

Panayotis Stamatoopoulos, Maria Katsouraki
NCC "DIOCRITUS", Athens

George Fillokypros
University of Athens

A NATURAL LANGUAGE SYSTEM
FOR DATA BASE QUERIES

In this paper a natural language system is presented for processing types of sentences to be used on a data base environment. The system is implemented in PROLOG using Transformational Grammar rules and introduces a verb oriented logical structure for the output representation, contrary to existing natural language processors, which employ a quantifier based approach. As for the analysis stage of the implementation, efficiency is obtained by employing a two-pass processor. The implementation is adapted to a Greek language subset.

1. Introduction

Interacting directly with data bases can be achieved either by special structured query languages (SQL, QUE, SQUARE etc.) or with natural languages which conform to the user's native language. It is with this latter approach that we are dealing in this paper presenting the description of a natural language processor adapted to Greek for data base queries. Following some recent developments in natural language analysis [2,4,5,8,9], we have employed logic to deal both with the underlying formalism and the programming environment. Specifically the system is implemented entirely in Prolog [1,7], a powerful programming language based in logic. For these purposes a specially oriented Prolog system has been developed (Interpreter) implemented on a HP 21 MX minicomputer of our laboratory [6].

The general structure of the Natural Language Processor is shown in Fig.1. The processor handles restricted natural Greek language using rules of Heteromorphosis Grammar (HG) formalism [3] under normal form. The application domain, a geographical data base, is only for demonstration purposes given that by no means restricts the adaptation to a variety of applications. There are mainly two ways in which the proposed system differs from previous work [2,4,5,8,9] and compares favourably with them:

- The way the natural language analysis is performed in two phases rather than one [2,5] or three [9].
- The logical structure of the output tree representation which is "verb" rather than "quantifier" oriented.

Besides, it represents the first adaptation of Logic Natural Language framework to Greek.

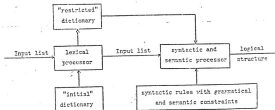


Fig.1 Natural Language Processor

2. The Natural Language World

2.1. Linguistic Coverage

The natural language subset which the system correctly interprets is determined by the Vocabulary and the Syntactic Structures it accepts. We present below the vocabulary scope and the types of sentences which are within the scope of our natural language subset.

2.1.1. Vocabulary

The vocabulary includes both general and application dependent words. Each word is classified according to its part of speech class, which in turn determines a combination of grammatical, syntactic and semantic feature-categories characteristic of that class. Below we present some of the vocabulary classes together with some examples. To this end we use the following abbreviations: g:gender, n:number, c:case, df:domain field, sct:standard output type, x:variable, ls:logical structure, d:degree. It is understood that certain of the above categories (e.g. gender, number, case) assume several values (e.g. masculine, feminine, neutral/singular, plural/nominative, genitive, accusative).

Class	Categories
a) Definite Articles e.g. the (O)	{g,n,c} {masculine,singular,nominative}
b) Relative Pronouns e.g. which (OPCIA)	{g,n,c} {feminine,singular,nominative}
c) Prepositions e.g. with (E)	{s} {s}
d) Nouns e.g. country (XPA)	{g,n,c,df,x,ls} {feminine,singular,accusative,country,x,country(x)}
e) Proper Nouns e.g. Sarah (DONTAMG)	{g,n,c,df,sct} {masculine,singular,nominative,river,Santa}
f) Comparative Adjectives e.g. biggest (NEGALSTFM)	{g,n,c,df,d} {feminine,singular,accusative,country,can}
g) Verbs e.g. borders (STORXVII)	{n,df,ls,df,ls,preposition,c,ls} {singular,country,ls,country,ls,E,accusative, border(ls,ls)}

Remarks: In the case of verbs, besides the domain field and the logical structure of the subject (df_1, ls_1), the domain fields (df_2), logical structures (ls_2), prepositions and cases (c_1) for each complement constitute distinct categories. Finally, ls_3 stands for the logical structure of the sentence.

2.1.2. Syntactic Structures

The system implemented recognizes affirmative sentences of the non-elliptic type, whose nouns can be either proper nouns or nouns, which in turn can be determined by comparative adjectives, quantifiers or relative sentences. Quant-

flers are the class of definite and indefinite articles, indefinite pronouns, numerals or structure forms of the kind "all the".

2.2. Application Domain

The application domain chosen is that of a wider district geography, specifically that of the Balkan States. It is essentially a relational data base organization with all information coded in the form of Prolog facts and rules, as shown in Appendix I. We believe that it should be relatively easy to adapt to different applications.

1. The Original Language Processor

3.1. The Lexical Processor

According to the previous section, all the information regarding the categories for each word must be included in the Prolog system to build an "initial" dictionary. But Greek language (and other languages too) have the characteristic that the values assumed by some linguistic categories for many word classes are too many to be efficiently manipulated as a whole. Indicatively we note that in Greek nouns can assume 6 different values, proper nouns 3, adjectives 18 etc.

The space problems resulting from such an implementation are faced by employing the Lexical Processor, which preprocesses the sentence to be analyzed consulting the "initial" dictionary and giving as output a "restricted" dictionary with entries which are relevant only to the words of the specific sentence. The "initial" dictionary entries are Prolog facts with the same predicate, taking the form:

Dictionary(Words Class,Word,[Category₁,Category₂,.....,Category_n]) (1)

In Appendix II a small part of the "initial" dictionary is shown. The Lexical Processor assigns to each word of the sentence to be analyzed all the facts of the form (1) whose second argument coincides with that word. In this way a "restricted" dictionary is produced with entries facts of the form:

Word Class(Word,Category₁,Category₂,.....,Category_n) (2)

The above predication form is easily amenable to the syntactic and semantic processing of the next stage. The Prolog program implementing the Lexical Processor together with the "restricted" dictionary for the following sentence are given in Appendix III and IV respectively.

Sentence: O DONAVAZES PIAZEPETI THS MEGALYTERHS EKRA S OPOIA SYNDREKHTEI TA THS ELLADHS

I.e. Danube flows through the biggest country which borders Greece.

Given that the linguistic categories of a word can be singled out solely on the basis of its appearance, even the "restricted" dictionary will include facts which will be of no further use.

3.2. The Syntactic and Semantic Processor

Heteromorphosis grammar rules under normal form, directly translated into Prolog are the rules employed by the Syntactic and Semantic Processor for defining some

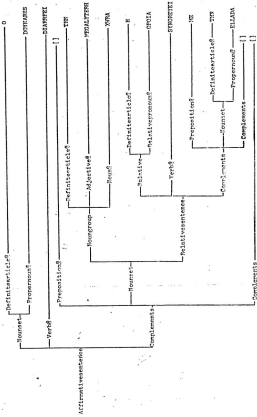


FIG. 2 Deep structure tree

formal properties of the natural language. The Syntactic and Semantic Processor, consulting the "restricted" dictionary output of the Lexical Processor, performs syntactic analysis while verifying validation of the natural language grammar rules and of the semantic features of the application domain. Thus, it verifies, for example, agreement between several linguistic categories (gender, number, case) among an adjective and a noun or (and) coincidence among a verb and its subject as far as the number category is concerned etc. Metamorphosis grammar rules under normal form can be directly translated into Prolog rules as it is shown below for the case of an affirmative sentence:

```
Affirmative sentence("Structure1") -> Nouncat("Number,nominative,"Structure2),
    Verb3("Number,"structure2,"Structure3,"Structure1),Complements("Structure1")
+Affirmative sentence("List1","List4,"Structure1)-Nouncat("List1","List2,"Number,
    nominative,"Structure2)-Verb("List2,"List3,"Number,"Structure2,"Structure3,
    "Structure1)-Complements("List3,"List4,"Structure3)
```

Within the above rule syntactic agreement restrictions and semantic constraints can be satisfied. In Appendix V we state all the rules which are necessary for analyzing the above mentioned sentence, while the corresponding deep structure tree is shown in Fig.2. Similar rules will enforce the analysis of other types of sentences.

3.2.1. Logical Structure

Confirming with the example-case of an affirmative sentence, the output logical structure after the application of the rules is of the form $P(\text{Set}_1, \text{Set}_2, \dots, \text{Set}_n)$ where P denotes the predicate associated with the main verb of the sentence, Set_2 denotes the set describing the subject and $\text{Set}_1, \dots, \text{Set}_n$ the sets describing the complements of the verb, in case they exist. The above sets can be either a list of proper nouns, a structure of the form $\text{set}(\text{variable}, \text{quantifier}, \text{property})$, or finally a variable previously introduced by some structure. The "quantifier" argument can be of the form $\text{the}(\text{number}, \text{howmany}, \text{cooperative})$, where $\text{howmany} \Rightarrow \delta, 1, 2, 3, \dots$ and $\text{cooperative} \Rightarrow \delta, \text{max}, \text{min}$ or $\text{some}(\text{number})$ or $n(\text{howmany})$ or all . Finally the "property" argument constitutes a set of properties associated with "and", which are satisfied by the variable introduced by the quantifier. Thus, the output structure of the example-case will take the form:

```
flow([Ismube], set(x, the(sin, delta, max), and(country(x), border(x, [Greece])))
```

Appendix I

```
+Area(Ugoslavia, 833504)
+Area(Greece, 131641)
+Length(Borako, 2850)
+Border([Greece, Ugoslavia])
```

+Flow(River, Uroslovia)
 +River(X)-Length(X, Y)
 +Country(X)-Area(X, Y)
 +Border(X, Y)-Border?(X, Y)
 +Border(X, Y)-Border?(Y, X)

Appendix II

+Dictionary(Definitearticle, O, [mas, sin, non])
 +Dictionary(Relativepronoun, OPODA, [fer, sin, non])
 +Dictionary(Propernoun, ILLADA, [fer, sin, acc, country, Greece])
 +Dictionary(Noun, XIPA, [fer, sin, acc, country, "X, country(X)])
 +Dictionary(Adjective, ISALYTESI, [fer, sin, acc, country, mas])
 +Dictionary(Verb, ENHOTEIPI, [sin, dom(country, "S1"), [agr(dom(country, "S2), "X, acc)],
 border("S1, "S2)])

Appendix III

+Process(sentence-Consult(Dictionary))-LineWrite("Give input sentence:")-Line
 -Readsentence("L")-Searchforword("L")
 +Readsentence(["X"|"L"])-Read("X")-Wordsentence("L")
 +Wordsentence([])
 +Searchforword([])
 +Searchforword(["X"|"L"])-Process("X")-Searchforword("L")
 +Process("X")-Dictionary("P, "X, "L")-Unit("Y, ["S", "X"|"L"])-Verb("Y")-Accus("P")-Fall
 +Process("X")
 +Get("X")-X/-Fall
 +Get("X")

Appendix IV

+Definitearticle(O, mas, sin, non)
 +Definitearticle(TEN, fer, sin, acc)
 +Definitearticle(U, fer, sin, non)
 +Relativepronoun(OPODA, fer, sin, acc)
 +Relativepronoun(OPODA, nest, pl, acc)
 +Preposition("E")
 +Propernoun(DOYRANOS, mas, sin, nom, river, Danube)
 +Propernoun(ILLADA, fer, sin, acc, country, Greece)
 +Noun(XIPA, fer, sin, acc, country, "X, country(X))
 +Adjective("ISALYTESI, fer, sin, acc, country, mas)
 +Verb(ENHOTEIPI, sin, dom(country, "S1"), [scr(dom(country, "S2), "X, acc)], border("S1, "S2))
 +Verb(DIARHESI, sin, dom(river, "S1"), [scr(dom(country, "S2), "X, acc)], flow("S1, "S2))

Appendix V

- +Affirmativesentence(*L1,*L4,*S1)-Nounset(*L1,*L2,*N,non,*S2)-Verb(*L2,*L3,*V,*S2,*S3,*S1)-Complement(*L3,*L4,*S3)
- +Nounset(*L1,*L3,*N,*C,dor[*D,[*S]])-Definitearticle(*L1,*L2,*C,*S,*C)-
-Pronoun(*L2,*L3,*C,*N,*C,*D,*S)
- +Nounset(*L2,*L3,*N,*C,*S1)-Noungroup(*L1,*L2,*C,*S,*C,*D,*S,*S2)
-Relative(*L2,*L3,*C,*N,*D,*S,*S,*S1,*S2)
- +Complement(*L1,*L1,[*])
- +Complement(*L1,*L1,[*S,*P,*C],[*N])-Preposition(*L1,*L2,*P)-Nounset(*L2,*L3,*C,*S)-Complement(*L3,*L4,*S)
- +Noungroup(*L1,*L4,*S,*C,*S,*S,tho[*S,*S],[*S],[*S])-Definitearticle(*L1,*L2,*C,*N,*C)-Relative(*L2,*L3,*C,*N,*C,*S,*S)-Noun(*L2,*L4,*C,*S,*S,*S,*S,*S,*S)
- +Relative(*L1,*L4,*N,*S,*S,*S,*S,rel[*S,rel[*S,*C,rel[*S1,*S2]])-Noun
-Relative(*L1,*L2,*C,*S,[*N,non])-Verb(*L2,*L3,*V,dor[*S,*S])*S3,*S2)
-Complement(*L3,*L4,*S3)
- +Relative(*L1,*L3,*C,*S,[*C])-Definitearticle(*L1,*L2,*C,*S,*C)-Relativepronoun
(*L2,*L3,*C,*S,*C)

References

1. M. Chockola, G. Wallish, "Programming in Prolog", Springer-Verlag, Berlin, Heidelberg, New York, 1981.
2. M. Coelho, "A program concerning in Portuguese providing a library service", Ph.D. Thesis, University of Edinburgh, 1979.
3. A. Colmerauer, "Les grammires de "etamorphose", Groupe d'Intelligence Artificielle, Universit  d'Aix-Marseille, 1975.
4. A. Colmerauer, "An interesting subset of natural language", in "Logic Programming", K. Clark, S.-A. Hirahara (eds.), Academic Press, 1982, pp.45-66.
5. V. Dehl, "Translating Spanish into logic through logic", American Journal of Computational Linguistics, vol.7, num.3, July-September 1981, pp.149-164.
6. M. Katsouraki, M. Gavallas, P. Stamatopoulos, A. Arvillias, "A flexible Prolog interpreter for minicomputers" (in Greek), Proc. of 1st Panhellenic Conference in Informatics, vol.2, Athens, 1980, pp.235-244.
7. P. Roussel, "PROLOG: Manuel de r f rence et d'utilisation", Groupe d'Intelligence Artificielle, Universit  d'Aix-Marseille, 1975.
8. S. Szpakowicz, "Syntactic analysis of written Polish", Inf UW Report, Nr 64, 1971.
9. D. Warren, P. Pereira, "An efficient easily adaptable system for interpreting natural language queries", BAI Research Report, No.155, University of Edinburgh, 1981.