

QUERYING DATABASES IN NATURAL GREEK

Panagiotis Stamatopoulos, Alexis Arviliadis, Maria Kaliziani, George Psilodimitrou (1)

Computer Center, NROPS Demokritos,
 Agria Paraskevi, Athens, Greece

(2) University of Athens,
 Panepistimioupolis, Athens, Greece

A natural language system for understanding sentences stated in Greek is presented in this paper. These sentences are considered to apply to a database environment, in which they act as database queries. The system is not restricted to a specific domain, as it consists of a general part and a special one. Adding different versions of the latter part makes it possible for the system to work on various applications. The accessible subset of Greek is described in the meta-morphosis grammar formalism. Every sentence belonging to the Greek subset is translated into a semantic structure whose interpretation, given a certain database, yields the appropriate answer to the user. The system has been implemented in Prolog, a programming language based on the principles of logic programming.

1. INTRODUCTION

Existing data manipulation languages (SQL, GUC, QUEL, QUELRE etc.) provide a formal way for retrieving information from databases. While these languages are desirable for programmers, a non-specialized person needs a more user friendly interface to interact with databases. Evidently, natural language understanding systems, which recognize the user's native language, fulfill the friendliness requirement.

Recently developed natural language understanding systems employ morphological [12], definite clause [9], extraposition [10] or other variations of logic grammars [18] to define the subset of the language being processed. Such systems have been created for English [5,7,8,14], French [4], German [11], Spanish [6] and Portuguese [3]. As far as the implementation is concerned the only programming tool used is the language Prolog [1,2].

All of the previously mentioned systems translate each natural language sentence into a formal structure, possibly through successive transformations to intermediate representations. This structure is either a Prolog goal [5,14] or a codification of the sentence according to a well defined semantic framework [3,4,6,7,8,11]. In order to yield an answer to the sentence, the surrounding Prolog system executes the Prolog goal in the first case, while in the second case a Prolog program is used to interpret the codification of the sentence.

Some systems [3,4,6] are based on the semantic framework proposed by Alain Colmerauer [13]. In this framework the quantifiers (articles, pronouns, numerals etc.) are of special importance. Moreover, a three valued logic is defined, where, besides the well known values "true" and "false", the value "undefined" is introduced. This value corresponds to sen-

tences that are neither "true" nor "false" because of the existence of false presuppositions.

Following the developments already reported, we designed and implemented a natural language understanding system for database access adapted to Greek. For this purpose we introduced an entirely new semantic framework, to define the translation of every acceptable Greek sentence into a semantic structure. The most remarkable characteristics of our semantic framework is the distinction between definite and indefinite quantifiers, a feature which does not appear in Colmerauer's framework. This distinction leads on the one hand to the correct interpretation of sentences containing false presuppositions and on the other to the efficient implementation of the natural language understanding system.

It is worthwhile noting that the design of the system was also highly affected by the specific natural language being analyzed. We mention the delimitation system of Greek, the partial freedom related to word order and the semantic differences existing between Greek and other natural languages. For the definition of the accessible Greek subset we used metamorphosis grammars. There has been no attempt to handle unrestricted Greek language, since, as it is generally admitted for every natural language, such a goal is unfeasible, at least in the foreseeable future.

The system was implemented using a Prolog interpreter developed on an IBM AT (8 MHz) of our laboratory and was functionally tested on a couple of application domains. The first domain is a subset of the Balkan States geography and the second one is related to radioactivity measurements. Given the portability of the system it could be equally reapplied to any other world.

2. THE PROPOSED SEMANTIC FRAMEWORK

2.1. The model of the surrounding world

Natural language understanding is meaningful only when a specific world domain is considered, whose data should be organized according to a given model.

In the model we propose the atoms of the world are classified into groups. There are also relations between atoms belonging to certain groups. A subgroup is a subset of a group whose atoms relate to a given atom of another group.

The groups of the surrounding world are organized in a group tree. Figure 1 shows the group tree for the Balkan States geography world.



Figure 1

2.2. Semantic structures

Every sentence of the considered Greek subset is translated into a Finite ordered tree, called the semantic structure of the sentence. The nodes of a semantic structure are special keywords, codifications of atoms, groups or relations and names of variables. The syntactic components of a sentence are mapped into semantic structures, which, when properly arranged, form the semantic structure of the sentence.

In what follows we present informally how different types of syntactic expressions correspond to semantic structures.

a) The conjunction of common or proper nouns corresponding to atoms is mapped into a collection.

Example :
 ... η Ελλάδα, η Βουλγαρία και η Αλβανία ...
 (... Greece, Bulgaria and Albania ...)

Semantic structure :
 collection([["Ελλάδα"], ["Βουλγαρία"],
 ["Αλβανία"]])

b) Common nouns corresponding to groups are mapped into group structures.

Example :
 ... areas ...

[... areas ...]

Semantic structure :
 group_structure([numeral, area], X)

c) Common nouns which are preceded by definite quantifiers are mapped into sets. The absence of relative sentences designating a common noun is coded as an empty structure. As definite quantifiers are considered, the definite articles possibly followed by a cardinal, an ordinal, a comparative adjective, a cardinal and a comparative adjective, the adjective "uberoi" (mean) or the adjective "kovvaxia" (total). Moreover, the adjective "oloi" (all) followed by a definite article form a definite quantifier.

Example :
 ... οι 5 πρωτεύουσες ...
 (... the 5 capitals ...)

Semantic structure :
 set([the(N), X, group_structure([region, city,
 capital], X), empty_structure])

The collections and the sets are called complement structures.

d) Common nouns designated by anything coded as a complement structure are mapped into subgroup structures.

Example :
 ... πρωτεύουσα της Αθήνας ...
 (... population of Athens ...)

Semantic structure :
 subgroup_structure([numeral, population], X,
 collection(["Αθήνα"]])

The group and subgroup structures are called selection structures.

e) Declarative sentences whose quantifiers (in case they exist) are definite, are mapped into relational structures.

Example :
 Η μεγαλύτερη χώρα συνορεύει με τη Βουλγαρία.
 (The biggest country borders Bulgaria.)

Semantic structure :
 relational_structure(call_is_all, yes, borders,
 [set([min([1, bigger, [numeral, area]], X,
 group_structure([region, country], X),
 empty_structure), collection([
 ["Βουλγαρία"]])])])

f) Declarative sentences containing at least one indefinite quantifier are mapped into logical structures. Indefinite quantifiers are the cardinals, the indefinite pronouns "adexia" (a), "adexini" (some), "adexos" (two) and "adexi" (each) or one of the expressions "kovvaxi-drovo" (at least), "to oloto" (at most), "horo-drovo" and "horo" (more than) or "horo-drovo" and "horo" (less than) followed by a cardinal.

Example :
 3 ποταμοί διασχίζουν τον Τσεχικό.
 (3 rivers flow through Turkey.)

Semantic structure :
 logical_structure(quantifier(3),X,group_structure(
 Driver),X,empty_structure,
 relational_structure(all_to_all,yes,flows,
 collection(3X),collection(
 [{"Turkey"}]))))

In case of more than one indefinite quantifier in a sentence the following rules apply.

1. Every indefinite quantifier belonging to the main part of a sentence dominates all the indefinite quantifiers of possible relative sentences.
2. When a noun preceded by indefinite quantifier is followed by a designator containing another indefinite quantifier, then the designator's quantifier dominates the noun's quantifier.
3. If the subject and the objects of a verb are introduced by indefinite quantifiers, then the quantifiers' hierarchy decreases from left to right, according to the natural order of their appearance.

g) **Compound declarative sentences** are mapped into conjunction structures.

Example :
 Η Ελλάδα είναι η τέταρτη χώρα με το μεγαλύτερο
 σε μέγεθος.
 (Greece is the fourth country and does not border
 Russia.)

Semantic structure :
 conjunction_structure(relational_structure(
 one_to_one,yes,is,collection(["Ελλάδα"])),
 self_order(4,bigger,[natural,area]),X,
 group_structure([region,country],X),
 empty_structure()),relational_structure(
 all_to_all,no,borders,collection(
 [{"Ελλάδα"}]),collection(["Ρωσία"])))

The relational, logical and conjunction structures are called declarative structures.

h) **Relative sentences** are mapped into declarative structures.

Example :
 ... οι πόλεις οι οποίες βρίσκονται στα Αλβανικά
 ...
 (... the cities which belong to Albania ...)

Semantic structure :
 self(1,X,group_structure([region,city],Y),X),
 relational_structure(all_to_all,yes,belongs,
 collection(3X),collection(
 [{"Αλβανία"}]))))

i) **Interrogative sentences** introduced with the pronoun "ποιος" (which) are mapped into "which_ane"

or into "which_ane" structures.

Example :
 Ποιο ποτάμι διασχίζει τον Γιουγκοσλαβικό;
 (Which rivers flow through Yugoslavia?)

Semantic structure :
 which(X,group_structure([river],X),
 empty_structure,relational_structure(
 all_to_all,yes,flows,collection(3X),
 collection(["Γιουγκοσλαβία"])))

j) **Interrogative sentences** introduced with the pronoun "ποσοί" (howmany) are mapped into "howmany_ane" or into "howmany_ane" structures.

Example :
 Πόσοι ποταμοί διασχίζουν την Ελλάδα;
 (Howmany are the cities of Greece?)

Semantic structure :
 howmany_ane(quant(all,X,subgroup_structure(
 [region,city],Y),X,collection(["Ελλάδα"])),
 empty_structure))

Every "which", "which_ane", "howmany" and "howmany_ane" structure is called **interrogative structure**.

2.3. Values of semantic structures

A variable appearing in a semantic structure is **closed** if it is introduced by a selection structure which is a sub-tree of the semantic structure, otherwise it is **free**.

A semantic structure is **closed** if it contains only bound variables, otherwise it is **open**. Every complete sentence is mapped into a closed semantic structure, while parts of the sentence may be represented by either open or closed semantic structures.

For a given status of the surrounding world a value can be assigned to every closed semantic structure. The possible values of a closed semantic structure are the following.

1. null_value()
2. where: to: true, false or undefined
3. collection()
 - where: c: a list containing lists of atoms
 - number(n):
 - where: n: a non negative integer number

2.4. On Galsterauer's semantic framework

In the semantic framework proposed by Aigun Galsterauer [12] there is no distinction between definite and indefinite quantifiers, although definite quantifiers introduce in most cases independently computable sets. This unified approach of quantifiers fails to interpret correctly sentences belonging to a certain class. For example, consider the sentences:

Η μάχη του Σίου πραγματοποιήθηκε τον Απρίλιο
 διαρκώντας 3 μέρες.

(the country which is smaller than Albania borders 3 countries.)

and

3 ποταμοί διαρρέουν το κράτος που είναι μικρότερο από την Αλβανία.

(3 countries border the country which is smaller than Albania.)

Given that in the Balkan States geography there is no country smaller than Albania, according to Colmerauer's semantic framework the formal representation of the first sentence is assigned the value "undefined", while that of the second one is assigned the value "false", although the two sentences are semantically equivalent. This phenomenon appears between sentences having the following features. Both have the same main verb representing a generic relation. In the first sentence the main quantifier of the subject is definite, introducing a false presupposition, and that of the object is indefinite. The second sentence is derived by mutually interchanging the subject and the object of the first one.

The system we implemented tries to compute for both sentences the set of "the country which is smaller than Albania". As this attempt fails, because of the existence of a false presupposition, both semantic structures of the previous sentences are assigned the value "truth,value undefined", thus yielding the correct answer.

Moreover, Colmerauer's semantic framework fails to lead to efficient implementation of natural language understanding systems in terms of evaluating the computable sets introduced by definite quantifiers. Let's examine the sentence :

2 ποταμοί διαρρέουν το κράτος που συνορεύει με τη Βουλγαρία, οι ποταμοί των οποίων οι εκβολές είναι μεγαλύτερες από τις εκβολές της Βουλγαρίας.
(2 rivers flow through the country which borders the countries whose capitals are bigger than the capital of Bulgaria.)

In order to assign a value to the semantic formula of the above sentence, any natural language understanding system following Colmerauer's proposal has to compute R times the set of "the country which borders the countries whose capitals are bigger than the capital of Bulgaria", $R \times R$ times the set of "the countries whose capitals are bigger than the capital of Bulgaria" and $R \times R \times R$ times the set of "the capital of Bulgaria", where R is the number of rivers and R the number of countries in the surrounding world.

Instead, in the natural language understanding system we implemented, based on the semantic framework we propose, each of the previous sets is computed only once. This is derived from the fact that the assignment of values to semantic structures containing nested sets fol-

lows a bottom up procedure. Considering the specific example, the third set is initially computed and the result is used for the computation of the second set. The new value obtained is passed to the first set, whose computation allows for a value to be assigned to the semantic structure of the original sentence.

3. IMPLEMENTATION OF THE SYSTEM

3.1. Organization

Figure 2 shows the general outline of the natural language understanding system we designed and implemented. The system consists of three main parts, the supervisor, the natural language processor and the semantic structure processor.

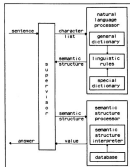


Figure 2

Each sentence posed to the system is transformed by the supervisor to a list containing the characters of the sentence. This list is passed to the natural language processor which creates the corresponding semantic structures, if the sentence belongs to the acceptable subset of Greek. The semantic structure processor analyzes the semantic structures just created and assigns to it a value which is given to the supervisor. Finally, the supervisor formulates the appropriate answer to the user's original sentence.

3.3. Semantic structure processor

The semantic structure processor consists of the semantic structure interpreter and the database.

The Prolog clauses composing the semantic structure interpreter define logically a set of procedures which are responsible for the evaluation of the semantic structures produced by the natural language processor. A semantic structure may be considered as a program which is processed by the semantic structure interpreter to derive a value as output, while abstracting the present status of the surrounding world as input. The procedure which is initially activated when the supervisor triggers the semantic structure processor is defined as follows.

- ```

1. value(Semantic_structure, Value) :-
 interpretive_structure(
 Semantic_structure, Value),
 !.
2. value(Semantic_structure, Value) :-
 declarative_structure(
 Semantic_structure, Value),
 !.
3. value(Semantic_structure, truth_value(
 undefined)).

```

The database is the formal description of the surrounding world. It is in this part of the system that the relations, the groups and the subgroups of the world are defined.

The relations are distinguished into actual and virtual, depending on the way of their implementation. An actual relation is defined by a set of Prolog facts, while for the definition of a virtual relation a set of Prolog rules is required. For each relation there exists a certain schema, which is defined by the name of the relation and the names of the groups participating in the relation. We consider the following relations for the implementation of the Balkan States geography.

- ```

1. belongs(City, Country)
2. is_region(Region, Population)
3. is_city(City, Population, Country)
4. borders(Country, Country)
5. flows(River, Country)
6. is_country(Country, Area, Population)
7. is_capital(Capital, Population,
  Country)
8. is_sample_city(Sample_City, Population,
  Country)
9. is_river(River, Length)
  
```

Relations 1 to 3 constitute the virtual database and relations 4 to 9 the actual one.

A small sample of Prolog clauses implementing relations of the database is as follows.

- ```

1. is_city(City, Population, Country) :-
 is_capital(City, Population, Country).

```

- ```

2. is_city(City, Population, Country) :-
  is_sample_city(City, Population, Country).
3. borders('Bulgaria', 'Rumania').
4. flows('Danubius', 'Rumania').
5. is_country('Etiopia', 10144, 890000).
6. is_capital('Addis', 100100, 'Etiopia').
7. is_sample_city('Sofia', 10100, 'ARBUL').
8. is_river('Aisou', 200).
  
```

Every group corresponding to a leaf of the group tree is defined in the database in a way shown by the examples below.

- ```

1. group([region, city, capital], (Cap, Val)) :-
 is_capital(Capital, Population, Country).
2. group([numeral, population], (Population,
 Region)) :-
 is_region(Region, Population).

```

The subgroups of the surrounding world are also defined in the database in a way similar to the definition of groups.

- ```

1. subgroup([river], (River, Country)) :-
  flows(River, Country).
2. subgroup([numeral, area], (Area, Country),
  Country) :-
  is_country(Country, Area, Population).
  
```

3.4. Portability of the system

The modular design of the natural language understanding system we presented allows for an easy implementation of various applications. Four parts of the system, i.e. the supervisor, the linguistic rules, the general dictionary and the semantic structure interpreter, constitute the nucleus as they are application independent. The sizes of these parts are 48, 285, 265 and 109 Prolog clauses respectively, giving a total of 607 clauses. The only parts of the system that need to be created in case we adopt another domain are the special dictionary and the database. These latter parts form the surroundings of the nucleus. In the Balkan States geography domain the special dictionary and the database contain 81 and 88 Prolog clauses respectively. In the radioactivity measurements domain the special dictionary consists of 295 clauses and the database of 543 clauses.

4. SAMPLE QUERIES

Indicative queries posed to the system are presented below along with the answers obtained and the execution times needed for the IBM AT configuration employed.

- ```

#| User :
#| Prolog>country(whatever) at its greatest or
 smallest pop.
 (Rumania borders the smallest country in
 terms of population.)

```

```

System :
#| system>city(what)
 (The sentence is correct)

```

Natural language processor time : 3.08 sec  
 Semantic structure processor time : 0.80 sec  
 Total time : 3.88 sec

b) User :  
 Ποια ποταμοί της Βουλγαρίας διαρρέουν 3 χώρες;  
 (Which rivers of Bulgaria flow through 3 countries?)

System :  
 Πάρισις 2  
 (There exist 2)

Natural language processor time : 3.90 sec  
 Semantic structure processor time : 2.03 sec  
 Total time : 5.93 sec

c) User :  
 Ποια πόλεις (αποτελούν) και είναι της χώρας η  
 πρωτεύουσα είναι πρωτεύουσα και όχι πρωτεύου-  
 σα της Ρουμανίας;  
 (Which cities belong to the country whose capi-  
 tal is bigger than the capital of Romania?)

System :  
 Αθήνα, Θεσσαλονίκη  
 (Athens, Thessaloniki)

Natural language processor time : 9.78 sec  
 Semantic structure processor time : 5.27 sec  
 Total time : 15.05 sec

d) User :  
 Ποια είναι η γηνη έκταση του κράτους του Δουβλίου  
 σε τετραγωνικά μίλια;  
 (Which is the mean area of the countries that  
 border the third country?)

System :  
 683388

Natural language processor time : 23.88 sec  
 Semantic structure processor time : 2.91 sec  
 Total time : 26.79 sec

e) User :  
 Ποια είναι τα μήκη των ποταμών τα οποία διαρρέουν  
 μέσω των χωρών της Γαλλίας και της Γερμανίας;  
 (Which are the lengths of the rivers which flow  
 through a country that is smaller than Rumania  
 in terms of area?)

System :  
 389(A) μήλ., 2850; 420(Β) μήλ., 530; (Γ) μήλ.  
 380; Αίνας, 2850; Δούναβη, 530; (Υ) μήλ.

Natural language processor time : 28.31 sec  
 Semantic structure processor time : 4.67 sec  
 Total time : 32.98 sec

f) User :  
 Ποια είναι οι εθνότητες των οποίων τον πληθυσμό  
 ο Εύρος και οι εθνότητες Ουγγαρίας υπέρ των οποίων  
 τον πληθυσμό ο Γαλλικός είναι ο ίδιος με τον  
 πληθυσμό της Γερμανίας και πρωτεύουσας και  
 500000 κατοίκους και πρωτεύουσας και 1000000  
 κατοίκους;

(Which are the areas of the countries that  
 Έυρος flows through and border the country  
 whose capital is the city whose population  
 is greater than 500000 citizens and less  
 than 1000000 citizens?)

System :  
 631944; (Α) μήλ., 280576; Τσεχία  
 (131944; Γερμανία, 280576; Τουρκία)

Natural language processor time : 44.38 sec  
 Semantic structure processor time : 11.20 sec  
 Total time : 55.58 sec

## 5. CONCLUSION

We have described a natural language understand-  
 ing system which makes it possible for a  
 subset of Greek to act as a query language for  
 databases. This system is based upon a novel  
 semantic framework, which seems to be more  
 advantageous compared to that of Alain Colmerauer  
 in terms of both correctness and implementa-  
 tion. As far as the efficiency of the system  
 is concerned, it is sufficiently satisfactory  
 as even relatively complex queries are answered  
 in less than a minute on an ordinary personal  
 computer using a Prolog interpreter without  
 special performance characteristics. The execu-  
 tion times could be dramatically improved if  
 a Prolog compiler were used instead, which  
 would result in a 2-3 orders of magnitude bet-  
 ter performance.

## REFERENCES

- [1] Glockoin, W. and Mellish, C., Programming in Prolog, Springer Verlag, Berlin, 1981.
- [2] Sterling, L. and Shapiro, E., The art of Prolog, MIT Press, Cambridge, 1986.
- [3] Coelho, H., A program conversing in Portuguese providing a library service, Ph.D. Thesis, University of Edinburgh, 1979.
- [4] Perlehand, L., Consultation en Français d'une banque de données sur fichiers et mise en place du système Prolog nécessaire, Thèse de 3ème cycle, Université Aix-Marseille, 1985.
- [5] Warren, S. and Pereira, F., An efficient easily adaptable system for interpreting natural language queries, IAI Research paper, No. 158, University of Edinburgh, 1981.
- [6] Dahl, V., Translating Spanish into logic through logic, American journal of computational linguistics, Vol. 3, No. 3, pp. 149-164, 1985.
- [7] McCord, R., Using slots and modifiers in logic grammars for natural language, Artificial intelligence, Vol. 18, No. 3, pp. 327-367, 1982.
- [8] Hadley, R., Shadow: a natural language query analyzer, Computers and mathematics with applications, Vol. 11, No. 3, pp. 481-504, 1988.
- [9] Pereira, F. and Warren, S., Definite

clause grammars for language analysis - A survey of the formalism and a comparison with augmented transition networks, *Artificial intelligence*, vol. 13, no. 3, pp. 291-376, 1980.

- [10] Pereira, F., Extrapolation grammars, *American journal of computational linguistics*, vol. 7, no. 4, pp. 243-256, 1981.
- [11] Schmidt, G., Logic based natural language processing, in: Datt, V. and Saint-Dizier, P. (eds.), *Natural language understanding and logic programming*, pp. 207-219, North Holland, 1989.
- [12] Colmerauer, A., Metamorphic grammars, in: Bolc, L. (ed.), *Natural language communication with computers, Lecture notes in computer science*, No. 83, pp. 133-169, Springer Verlag, Berlin, 1978.
- [13] Colmerauer, A., An interesting subset of natural language, in: Clark, G. and Tarslund S. (eds.), *Logic programming*, pp. 45-66, Academic Press, London, 1982.
- [14] Filgueiras, R., A kernel for a general natural language interface, *Proceedings of logic programming workshop*, pp. 419-436, Albufeira, Portugal, 1983.
- [15] Datt, V., Current trends in logic grammars, *Proceedings of logic programming workshop*, pp. 578-607, Albufeira, Portugal, 1983.