
A CONSTRAINT BASED SYSTEM FOR LEXICAL KNOWLEDGE RETRIEVAL*

STOYAN MIHOV

This paper is concerned with the lexical knowledge retrieval system created at the Linguistic Modelling Laboratory. The main goal of the system is to provide a powerful and comfortable interface for lexical knowledge retrieval from large morphological dictionary. To achieve this a constraint based approach is applied that leads to a very effective algorithm. The algorithm for query building, which is also used for retrieving of general grammatical knowledge, is presented in details. In our opinion this method is very suitable for knowledge retrieval in domains with complex and irregular classifications.

Keywords: lexical knowledge retrieval, morphological dictionary, query building

1991/95 Math. Subject Classification: 68T50

1. INTRODUCTION

Recently, many systems containing large amount of lexical knowledge have been built. They make use of different approaches for processing and knowledge representation. Detailed study of the problem is presented in [4] and [1].

The Linguistic Modelling Laboratory is working on a Large Morphological Dictionary that will cover most of the wordforms in modern Bulgarian (see [3]). Now the system contains grammatical information for more than 500 000 wordforms and is systematically upgraded. The grammatical information is structured in Features Structures (refer to [5] for a good introduction). It is known that Feature Structures (FSs) are *de facto* standard for representing linguistic information. That is because

* Lecture presented at the Fourth Logical Biennial, Gjuleitchitza, September 12–14, 1996.

they allow comfortable description of knowledge with complex classifications and many irregularities.

The program which interfaces the lexical knowledge in our system is called kernel program. We can think the knowledge as a set containing all 500 000 feature structures corresponding to the information of wordforms. The kernel input is a Feature Structure — the query constraint that should be satisfied. The output contains all Feature Structures included in our knowledge base, unifiable with the input. See Fig. 1 for some examples of the kernel functioning. It is clear that it is very inefficient to keep all 500 000 FSs in the memory and to check every one for the query constraint satisfaction. The kernel program uses a synthesizing algorithm. The grammatical information of the output Feature Structures is built by the unification of certain *basic constraints*. Each basic constraint is a Feature Structure which corresponds to a feature-value pair of the grammatical information. In Fig. 2 there are some examples of Feature Structures¹ which present basic constraints corresponding to feature-values.

This paper reveals in details our approach for query building. It is shown how general grammatical knowledge is retrieved by our procedure. In Section 2 we describe the problem for query building. Then in Section 3 the algorithm is presented and explained. In Section 4 some details about the implementation are described. Some comments on the possibilities for generalization are given in the conclusion.

2. PROBLEM DESCRIPTION

The main goal of the system is to provide an appropriate way for knowledge retrieval. There are very different necessities concerning system interaction. For example, in syntactic analysing systems the grammar information should be extracted after the input of wordform, as it is shown in the first example in Fig. 1. In Natural Language Generation the system should find the wordform by processing input of a stem and some grammatical information. The professional linguist should be able to extract all wordforms (or stems, endings...) which satisfy a certain constraint. On the one hand, there should be no limitations on the contents of a query constraint. But, on the other hand, many of the queries are inconsistent with the grammatical knowledge. For example, there are no items in the knowledge base, which satisfy the following FS:

$$\left[\begin{array}{ll} \text{lexeme info} & \left[\text{gradability} \quad \text{gradable} \right] \\ \text{grammeme info} & \left[\text{tense} \quad \text{imperfect} \right] \end{array} \right]$$

That is because only adjectives and adverbs are gradable and they do not allow a tense characterization. It is meaningless to process a query which we can a priori

¹ For a clearer presentation, in the example we are noting the unbound anonymous variables explicitly with ‘_’.

Input FS (Query)	Output FSs
<p>[wordform 'BOEHHA']</p>	<p>wordform stem 'BOEHHA'</p> <p>lexeme info 'BOEHEH'</p> <p>grammeme info [part of speech adjective ungradable]</p> <p>ending 'A'</p> <p>alternation [degree positive] [gender feminine] [number singular] [definiteness indefinite]</p> <p>[from 'E'] [to ''] [pos 2]</p>
<p>[stem 'CTABA']</p> <p>grammeme info [tense present] [person 3]</p>	<p>wordform stem 'CTABAT'</p> <p>lexeme info 'CTABA'</p> <p>grammeme info [part of speech verb] [tense present] [person 3] [number plural]</p> <p>...</p>

Fig. 1. Examples of the kernel functioning

Feature-value	Corresponding basic constraint																																																		
part of speech — noun	<table border="0"> <tr> <td data-bbox="644 129 820 165">wordform</td> <td data-bbox="911 138 927 165">-</td> <td data-bbox="927 129 1366 165"></td> </tr> <tr> <td data-bbox="644 174 719 210">stem</td> <td data-bbox="911 183 927 210">-</td> <td data-bbox="927 174 1366 210"></td> </tr> <tr> <td data-bbox="644 219 746 255">ending</td> <td data-bbox="911 228 927 255">-</td> <td data-bbox="927 219 1366 255"></td> </tr> <tr> <td data-bbox="644 479 820 515" rowspan="10">lexeme info</td> <td data-bbox="911 273 1145 309">part of speech</td> <td data-bbox="1187 273 1257 309">noun</td> </tr> <tr> <td data-bbox="911 318 1086 353">noun type</td> <td data-bbox="1187 318 1203 353">-</td> </tr> <tr> <td data-bbox="911 362 1118 398">animateness</td> <td data-bbox="1187 362 1203 398">-</td> </tr> <tr> <td data-bbox="911 407 1102 443">humanness</td> <td data-bbox="1187 407 1203 443">-</td> </tr> <tr> <td data-bbox="911 452 1034 488">gender</td> <td data-bbox="1187 452 1203 488">-</td> </tr> <tr> <td data-bbox="911 497 1082 533">verb type</td> <td data-bbox="1187 497 1347 533">not defined</td> </tr> <tr> <td data-bbox="911 542 1066 577">gradable</td> <td data-bbox="1187 542 1347 577">not defined</td> </tr> <tr> <td data-bbox="911 586 1102 622">transitivity</td> <td data-bbox="1187 586 1347 622">not defined</td> </tr> <tr> <td data-bbox="911 631 1134 667">numeral type</td> <td data-bbox="1187 631 1347 667">not defined</td> </tr> <tr> <td data-bbox="911 676 959 712">...</td> <td data-bbox="1187 676 1203 712"></td> <td data-bbox="927 676 1366 712"></td> </tr> <tr> <td data-bbox="644 882 884 918" rowspan="6">grammeme info</td> <td data-bbox="911 743 1050 779">number</td> <td data-bbox="927 743 1366 779"></td> </tr> <tr> <td data-bbox="911 788 1114 824">definiteness</td> <td data-bbox="1145 788 1161 824">-</td> </tr> <tr> <td data-bbox="911 833 1114 869">article form</td> <td data-bbox="1145 833 1161 869">-</td> </tr> <tr> <td data-bbox="911 878 1018 913">tense</td> <td data-bbox="1145 878 1305 913">not defined</td> </tr> <tr> <td data-bbox="911 922 1034 958">person</td> <td data-bbox="1145 922 1305 958">not defined</td> </tr> <tr> <td data-bbox="911 967 1002 1003">case</td> <td data-bbox="1145 967 1305 1003">not defined</td> </tr> <tr> <td data-bbox="911 1012 959 1048">...</td> <td data-bbox="1187 1012 1203 1048"></td> <td data-bbox="927 1012 1366 1048"></td> </tr> <tr> <td data-bbox="644 1079 692 1115">...</td> <td data-bbox="927 1079 1366 1115"></td> <td data-bbox="927 1079 1366 1115"></td> </tr> </table>	wordform	-		stem	-		ending	-		lexeme info	part of speech	noun	noun type	-	animateness	-	humanness	-	gender	-	verb type	not defined	gradable	not defined	transitivity	not defined	numeral type	not defined	...			grammeme info	number		definiteness	-	article form	-	tense	not defined	person	not defined	case	not defined		
wordform	-																																																		
stem	-																																																		
ending	-																																																		
lexeme info	part of speech	noun																																																	
	noun type	-																																																	
	animateness	-																																																	
	humanness	-																																																	
	gender	-																																																	
	verb type	not defined																																																	
	gradable	not defined																																																	
	transitivity	not defined																																																	
	numeral type	not defined																																																	
	...																																																		
grammeme info	number																																																		
	definiteness	-																																																	
	article form	-																																																	
	tense	not defined																																																	
	person	not defined																																																	
	case	not defined																																																	
...																																																			
...																																																			
non-finite form — participle	<table border="0"> <tr> <td data-bbox="644 1140 794 1176">wordform</td> <td data-bbox="911 1149 927 1176">-</td> <td data-bbox="927 1140 1366 1176"></td> </tr> <tr> <td data-bbox="644 1184 719 1220">stem</td> <td data-bbox="911 1193 927 1220">-</td> <td data-bbox="927 1184 1366 1220"></td> </tr> <tr> <td data-bbox="644 1229 746 1265">ending</td> <td data-bbox="911 1238 927 1265">-</td> <td data-bbox="927 1229 1366 1265"></td> </tr> <tr> <td data-bbox="644 1422 820 1458" rowspan="7">lexeme info</td> <td data-bbox="911 1283 1145 1319">part of speech</td> <td data-bbox="1187 1283 1257 1319">verb</td> </tr> <tr> <td data-bbox="911 1328 1082 1364">verb type</td> <td data-bbox="1187 1328 1203 1364">-</td> </tr> <tr> <td data-bbox="911 1373 1086 1408">noun type</td> <td data-bbox="1187 1373 1347 1408">not defined</td> </tr> <tr> <td data-bbox="911 1417 1066 1453">gradable</td> <td data-bbox="1187 1417 1347 1453">not defined</td> </tr> <tr> <td data-bbox="911 1462 1102 1498">transitivity</td> <td data-bbox="1187 1462 1347 1498">not defined</td> </tr> <tr> <td data-bbox="911 1507 1134 1543">numeral type</td> <td data-bbox="1187 1507 1347 1543">not defined</td> </tr> <tr> <td data-bbox="911 1552 959 1588">...</td> <td data-bbox="1187 1552 1203 1588"></td> <td data-bbox="927 1552 1366 1588"></td> </tr> <tr> <td data-bbox="644 1767 884 1803" rowspan="6">grammeme info</td> <td data-bbox="911 1619 1082 1655">finiteness</td> <td data-bbox="1187 1619 1347 1655">non-finite</td> </tr> <tr> <td data-bbox="911 1664 1161 1700">non-finite form</td> <td data-bbox="1187 1664 1337 1700">participle</td> </tr> <tr> <td data-bbox="911 1709 1050 1744">number</td> <td data-bbox="1187 1709 1203 1744">-</td> </tr> <tr> <td data-bbox="911 1753 1018 1789">tense</td> <td data-bbox="1187 1753 1203 1789">-</td> </tr> <tr> <td data-bbox="911 1798 1114 1834">definiteness</td> <td data-bbox="1187 1798 1347 1834">not defined</td> </tr> <tr> <td data-bbox="911 1843 1002 1879">case</td> <td data-bbox="1187 1843 1347 1879">not defined</td> </tr> <tr> <td data-bbox="911 1888 959 1924">...</td> <td data-bbox="1187 1888 1203 1924"></td> <td data-bbox="927 1888 1366 1924"></td> </tr> <tr> <td data-bbox="644 1955 692 1991">...</td> <td data-bbox="927 1955 1366 1991"></td> <td data-bbox="927 1955 1366 1991"></td> </tr> </table>	wordform	-		stem	-		ending	-		lexeme info	part of speech	verb	verb type	-	noun type	not defined	gradable	not defined	transitivity	not defined	numeral type	not defined	...			grammeme info	finiteness	non-finite	non-finite form	participle	number	-	tense	-	definiteness	not defined	case	not defined								
wordform	-																																																		
stem	-																																																		
ending	-																																																		
lexeme info	part of speech	verb																																																	
	verb type	-																																																	
	noun type	not defined																																																	
	gradable	not defined																																																	
	transitivity	not defined																																																	
	numeral type	not defined																																																	
	...																																																		
grammeme info	finiteness	non-finite																																																	
	non-finite form	participle																																																	
	number	-																																																	
	tense	-																																																	
	definiteness	not defined																																																	
	case	not defined																																																	
...																																																			
...																																																			

Fig. 2. Example of feature-value pairs and the corresponding *basic constraints*

consider as inconsistent. That is why we have to implement a more sophisticated query building algorithm.

We think that the best way is to build the query incrementally. That means that the user should be able to specify the query step by step, selecting a feature-value pair. The purpose is to receive information which feature-value pairs are acceptable (do not lead to an inconsistent query) after each step. If we switch those feature-values as unacceptable, the user would not be able to build an inconsistent query. Moreover, the user will receive general grammatical information about the consistency of some constraint combinations.

Another advantage of the system would be if the algorithm automatically selects some feature-value pairs which are derivable from the existing query specification. For example, if the user selects 'part of speech – verb' and 'gender – masculine', then the 'finiteness' has to be 'non-finite'.

To achieve the above mentioned requirements for the query building procedure, a deduction procedure has to be implemented. This procedure should check on each step the consistence of every feature-value pair with the query and should select the feature-values which are deducible from the specification of the query.

3. THE QUERY BUILDING ALGORITHM

The simplest way to fulfil the above algorithm specification is to generate a list of all possible combinations of feature-values. Unfortunately, there are thousands of possible combinations of grammatical feature-values in the system. That is why an algorithm based on this information would be very inefficient. Our algorithm is based on the *basic constraints* corresponding to feature-values which are only about 150 in the system. Those constraints are already defined in the system, because the kernel program is producing the result via their unification.

The FSs used in our application are of the classic type. That means that they are not sorted and we do not allow negation and disjunction inside the FSs. The generalization of the algorithm in order to use disjunctive FSs is not a serious problem, but if we want to use negation inside the FSs, the algorithm should be generally revised. Using negation, we loose the nice classic semantic about FS's — the interpretation that a FS represents partial knowledge. For a comprehensive study of FS semantics see [2].

Some notion preliminaries: when we write a feature-value pair, in fact we mean a pair of feature path and value, where the feature path is a list of features. In our application we are interested only in features carrying grammatical information (other features like 'stem', 'wordform', etc. could not be classified). That is why we note only the last feature in the path and the value and call this a feature-value pair. In the application there are about 60 features (feature paths) carrying grammatical information. All feature-value pairs are about 150. This is a rather small number, hence the algorithm based on this information will be comparable effective. Bellow we present our algorithm which satisfies all requirements mentioned above.

Algorithm 1. *The Query Building Algorithm.*

- Step 1. Set the initial query FS to the empty FS.
Set all Feature-value pairs to 'acceptable'.
- Step 2. Wait the user to select a ('acceptable') feature-value pair.
- Step 3. Unify the query FS with the basic constraint corresponding to the selected feature-value pair and mark it as 'selected'.
- Step 4. For every 'acceptable' basic constraint
check the unifiability with the query
if the basic constraint does not satisfy the query,
then set it to 'unacceptable'.
- Step 5. For every feature check
if exactly one value for this feature is acceptable, then
select this feature-value and go to Step 3.
- Step 6. If there are no more acceptable pairs or the user has finished,
then go to Step 7, else go to Step 2.
- Step 7. Call the kernel program with the query as input.

This algorithm is almost self-explaining. In Step 1 the initialization is made. Step 2 and Step 3 build the query by unification of the basic constraint corresponding to the selected feature-value pair. Step 4 checks all other basic constraints for consistence with the query and switches all feature-values which are inconsistent with the query to an 'unacceptable' state. Step 5 is responsible for the automatic deduction of feature-values. Step 6 and Step 7 close the loop and invoke the kernel program respectively. We omit a detailed proof about the correctness of this algorithm, which in our point of view, is rather obvious. Maybe the only non trivial problem is the termination. The next lemma is concerned about that.

Lemma 1. *Algorithm 1 is always terminating.*

Proof. The two loops are always passing through Step 3, where a feature-value pair is set to 'selected'. There are only a finite number of pairs. Hence, after a finite number of iterations there will be no more 'acceptable' pairs, which guarantees the termination of the algorithm.

It has to be noted also that in the Algorithm is essentially used the 'Closed World Assumption'. In Step 5 we assume that there are no other values possible for the feature. In our application, where the general grammatical knowledge is fixed, this assumption is generally true. In fact, exactly this step is responsible for the automatic feature-value selection, when the value of a feature can be deduced from the information already specified by the query. Without the 'Closed World Assumption' we could not deduce anything new in our system.

4. IMPLEMENTATION

We have created a prototype version in Sicstus Prolog with Tcl/Tk on Windows environment, which is able to extract about 20 FS per second. The system is

supplied with a very friendly user interface. The query is created using the Windows Graphical User Interface. By clicking on a feature in the list box, another list box is displayed, where the acceptable values for this feature are listed. By selecting a value, the unification with the corresponding constraint is invoked. If some feature or value is disabled, then the corresponding entry in the list box will be switched 'gray' (unacceptable). There are several options for the output format. The user can choose between the output of the whole FSs or only the values of some features (e.g. wordform, stem, etc.).

A more faster version will be created using C/C++ language soon. This implementation will provide a retrieval speed of about 200 FS per second. There will be no other differences between the C/C++ and the Prolog version. We hope that this system will be widely used for Bulgarian language education and research purposes.

Also, there is a World Wide Web version planned. The idea is to specify the query using the form options in HTML. Then the query will be passed to the knowledge retrieval system using a CGI-script. In this way the resources will be accessible through INTERNET.

5. CONCLUSION

The most interesting part of the algorithm, in our opinion, is the untraditional deduction procedure. It is clear that the application is very simple. That is why in fact the loop Step 3 – Step 5 will not make new changes to the acceptability of feature-value pairs. We think that this deduction procedure could be classified as a new approach to certain problems. It leads to a very elegant and effective algorithm for deduction in domains with complex classifications. In current version only classic feature structure interpretation and unification in empty theory are applied. At the moment we are working on the generalization of this approach to allow more powerful constraint based technics.

ACKNOWLEDGEMENTS. The author of this paper wants to thank all members of the team of Linguistic Modelling Laboratory for the support and fruitful atmosphere, without which this work could not exist.

REFERENCES

1. Briscoe, T., A. Copestake, V. de Paiva (eds.). Default Inheritance in Unification Based Approaches to the Lexicon. *ESSLLI'92 readings*, 1992.
2. King, P. J. A logical formalism for head-driven phrase structure grammar. Doctoral dissertation. Manchester University, Manchester, England, 1989.
3. Paskaleva, E., K. Simov, M. Damova, M. Slavcheva. The long journey from the core to the real size of a large LDB. In: *Acquisition of Lexical Knowledge from Text*, Boguraev, Pustejowski (eds.), Columbus, Ohio, 1993, 161–169.

4. Pustejovsky J. *The Generative Lexicon*. MIT Press, 1995.
5. Shieber, S. *An Introduction to Unification-Based Approaches to Grammar*. *CSLI Lecture Notes*, 4, 1986.

Received on September 12, 1996

Linguistic Modelling Laboratory
Laboratory for Parallel and Distributed Processing
Bulgarian Academy of Sciences
E-mail: `stoyan@lml.acad.bg`