

A DETERMINISTIC ANALYZER
FOR THE INTERPRETATION OF NATURAL LANGUAGE COMMANDS

Leonardo LESMO, Daniela MAGNANI, Piero TORASSO

Istituto di Scienze dell'Informazione - University di Torino
C.so Massimo D'Aseglia, 42 - 10125 TORINO - ITALY

Abstract

This paper describes a system which translates a query in Italian language into a representation which can be immediately interpreted as a sequence of algebraic operations on a relational data base. The use of a lookahead buffer allows the system to operate deterministically. Different knowledge sources are used to cope with semantics (associated with the lexicon) and syntax (represented as pattern-action rules). These knowledge sources cooperate during the query translation so that independent translation steps and intermediate representations of the command are avoided. Therefore the term "determinism" is used to mean that all the structures built during the process concur to build the final command representation.

Introduction

One of the most popular new ideas in the field of natural language processing is the "determinism" introduced by Marcus /1/. The system presented in this paper accepts Marcus' proposals for the syntactic analysis of the text, but extends that analysis to a real translation of an input query in Italian language by means of a more semantical structure of the actions associated to the pattern-action syntactic rules. These actions have the purpose of building representations of the constituents with which the syntactic rule is associated, in a form which allows the system to perform semantic checks by using the semantic information available from the lexicon.

One of the main features of the system is a cooperation between declarative (lexicon and pattern part of the syntactic rules) and procedural (procedures to perform semantic checks and the action part of the syntactic rules) knowledge sources such that the control flow is automatically transferred from one knowledge source to another without the intervention of a central controller.

As a benchmark, the system has been specialized to act as an interface to a medical consultation program for the analysis of the liver functional assessment /2/.

Command representation

The interpretation process translates the input command into a set of frame instantiations linked together. By "frame" we mean a collection of slots, each of which has a name, may contain a value and may be associated with a procedure which checks the semantic correctness of the value inserted into the slot by the grammar rules.

The system contains a set of prototype frames (currently 4: ACTIONFR, COMMANDFR, RESTRICTIONFR, CONCEPTFR) which are instantiated under request of the grammar rules. Notice that the checks associated with the slot of a given frame F may operate on data not necessarily stored in F, but in other frames needed to F. The check procedure associated with a frame slot is executed when the slot is filled: it may happen that, at that time, some of the slots (in the same frame or not) needed to perform the checks are still empty. In this case the procedure is suspended and its suspension point is stored in a list together with information which allows the procedures associated with slot filling and frame linking to resume its execution as soon as possible.

Lexicon

The lexicon contains noise words, non-content words (e.g. prepositions) and content words. The latter are used to build the final representation of the command and the information associated with them concerns their syntactic category and (for some of them) semantic information stored as links to other content words /3/, so that the lexicon may be considered as a network which allows the system to verify the semantic correctness of the command. A spelling correction is performed to detect and correct possible typing errors.

A more detailed description of the lexical knowledge source is reported in /4/.

Syntax

The syntactic knowledge is represented by means of Syntactic Rule Packets (SRP) each of which contains a set of pattern-action rules which specify the different forms a particular constituent may assume. The pattern part of the rules allows the syntactic

tem to deterministically identify the rule of the Packet which has to be applied to analyze the input. The pattern consists in a sequence of syntactic categories or constituents which have to be matched on the lookahead buffer, which is filled under request of the syntactic knowledge (CAT statements; see below) with the exact portion of input string needed to discriminate among the rules of the Packet. For the sake of efficiency, the set of patterns concerning the same SRP have been represented by means of a discriminant net which guides the filling of the buffer.

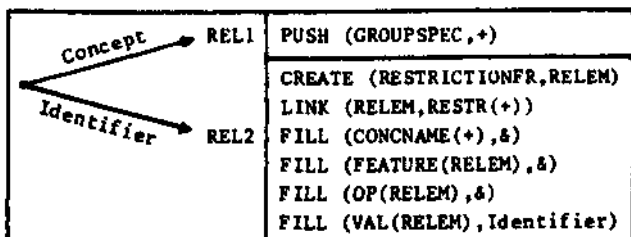
The action part of the rules is a procedure involving imperative statements:

- FILL. It fills a slot of a frame and triggers the checks associated with that slot.
- CREATE. It builds a new instantiation of a frame.
- LINK. It connects a frame instantiation to a slot of another one.
- CAT. It checks the category of the current word of the lookahead buffer. If all the elements in the buffer have already been scanned, a new word is transferred into it from the input string. Incidentally, the categories used in the grammar (and the constituents as well) are not the classical categories (noun, adjective, etc.) but are more semantically biased (concept, property, etc.).
- PUSH. It involves the triggering of the SRP associated with the embedded constituent that has to be analyzed and the transfer of control to the SRP handler.
- DEL. It deletes a word from the buffer.

When a FILL or LINK operations is executed, the list of suspended checks is scanned to decide whether some of them may be resumed. After the completion of the action part of a rule, the control is returned to the action which activated the SRP containing the rule whose action has been completed.

As an example of Syntactic Rule Packet, consider the constituent ELEMENT: it refers to a group of objects defined by a concept name (e.g. Patient) followed by the conditions selecting the individuals that should belong to the group, or, in case the group is restricted to a single individual, by its name (Identifier). The structure of the SRP is as follows:

ELEMENT

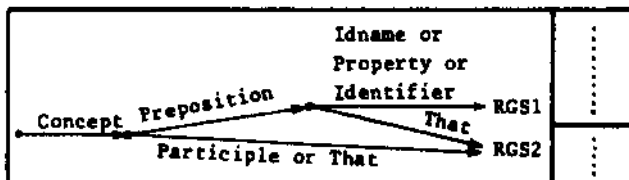


The symbol + refers to the current frame instantiation (in this case a copy of CONCEPTFR), whereas the symbol & means that the filler of the slot is

not explicitly present in the input sentence; in this example all the slots containing & will be filled by the semantic check procedure associated with the slot VAL, by inspecting the semantic information stored in the lexical entry corresponding to the current identifier. As regards RELEM, it is a variable which is filled by the CREATE operation with the actual name of the newly created instantiation.

It may be observed that, in this case, the discriminant net is very simple because only one word has to be examined in order to determine which rule has to be applied. In the GROUPSPEC discriminant net it is apparent that two or three lookahead words are needed to deterministically decide what rule has to be applied; for this reason the lookahead buffer size is not fixed, even if our grammar never requires more than four cells.

GROUPSPEC



Notice that the travelling on the discriminant net does not imply that an input word is consumed (for example, the first word in the PUSHED GROUPSPEC is the same Concept tested by ELEMENT).

Examples and comments

One of the most important features of the system is its ability to obtain the same representation for sentences which have the same meaning but completely different surface structures. Consider for example the sentences:

- S1) Dammi il valore della bilirubina totale per i pazienti con ittero (Give me the value of total bilirubin for the patients with jaundice).
- S2) Per i pazienti che hanno ittero dimmi quanto vale la bilirubina totale (For the patients who have jaundice tell me what the total bilirubin is).
- S3) Quanto hanno di bilirubina totale i pazienti con ittero presente? (How much total bilirubin have the patients with jaundice present?).

For all of them the final representation is the same and is reported in fig.1.

It should be pointed out that information which is implicit in the input command has been made explicit (e.g. in RESTR3 the fact that "jaundice" is the "name" of the symptom and in RESTR2 that the "value" of the jaundice used for patient selection should be "present", even if in the sentences S1 and S2 it is understood).

The resulting representation may seem rather awkward. For example, it may not be apparent why RESTR2 is linked to an action (named SHOW) instead

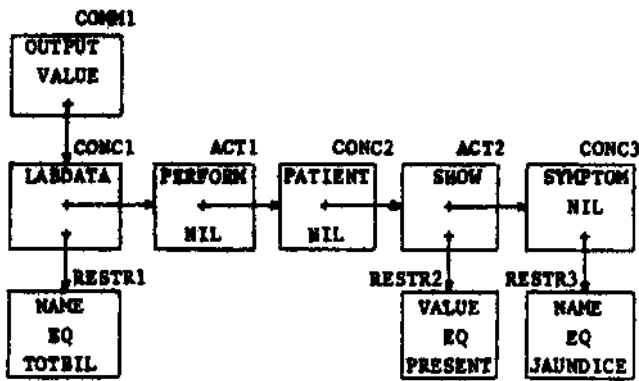


Fig.1 - Final representation of S1, S2 and S3

of the concept it qualifies (jaundice). The advantage of this representation is the one-to-one mapping between the pointers which link the frames and the algebraic operators of the relational approach to database construction. For example, each pointer to a restriction implies a SELECT operation on the relation to which the concept (or the action) is associated, each pointer connecting an action and a concept implies a JOIN operation on the two involved relations and so on.

It is worth noticing the difference between the representation of the sentence S1 (fig.1) and the representation (fig.2) of the sentence:

S4) Dammi il valore della bilirubina totale per i pazienti di sesso maschile (Give me the value of total bilirubin for the patients of male sex). The difference is due to the fact that "sex" is a property of a patient and for this reason it is stored in the data base relation PATIENT, whereas "jaundice" is an identifier of the concept "symptom" (i.e. a key of the data base relation SYMPTOM), so that the value of jaundice for a given patient is stored in the connection relation named SHOW.

In fig.1 and 2 the frame structure was heavily simplified for the sake of readability; let us consider as an example the actual structure of RESTRICTIONFR:

RESTRICTIONFR (FEATURE, OP, VAL, FEATURE!, OP!, VAL!, BACK, RESTRCHAIN)

The slots composing the frame are of three kinds: External, Internal, Link; the External slots (e.g. FEATURE, OP, VAL) contain the actual lexical entries appearing in the input sentence; the Internal ones (e.g. FEATURE!, OP!, VAL!) concern the corresponding internal representation, which will be used to access the data stored in the data base relations; the Link slots (e.g. RACK, RESTRCHAIN) contain the pointers to other instantiations.

In RESTRICTIONFR the slot FEATURE contains the name of the attribute used to perform the selection, whereas the slots OP and VAL contain the operator and the linguistic or numeric value involved in the

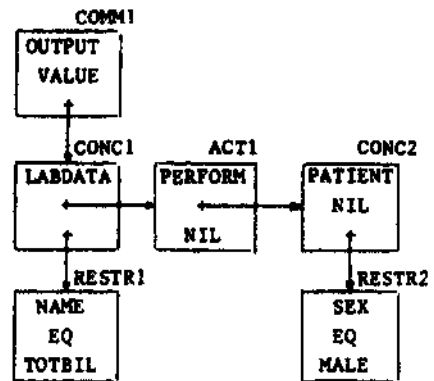


Fig.2 - Final representation of the sentence S4

SELECT operation. When the string "di sesso maschile" in the sentence SA is analysed, the slots FEATURE and VAL are filled with "sesso" and "maschile" respectively, whereas the slot OP is filled with & (undefined). The semantic checks associated with FEATURE and VAL provide not only the filling of FEATURE! and VAL! with the internal representation of "sesso" and "maschile" (i.e. SEX and MALE), but also the filling of OP! with EQ (the choice of EQ depends on the information associated with "maschile") thus leading to the situation of RESTR2 shown in fig.2.

The slot BACK contains a backward pointer to the CONCEPTFR (or ACTIONFR) instantiation for which the RESTRICTIONFR instantiation acts as a condition. Finally, the slot RESTRCHAIN may contain a pointer to another instantiation of RESTRICTIONFR which acts as a further condition for the selection (i.e. the two - or more - conditions are ANDed).

The system has been implemented in LISP on a DEC-10; the lexical and syntactic processes have been tested fully, whereas the procedures associated with the semantic checks are currently under debug.

References

- /1/ Marcus, M.P.: A Theory of Syntactic Recognition for Natural Language. MIT Press, Cambridge, Mass., 1980.
- /2/ Lesmo, L.; Saitta, L.; Torasso, P.: Computer-Aided Evaluation of Liver Functional Assessment, Proc. 4th Symposium on Computer Applications in Medical Care, Washington, D.C., 1980, pp.181r189.
- /3/ Heidorn, G.E.: Augmented Phrase Structure Grammars, in Schank, R.C. & Nash-Webber, B. (eds.): Proc. Theoretical Issues in Natural Language Processing, Cambridge, Mass., 1975, pp.1-5.
- /4/ Lesmo, L.; Magnani, D.; Torasso, P.: Lexical and Pragmatic Knowledge for Natural Language Analysis, 1981 Int. Conf. on Cybernetics and Society, Atlanta, Georgia, Oct.1981.