

SEOUL94.DOC/940922

*Representation of finite utterances  
and the automatic parsing of texts*

Maurice Gross

Laboratoire d'Automatique Documentaire et Linguistique

University Paris 7

1. Models of grammar
2. Finite-state graphs
3. Finite constraints
4. Inserts and non-finite constraints
  - 4.1. Adverbial inserts
  - 4.2. Sentential determiners
5. Parsing

References

- |         |                       |
|---------|-----------------------|
| Annex 1 | Table 4               |
| Annex 2 | Automaton for table 4 |

We describe the use of finite state automata for the description of natural languages. We demonstrate the use of this model of grammar through linguistically varied examples, from time adverbials and sentential determiners to elementary sentences of a lexicon-grammar.

## 1. Models of grammar

N. Chomsky 1955-6 gave a discussion of formal models of grammars and concluded that neither finite-state grammars nor phrase structure grammars (context-free or context sensitive) were adequate to describe natural languages. N. Chomsky's mathematical 'proof' proceeds by showing that the description of certain syntactic phenomena requires formal devices that are beyond the power of those he criticized. Chomsky used examples that he singled out for the purpose of the discussion. However, a careful analysis of these examples indicates that they can well be considered as exceptional linguistic structures, hence they could be treated independantly of the bulk of syntactic phenomena.

To show the inadequacy of finite-state grammars, Chomsky invokes the phenomenon of self-embedding, that is, the relative clause embedding of the examples:

*The cake was stale*

*The cake (that the rat ate) was stale*

*The cake (that the rat (that the cat killed) ate) was stale*

It is true that the rule that embeds relative clauses whose pronoun is an object is recursive. But it is also clear that with respect to understanding, embedding has to be limited to depth 3 at most. What is more interesting is that this recursive phenomenon seems unique: outside of this particular type of relative clause embedding, it is hard to find another clear-cut example. On the contrary, we mostly observe finite-state structures such as:

*The cat killed the rat that ate the cake that was stale*

We can set aside the self-embedding mechanism, either by considering it as an exception to be treated by a special device or by limiting arbitrarily the depth of embedding.

To show that context-free grammars are inadequate, Chomsky used the same type of argument, observing that coordinations involving the adverb *respectively* cannot be correctly described by phrase-structure grammars. But again, when one investigates the structures of English (and of other well-described languages), one finds practically no other phenomena of this type, except for the construction:

*Bob will work, leave or stay according to whether Jo will stay, leave or sing*

where the verbs of each half are paired in a way that generates an unbounded number of 'crossing' constraints, as shown by the paraphrase:

*If Bob works, Jo will stay, if he leaves, she will leave, if he stays, she will sing*

As a consequence, the transformational model remains the only adequate candidate for the description of these phenomena. We won't discuss how this conclusion is logically entailed from such examples (M. Gross 1972), we will just insist on the fact that syntactic phenomena present a large variety and that only very few of them, those N. Chomsky pointed out, escape the range of application of the weakest models. Along the same line of discussion, G. Harman 1963 has provided convincing arguments running against Chomsky's conclusion.

## 2. Finite state graphs

Finite state automata are by now a familiar object in computational linguistics. Among the well-known uses of this model is the ATN system (Augmented Transition Network, W.A. Woods 1970) and its variants, used for specific applications. From a theoretical point of view, the variety of notational variants can be reduced to a minimal set of algebraic structures (e.g. D. Perrin 1994).

Linguistic phenomena are represented in a natural way by the formalism of graphs. Other formalisms such as triples (*State, symbol, State*), rewriting rules:  $S_i \rightarrow a_j S_k$ , regular expressions or algebraic systems do not reflect as directly as graphs the word sequences to be described.

We illustrate the use of graphs<sup>1</sup> by two examples of a different formal nature:

Example 1: Adverbial expressions that correspond to rounded dates such as in the example:

*(It happened) in the early twenties*

---

FIGURE 1

---

In this example, the family of adverbs corresponds exactly to all sequences that can be read from the initial (left-most) state to the final (right-most) state. The number of phrases is strictly finite (equal to 244 here)<sup>2</sup>.

Example 2: Double conjunctions such as:

*On the one hand, Bob is wrong, but on the other, one should listen to him*

---

---

1. M. Silberztein 1993 has design a graphic tool for the construction of such finite-state graphs FSGRAPH and of associated parsers.

2. To be complete, one should append to this graph productive forms such as *in the 1970s*.

FIGURE 2

---

In figure 2, we have represented a set of adverbial conjunctions *CONJ* that build conjoined structures of two sentences  $S_1$  and  $S_2$ . The conjunction has two parts (at least)<sup>3</sup>, hence the complex sentence forms that we represented:

$$CONJ_1 S_1 CONJ_2 S_2$$

Moreover, the part  $CONJ_1$  has adverbial mobility in  $S_1$  and so has  $CONJ_2$  in  $S_2$ :

*Bob, on the one hand, is wrong, but we should, on the other, listen to him*  
=  
*On the one hand, Bob is wrong but one should listen to him on the other.*

In figure 2, we did not attempt to represent the exact sentence structures. The graph simply indicates that both parts  $CONJ_1$  and  $CONJ_2$  can be separated by an arbitrary number of words, a feature represented by a loop (or cycle) on the variable *MOT* i.e. *WORD*). Moreover, we gave no indication in the graph about adverbial mobility, the reason being that the formalism of automata is not well adapted to the description of sentences that differ by a permutation of some of their parts.

The main difference between graphs 1 and 2 is that graph 1 is strictly finite. Such finite graphs are called DAGs (directed acyclic graphs), in contrast, graph 2 contains one cycle. Graphs without cycles (DAGs) can be seen as a natural extension of a text. A text can be considered as a flat graph, read from left to right, as in figure 3:

---

FIGURE 3

---

A non trivial DAG is read in the same way, but contains possible options in the reading process: at each branching point, several texts are possible. This remark<sup>4</sup> is used to represent ambiguities and variants of texts<sup>4</sup>.

The difference between strictly finite and cyclic structures can be used to classify syntactic phenomena. For example, a good deal of the structure of noun phrases is strictly finite. Consider the general form:

(1) *Prep Det N*

---

3. There are examples with unbounded number of parts:

*Firstly S<sub>1</sub>, secondly S<sub>2</sub>, thirdly S<sub>3</sub>, etc.*

4. E. Roche 1993 has represented in this way the ambiguities of texts to be parsed automatically.

where the preposition *Prep* and the determiner *Det* can be 'zero'. This oversimplified global structure corresponds to a large variety of complex forms:

- *Prep* can be a complex form such as: *on behalf of*,
- *Det* can also be a complex determiner, such as *a large number, forty of fifty*.

Hence, (1) can correspond to the phrase:

*on behalf of a large number of players*

Moreover, the noun can be preceded by adjectives, themselves modified by adverbs:

*on behalf of a large number of very well motivated players*

In the absence of a detailed analysis of the sequence of modifiers that can precede a noun, a loose way of representing the structure is by means of the cyclic graph of figure 4.

---

FIGURE 4

---

However, more refined studies of the compounding process of modifiers (e.g. Z.S. Harris 1976) show that the sequence of pre-nominal modifiers is strictly finite, this result eliminates all loops in the graph of figure 4. Instead, strictly finite graphs have to be built, they are much more complex, but much more precise.

### Remarks

1. In post-nominal positions, conjoined sequences of modifiers are common, less so in pre-nominal positions. Since, constraints on conjoined units are not describable by linguistic tools, one must use loops to represent them.

2. Inserts may occur in structure (1), such as in the following form:

*on behalf, we think, of forty of fifty players*

The insert *we think* is of a sentential nature, hence its length is unbounded, for example it could be replaced by the longer insert: *we are absolutely sure of this fact*. Longer inserts can be stylistically awkward, but they are still grammatical. It is clear that such inserts, do not belong to the structure of noun phrases. We will discuss them in a general way below in 4.

3. Finite constraints

The original model of transformational grammar proposed by Z.S. Harris 1952 and the first model of generative grammar (N. Chomsky 1955) both make a clear separation between two sentence types:

- elementary, simple or kernel sentences which constitute generators, for
- complex sentences.

In these models, unary transformations affect the elementary structures and binary transformations combine simple structures into complex ones. This natural schema is also present in traditional textbooks, but has disappeared from the later models of generative grammar.

The study of elementary sentences can be performed in a way totally independent of the complex structures. It amounts to determining the argument structure of sentences and the possible modifications of basic argument structures by unary transformations. Descriptions of elementary structures have been systematically performed for several languages within the theory of lexicon-grammar. One important empirical result then obtained is that the maximum number of arguments of verbs is three, as for example in a sentence such as:

*Bob gave a ring to Jo*

Forms with more arguments can be observed, but they are quite restricted and may be subject to reanalysis with fewer arguments:

- there can be true exceptions such as the French idiomatic form with five arguments (all obligatory):

*(Luc)<sub>0</sub> a tourné (sa langue)<sub>1</sub> (sept fois)<sub>2</sub> (dans sa bouche)<sub>3</sub> (avant de parler)<sub>4</sub>*

- there are remaining theoretical difficulties in separating the essential arguments of a given verb from its circumstantial ones. The latter ones are brought, in principle, into the simple sentence through binary transformations of the type:

*Bob gave a ring to Jo yesterday*  
=  
*Bob gave a ring to Jo, this happened yesterday*

But in the following sentences with four arguments, the argument status of *for ten dollars* and of *for this ring* is not so clear:

*Bob paid ten dollars to Max for this ring*  
*Bob bought this ring from Max for ten dollars*

Both *for*-complements may seem circumstantial, however their *NP* part may occur in a direct object position which is definitely an argument position of the verb. In the same way, in the sentence:

*Bob wasted ten hours on this report*

*ten hours* is a direct object but is transformationally related to the duration complement of *write* in the complex sentence:

*Bob wasted ten hours writing this report*

- certain unary transformations may change the number of arguments of a sentence. The Passive transformation leaves invariant the number of arguments:

=  $(\text{Bob})_0 \text{ attacked } (\text{the fort})_1$   
 =  $(\text{The fort})_1 \text{ was attacked by } (\text{Bob})_0$

but the nominalization:

=  $(\text{Bob})_0 \text{ attacked } (\text{the fort})_1$   
 =  $(\text{Bob})_0 (\text{launched} + \text{made}) (\text{an attack})_1 \text{ against } (\text{the fort})_2$

increases by one the number of arguments. However, the main verbs are of a very different nature in such paired sentences: *to attack* is a distributional verb which constrains semantically its subject and object, whereas *to launch* is a support verb, namely a grammatical auxiliary with limited semantic role. Nominalizations with support verbs do not always increase by one the number of arguments, in many cases they modify the role of arguments. For example, in the relation with support verb *to put*:

=  $(\text{Bob})_0 \text{ coated } (\text{the cake})_1 \text{ with } (\text{chocolate})_2$   
 =  $(\text{Bob})_0 \text{ put } (\text{a coating of chocolate})_1 \text{ on } (\text{the cake})_2$

*coating*, the nominal form of the verb, has for noun complement the instrument complement of the verb, that is the noun *chocolate*. From a syntactic point of view *coating of chocolate* is a single noun phrase, hence it should be counted as a single argument; consequently, both the nominal and the verbal sentences have three arguments. In the process of nominalization, an argument of a verb has become a modifier of a noun, which could be seen as having a non essential role in a sentence. Such changes in the syntactic properties of the various arguments show the complexity of the correspondance between syntactic structures and argument structures that are closer to semantic interpretation.

After a systematic study of the French lexicon, the set of kernel sentence forms appears to be the following<sup>5</sup>:

$N_0 V$	intransitive forms
$N_0 V \text{ Prep } N_1$	2 arguments, <i>Prep</i> can be 'zero'.
$N_0 V \text{ Prep } N_1 \text{ Prep } N_2$	3 arguments, <i>Prep</i> can be 'zero'

and marginally:

---

5. In English and other languages, the structures and even their numerical proportions in the lexicon do not seem to be essentially different.

$N_0 V(\text{Prep } N_i)^n$ , with  $n$  no larger than 4.

Such a set of structures is thus strictly finite and is described in a very natural way<sup>6</sup> by the finite automaton of figure 5.

---

FIGURE 5

---

The same form of automaton can be used for a different purpose. Consider the sentence with three arguments:

*(Bob)<sub>0</sub> talked to (Jo)<sub>1</sub> about (the ring)<sub>2</sub>*

the complement arguments are not obligatory, and the following forms are also accepted as sentences:

*(Bob)<sub>0</sub> talked to (Jo)<sub>1</sub>*  
*(Bob)<sub>0</sub> talked about (the ring)<sub>2</sub>*  
*(Bob)<sub>0</sub> talked*

The automaton of figure 5 can represent this set of four sentences. However, this set is only valid for *to talk*, we need a different automaton for *to mention*, which has the different paradigm:

*Bob mentioned the ring to Jo*  
*Bob mentioned the ring*  
*\*Bob mentioned to Jo*  
*\*Bob mentioned*

As a consequence, to represent the optional or obligatory status of arguments of verbs, the general automaton of figure 5 must be lexically specified: the verb and the prepositions must be made explicit and the nature of the arguments clearly specified, which is the case in the matrix representations of the lexicon-grammar (M. Gross 1975). This method of representation can be extended to other structures, for example to the structures obtained through transformations. This possibility directly derives from the nature of lexicon-grammar. Let us recall the principle of the matrix representations (annex 1). A row of a matrix is an entry, for example a distributional verb. It is important at this stage that the various meanings of the entry word, that is the word form appearing in editorial dictionaries, have been clearly separated<sup>7</sup>. The argument structure of verbs has been used to establish a classification. For

---

6. It should be noted that the graph makes explicit the structural invariance of the sequence  $N_0 V$ , common to all sentences. This observation should be opposed to the insistence of linguists to consider the  $VP$  structure (verb phrases) as a universal invariant.

7. For example figurative and proper meanings of a word often constitute separate entries, since in general for each meaning the set of syntactic properties differs (J.-P. Boons 1971).

12.000 French verbs we have defined about 50 classes (C. Leclère 1991). Each class is represented by a specific matrix. The rows of a matrix correspond to the entries (e.g. the verbs). Columns are sentence forms, for example:

- the Passive form:  $N_1$  *be V-ed by*  $N_0$
- the Impersonal form: *it V*  $N_0$  *Prep*  $N_1$

Hence, a transformation is a pair (unordered) of columns. The Extraposition transformation can then be written:

$$N_0 V Prep N_1 = it V N_0 Prep N_1$$

*That Bob would fail occurred to Jo*

= *It occurred to Jo that Bob would fail*

At the intersection of a row (entry) and a column (sentence form), we place a '+' sign if the entry is compatible with the sentence form, a '-' sign otherwise. In this way, we associate to a given entry a set of compatible sentence structures. In exactly the same way we associated above the substructures of the verbs *to talk* and *to mention* to finite automata, we can construct all the automata corresponding to all the entries of the lexicon-grammar. E. Roche 1993. has effectively constructed such automata in a highly formalized way, to the point where the automata he built can be used in automatic syntactic analysis (annex 2).

#### 4. Inserts and non-finite constraints

If we attempt to match the basic structures described in the lexicon-grammar with sentences found in texts, many questions arise. One set of questions relates to complex sentences, answers to these questions lie in the detailed description of coordination and subordination, that is of binary transformations. Many questions are still open in this active area of research, in particular the role of the lexicon-grammar has to be determined (M. Mohri 1993, M. Piot 1991).

Another series of discrepancies between theoretical and observed forms is related to inserts of the type exemplified in 2.

##### 4.1. Adverbial inserts

Let us consider an elementary structure of a general type:

(1)  $N_0$  *Aux V Prep*  $N_1$  *Prep*  $N_2$  =:

*Bob has given a ring to Jo*

and any type of adverbial, namely *three days ago*, *generously*, *in a bar*, etc. Such adverbials may systematically occur at the juncture of the units of (1), that is next to any of the noun phrases or of the verbs. We mark these positions by a \$-sign in:

(2) \$  $N_0$  \$ *Aux V* \$ *Prep*  $N_1$  \$ *Prep*  $N_2$  \$ =:

*Three days ago, Bob has given a ring to Jo*

*Bob, three days ago, has given a ring to Jo*  
*Bob has, three days ago, given a ring to Jo*  
*Bob has given, three days ago, a ring to Jo*  
*Bob has given a ring, three days ago, to Jo*  
*Bob has given a ring to Jo, three days ago*

In general, Adverbial inserts are not permitted inside noun phrases. Some inserts are not allowed in all the \$-positions<sup>8</sup>.

Adverbials have unbounded length, as in:

*the day they had decided to go to the beach*  
*in the generous way his parents had always taught him*  
*in a bar where several extremely serious accidents had occurred*

as a consequence, a relation between two of the sentence units of (1) can hold at any distance. For example, matching the person-number of the subject with the person-number of *Aux* may require that one the preceding lengthy insert has been recognized in the substructure  $N_0 Adv Aux$ <sup>9</sup>.

Performative inserts such as *I think, God knows why, as I just told my sister*, are also allowed in the same positions (M. Gross 1990):

*God knows why, Bob has given a ring to Jo*  
*Bob, God knows why, has given a ring to Jo*  
*Bob has, God knows why, given a ring to Jo*  
 etc.

#### 4.2. Sentential determiners

Another syntactic process that can keep apart noun phrases from their verbs is an extension of the determiners of nouns. Common determiners such as articles (definite, indefinite), demonstrative and possessive provide a picture *Det N* of the noun phrase where a short *Det* can only be separated from its *N* by adjectives (cf \_2):

*Bob bought (the + a + this + my) car*  
*Bob bought (the + a + this + my) extremely nice and inexpensive car*

In the following examples of *Det* are of a different nature:

*Bob bought God knows exactly how many cars*  
*Bob bought I cannot tell you what brand of car*

8. The acceptability of Inserts may vary according to stylistic features. But all \$ positions are in principle grammatical. An exception is observed with *barely*:

*Bob barely reads*

*Bob reads barely*

\**Barely, bob reads*

9. Moreover, *Adv* may stand for more than one adverbial sequence.

The determiner sequence is sentential, and as such, it can be of any length. It is interesting to compare such determiners to the performative inserts, they are lexically related in the sense that it is the same types of main verbs that are found in both structures. But the structures are quite different, performative inserts can move freely between the phrases of the main structure, whereas the sentential determiner is fixed in the pre-nominal position *Det* of a noun phrase.

Another type of determiner generates sequences of unbounded lengths too. In principle, nominal determiners compound recursively:

*Bob bought a large number of books*

*Bob bought a large number of a certain kind of books*

However, very much as in the case of pre-nominal adjectives, the allowed combinations of nominal determiners are quite limited<sup>10</sup> and even if we set aside the stylistic problem of length, it is difficult to find interpretable examples with more than three levels. The sentence:

*Bob bought a subset of a collection of a certain kind of books*

is both logically correct and grammatically acceptable, but its set-theoretic relations which can be extended indefinitely do not translate into normal human discourse; the corresponding sentences belong to the language of set theory and are best phrased and interpreted by using the mathematical notations of the domain.

## 5. Parsing

The \$-positions of (2) in 4.1. introduce a difficulty in the analysis of (1). It is clear that if inserts could be recognized **first**, structure (1) would compare much more easily to the entry of *to give* in the lexicon-grammar. We advocate such a strategy of parsing, although it runs against the current attitude. Today, specialists are devising general processes as independently as possible of the specific grammatical features of the language to be parsed. Most parsers thus rely on a general model (usually some type of phrase-structure model) and algorithms that are applied (left-to-right, bottom-to top, etc.), are blind to the categorization of linguistic phenomena, even from the formal point of view we presented. For example, it is considered that phrase-structure parsing is general, powerful and efficient, because it treats in the same way finite and recursive constraints between words or phrases.

Our approach consists in using formal differences observed at the empirical level. For example, we saw in \_2 that sentence structures in languages that have fixed word-order can be modelled by finite-state automata in a very natural way. This is not the case for the structures with adverbial inserts we discussed in \_4.1. They are best described by means of a specific permutation device that acts on a finite-state representation. In other terms, we are making more specific the early transformational models:

---

10. Examples such as:

*\*Bob bought a certain quantify of a large amount of books*  
have to be blocked

- kernel sentences are described in terms of finite automata,
- kernel sentences are submitted to operations that transform the finite-state graphs into other finite-state graphs.

Transformations then appear to be highly specific, we have illustrated here this feature by examples as different as the adverbial permutation and the insertion of sentential determiners of nouns, the detailed grammar of many different languages provide many more examples supporting this view.

## References

Boons, Jean-Paul 1971. Métaphore et baisse de la redondance, in Syntaxe transformationnelle du français, *Langue française*, Paris: Larousse, pp.15-16.

Boons, Jean-Paul; Alain, Guillet; Christian, Leclère 1976. *La structure des phrases simples en français. I Constructions intransitives*, Geneva: Droz, 377p.

Boons, Jean-Paul; Alain, Guillet; Christian Leclère 1976b. *La structure des phrases simples en français, II Constructions transitives*, Paris: Rapport de recherches du LADL, N° 6, 85p., tables et index, 58p.

Chomsky, Noam 1956. Three models for the description of language, *IRE Transactions on Information Theory*, IT-2, pp.113-124.

Chomsky, Noam 1957. *Syntactic Structures*, The Hague: Mouton.

Gross, Maurice 1972. *Mathematical Methods in Linguistics*, Englewood Cliffs N.J.: Prentice Hall Inc., 159 p.

Gross, Maurice 1975. *Méthodes en syntaxe*, Paris: Hermann, 412 p.

Gross, Maurice 1990, Grammaire transformationnelle du français, III Syntaxe de l'adverbe, Paris: ASSTRIL, 670 p.

Guillet, Alain ; Christian Leclère ; Jean-Paul Boons 1992. *La structure des phrases simples en français. Verbes à complément direct et complément locatif*, Geneva: Droz, 445 p.

Harman, Gilbert H. 1963. Generative grammars without transformation rules, *Language* **33**, pp. 597-616.

Harris, Zellig 1952. Discourse Analysis, *Language* **28**, pp. 1-30.

Harris, Zellig 1976, *Notes du cours de syntaxe*, Paris: Le Seuil, 237 p.

Mohri, Mehryar 1993. Réduction des complétives à un nom et article défini générique, *Linguisticae Investigationes* XVII:1, Amsterdam-Philadelphia: J. Benjamins Pub. Co., pp. 83-98.

Leclère, Christian 1990. Organisation du lexique-grammaire des verbes français, in: Dictionnaires électroniques du français, *Langue française* N° 87, Paris: Larousse, pp. 112-122.

Perrin, Dominique 1994. Finite Automata, in *Handbook of Theoretical Computer Science*, Vol.B, Jan van Leuween ed., Amsterdam: Elsevier Science Publishers, pp. 1-57.

Piot, Mireille 1991. Quelques problèmes inédits de constructions avec des conjonctions 'conséquentielles', *Lingvisticae Investigationes* XV:2, Amsterdam-Philadelphia: J. Benjamins Pub. Co., pp. 285-304.

Roche, Emmanuel 1993. Une représentation par automate fini des textes et des propriétés transformationnelles des verbes, *Lingvisticae Investigationes* XVII:1, Amsterdam-Philadelphia: J. Benjamins Pub. Co., pp. 189-222.

Silberztein, Max 1993. *Dictionnaires électroniques et analyse automatique de textes : le système INTEX*. Paris: Masson, 233 p.

Woods, W.A. 1970. Transition network grammars for natural language, *CACM*, 13(10), pp. 591-606.

