general semantic considerations which allow us to reject *ice* as agent of *melt*), and thus such interaction between syntax and semantics may be possible with many verbs. Focusing on this problem has the dual advantage of extending the current coverage and also restricting the experimentation with the new control strategy to a carefully delimited arena.

## 1.4. Temporal Analysis

The module which performs temporal analysis has been restructured to provide a cleaner separation between the analysis of a predication, and the analysis of temporal adverbials which modify the predication. In the process, the existing analysis of temporal relations has been extended. Also, the restructuring will provide the means for extending the coverage to more temporal adverbials. In addition, the flow of information between the time component and the two semantics components (clause semantics and noun phrase semantics) has been augmented. A detailed report describing the temporal component of PUNDIT has been prepared and is being reviewed internally.

## 1.5. Semantic Interpreter

Semantics rules have been added for adjectives and prepositions. The interpreter has also been restructured in order to interpret adjectives and prepositions both as clausal predications (e.g., *The pressure was low*) and as noun modifiers (*The low pressure*). Finally, the clause semantics interpreter has been extended to cover a class of clausal constructions which follow the nouns they modify, e.g., the participial phrase *decreasing below 60 psig* in *Oil pressure decreasing below 60 psig caused sudden failure*.

## 1.6. Editor for Semantics Rules

The semantic interpreter makes use of three kinds of rules: rules associating semantic decompositions with predicating words; rules mapping the abstract semantic arguments of a decomposition to syntactic constituents; and rules checking the consistency of syntactic fillers for semantic arguments against the domain model. Work has begun on a new editing tool to enforce the consistency of these three rules sets with one another, to automate development of new rules, and to ease the process of porting to new domains. The design for the semantic rule editor has been completed and partially implemented.

## 1.7. Environment

Some switches (run/environment options) have been added to the PUNDIT system. We can now simply allow the parse to fail with unknown lexical items, as opposed to inv ing the lexical entry procedure. Another switch toggles the selection mechanism. A third switch allows the pri   of more detailed time information. The system has also been given the capability of running in unattended or    ch mode for testing large amounts of data.

## 1.8. Facilities

We have successfully ported the PUNDIT sy em to the Texas Instruments Explorer and the Sun 3 workstation. We have received Release 7.0 of the Symbolics operating system.

## 2. Change in Key Personnel

Bonnie Webber, a member of the computer science faculty at the University of Pennsylvania, has joined our group on a part-time (20%) basis. Dr. Webber is working on problems in time, events, and reference.

## 3. Summary of Substantive Information from Meetings and Conferences

### 3.1. Professional Meetings Attended

### 3.1.1. AAAI-86

Lynette Hirschman, Martha Palmer, Deborah Dahl, Francois Lang, Marcia Linebarger, and Leslie Riley attended the annual meeting of the American Association for Artificial Intelligence in Philadelphia. Dahl presented a paper, "Focusing and Reference Resolution in PUNDIT", describing the reference resolution component of PUNDIT.

### 3.1.2. Logic Programming Conference

John Dowding, Francois Lang and Lynette Hirschman attended the Symposium on Logic Programming held in Salt Lake City on September 22 - 25. John attended a tutorial on Building Prolog Interpreters and Compilers, and Francois one on Applications of Parallelism to Logic Programming. We were able to talk at length with several representatives from Quintus about possible extensions and changes to their Prolog and development environment.

### 3.1.3. Meeting with Robert Simpson

Robert Simpson visited SDC on August 13 for an in-depth review of the PUNDIT system. The presentation focused on the treatment of nominalizations and sentence fragments; it also included several demonstrations of the PUNDIT system running on a Xerox 1109.

## 4. Problems Encountered and/or Anticipated

The status of the follow-on contract needs to be confirmed.

## 5. Action Required by the Government

## 6. Fiscal Status

(1) Amount currently provided on contract:
    $ 672,833 (funded)                      $683,105 (contract value)

(2) Expenditures and commitments to date:
    $ 471,600

(3) Funds required to complete work:
    $ 201,233