

APPENDIX III

DECISION PROCEDURES FOR STRUCTURE IN NATURAL LANGUAGES*¹⁾

Y. BAR-HILLEL

Hebrew University, Jerusalem, Israel

The rules of formation of a logistic system are by definition¹⁾ such that the notion of formula, well-formed formula or sentence, determined by these rules, is effectively decidable. However, I am not convinced that the arguments brought forth by Church²⁾ to the effect that sentencehood has to be an effectively decidable notion for any system that may be used for communication purposes are conclusive. I therefore regard it to be a serious problem whether the syntactic structure of a natural language such as English can always be adequately described by a set of formation rules that guarantee the decidability of the notion of sentence or, for that matter, of any other syntactical structures such as phrases etc. Inasmuch as there exist good reasons for doubting whether the answer to this problem is affirmative, the prospects for fully-automatic, high-quality translation from one natural language into another natural language look dimmer than many workers in the field of machine translation would like to think. This is so since not even one necessary, though by no means sufficient, condition for this process, namely the mechanical determination of the syntactical structure of any given sentence in the source language, could possibly be completely fulfilled. Though applicability to machine translation is often in the back of my thinking on the description of the syntax of natural languages, I shall refer here no longer to this application, having dealt with it elsewhere at some length.³⁾

The seriousness of our problem has apparently not been sufficiently recognized so far because many linguists explicitly, and most if not all of them as well as most logicians implicitly, believed that the syntactical structure of natural languages is adequately describable by an immediate constituent model, or a phrase structure model according to the term recently introduced by Chomsky.⁴⁾ It is indeed true that if natural languages were adequately describable in terms of such a model, there would exist a decision procedure for structure, as I have shown in effect, though not with full rigor, in a paper published six years ago.⁵⁾

Before I proceed to present some arguments for the fact that the phrase structure model is not fully adequate, let me spend some time in presenting again,

*) A revised version of a talk given before the Colloque de Logique, Louvain, September, 1958. The present version was published in the Belgian journal *Logique et Analyse*, N.S., 2^e Année, No. 5, Janvier 1959. Since, however, this issue was sent to the printers only in the second week of February 1959, according to a communication from its editors, I decided not to wait for the arrival of the reprints and to reproduce it myself in the present form. So some minor discrepancies between the versions may be expected.

The reader will realize that the present paper overlaps with the one reproduced in Appendix II. After some hesitation, I decided nevertheless to include it here, as it is more elaborate in many points. A consolidation of my views on the theoretical aspects of MT is in preparation.

in briefer and, I hope, improved form, an informal outline of this proof. The basic idea behind the immediate constituent model is that every sentence can be regarded as a result of the operation of one continuous part of it upon the remainder such that those constituent parts which in general are not sentences themselves, but rather phrases, are themselves again the product of the operation of some continuous part upon the remainder, etc., until one arrives at the final constituents, say words or morphemes. To illustrate:

Young John slept soundly

would be regarded as the result of the operation of slept soundly upon young John; slept soundly in its turn would be considered the result of the operation of soundly upon slept and young John the result of the operation of young upon John. All this so far is nothing but reformulation in somewhat unfamiliar terms of the procedure well known from school days as parsing. As linguists put it, young John and slept soundly are the immediate constituents of the sentence under discussion, young and John the immediate constituents of the first immediate constituent of the sentence, slept and soundly the immediate constituents of the second immediate constituent. Hence altogether young, John, slept and soundly are the final constituents of the given sentence.

Another basic feature of the model is that all operator constituents must be contiguous with their argument constituents. Both these features are exemplified in our illustration, but this of course is by no means a proof that this model can be carried through all of language. On the contrary, linguists have realized that occasionally discontinuous constituents have to be taken into account, but they seem to have believed that these were exceptions which did not seriously affect the validity of the model with which they were used to work.

In most language systems invented by logicians, the two mentioned features were automatically incorporated into their respective rules of formation. The problems arising in connection with discontinuous expressions were, to my knowledge, never explicitly discussed by logicians.

According to the immediate constituent model, every word – and we shall for our purposes consider words to be the basic syntactical elements – of a natural language belongs to one or more syntactical category. Among these categories some will be pure argument categories, by which term I denote a category whose members always serve as arguments and never as operators, as well as operator categories whose members may operate upon other words though they may perhaps also be operated upon by other operator expressions. John, for instance, inasmuch as it belongs to the syntactic category of nominals, is always an argument and never an operator. Slept, inasmuch as it belongs to the category of intransitive verbals, may operate upon a nominal such as John to form the sentence John slept, but may also be operated upon by the adverbial soundly to form the intransitive verbal expression slept soundly. A word may belong to more than one category not only because it may be regarded as homonymous – as would be the case with regard to sleep, which clearly belongs to the category of nominals as well as to the category of intransitive verbals – but also because, for instance, many adverbials operate upon intransitive verbals as well as upon transitive verbals: soundly, for example in the sentence

Belgium soundly defeated the Netherlands

(in the last soccer game, of course), operates upon the transitive verbal defeated, forming the transitive verbal expression soundly defeated, and has therefore a different kind of argument as well as a different kind of value than has soundly when operating upon slept.

In order to exhibit the decision procedure for constituent structure let us denote, following Leśniewski and Ajdukiewicz, the category of nominals by 'n' and the category of declarative sentences by 's'. (Since I am engaged in presenting an outline only, I shall not go here into the very difficult question to what degree these two argument categories would have to be refined and expanded in order to get even the beginnings of a reasonably working model.) Operator categories will be denoted by symbols that will indicate both the categories of their arguments and the category of the resulting expression. In addition, since arguments may be positioned either at the immediate left or at the immediate right of their operator, these positions too will have to be indicated in the symbolism. Therefore, I shall, for instance, denote the category of slept by 'n\s' - read: n sub s - and the category of young by 'n/n' - read: n super n,⁶⁾ -- where the direction of the slash indicates in an obvious fashion whether the argument is to the left or to the right. And, for instance, qua sentence connective, will be assigned to the category s\s/s^{*} since in this function it is a word that out of a sentence to its immediate left and a sentence to its immediate right forms a sentence. Soundly will belong to the categories (n\s)\(n\s) - to be abbreviated in a self-explanatory way as n\s\(n\s) - and n\s/n/(n\s/n) - as well as to a few other categories.

Assume now that we have a complete category list of all English words, i.e. a list which gives all the syntactical categories to which every English word may belong. In order to arrive by a completely mechanical procedure at the constituent structure of any given English sentence, one would only have to copy from the category list the category symbols for all the words in this sentence, write them down in columns and go to work on them according to the following rule:

Replace a sequence of three symbols, having respectively the form α , $\alpha\beta/\gamma$ and γ with β . This rule comprises as limiting cases the following two subrules:

- (1) Replace the sequence of symbols of the form α and $\alpha\beta$ by β .
- (2) Replace the sequence of symbols of the form β/γ and γ by β .

Instead of going into a detailed but rather obvious description of the decision procedure let us illustrate through a somewhat more elaborate example. Assume that the word sequence to be tested for sentence-hood as well as for its constituent structure is

Paul thought that John slept soundly.

Assume further that copying from the category list yields the following result:

<u>Paul</u>	<u>thought</u>	<u>that</u>	<u>John</u>	<u>slept</u>	<u>soundly</u>
n	n	n	n	n\s	n\s\
	n\s	n\n			n\s/n/
	n\s/n	n\s			.
	n\s/s				.
	.				.
	.				.

(the three dots indicating that the complete list would probably contain further entries which shall, however, be here disregarded for the sake of simplification). The reader will do well to envisage contexts in which thought and that will belong to each of the given categories. He might as well try to find out to

*) [Added for the present version:] This notation may turn out to be too lax for certain purposes. A more strict notation is (s\s)/s. Similarly, the main rule given in the following paragraph should officially always be replaced by the two subrules given there. The explanation of a derivation, given below, is therefore somewhat inaccurate, and so are the examples. There should be no difficulty in introducing additional rigor, when required, in accordance with the procedure followed in Appendix II.

which categories thought would belong in such contexts as: John had thought of..., ...thought processes, and ...thought provoking... .

Now taking into account only the categories explicitly indicated we have twenty-four initial symbol sequences to which we will apply our rule. Starting for instance with

n n n n n\s n\s\\n\s

we see that subrule (1) can be applied for the fourth and fifth symbols yielding s. The resulting sequence is now

n n n s n\s\\n\s,

which obviously cannot be further operated upon. The same subrule operating upon the fifth and sixth symbols yields n\s, hence the sequence

n n n n n\s,

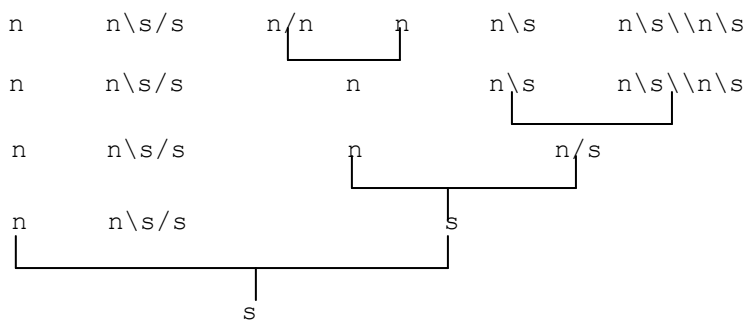
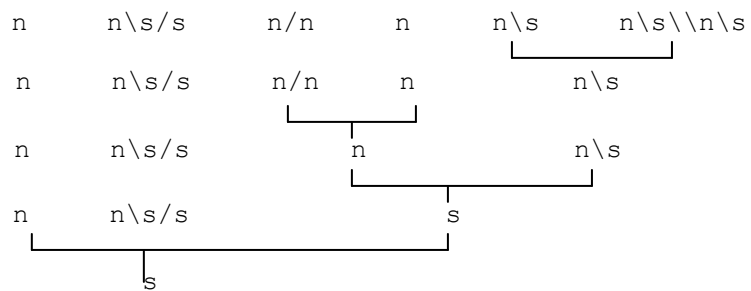
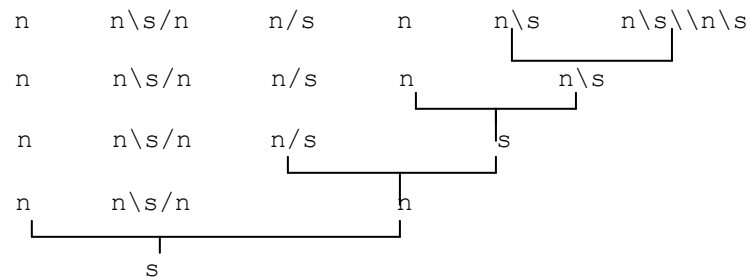
which has once more to be operated upon by the same subrule yielding

n n n s,

which cannot be processed any further.

Performing these operations upon all the twenty-four initial symbol sequences through all possible continuations, we would find that there exist exactly three derivations – as we shall call columns of symbol sequences each of which (with the exception of the first, of course) results from the preceding line by one application of the rule – whose final line, or exponent, consists of a single symbol which in both cases is 's'.

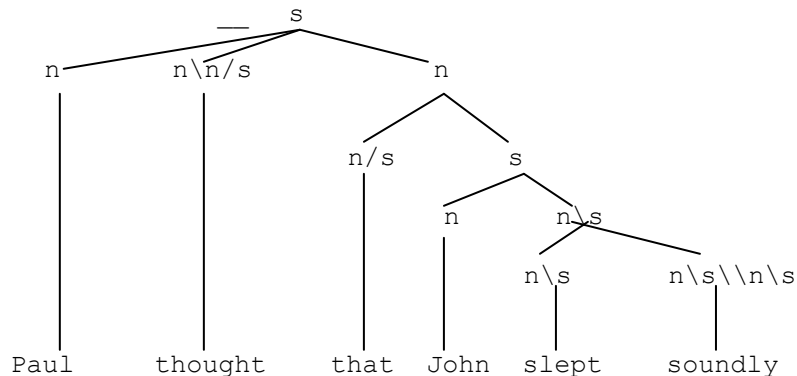
Here are the derivations:



The last two derivations being equivalent, in a rather obvious sense of the word, we have only two essentially different derivations before us, indicating, probably to the surprise of many readers – and to my own surprise some six years ago when I came across this situation simulating a machine processing of this illustration –, that the sentence under discussion is syntactically ambiguous or constructionally homonymous. The reader will do well to read out aloud this sentence according to its two essentially different constituent structures which in this case make the sentence also semantically ambiguous as such, though one constituent structure is much less likely to be used than the other.

I hope that this illustration is sufficient to show that under the essential and, as we shall see, highly problematic assumption that a complete and completely adequate category list is available, there exists indeed a wholly mechanical procedure to determine whether a given word sequence is a declarative sentence under one of its constituent structures as well as what all of its constituent structures are.

For certain purposes it is worthwhile to look upon our derivation procedure upside down, i.e. to deal with expansion rather than with derivation. The expansion corresponding to the first derivation exhibited above of our sample sentence would look like the following tree:



(Two derivations, by the way, are equivalent if they correspond to the same tree.)

How well then does the immediate constituent model work? Apparently quite well for relatively short sentences such as those discussed so far, but even there not too well. The number of categories to which the English words will have to be assigned to make the category list reasonably adequate will occasionally have to be rather large, and the categories themselves rather complex. In addition, it is quite clear that not only will one have to work with highly complex refinements of the categories mentioned so far in order to take care, for example, of the fact that John sleeps is a sentence but not John sleep, but that one will also have to refine the category of sentences and distinguish between declarative sentences, imperative sentences, yes-or-no question sentences, wh-question sentences, etc., these various types not being reducible to each other under our model. These refinements may result in such a piling up of category symbols assigned to the words occurring in a given sentence that the number of derivations would easily run into the trillions, hence be beyond the practical capacity of even the fastest electronic computers. For instance, if the average number of categories of the twenty words of a given English sentence is four, we will have up to 4^{20} initial lines and a still enormously higher number of derivations. This means, then, that the indicated method of mechanically resolving the syntactical structure of any given English sentence would certainly be impractical as such. However, were it the case that this is still a theoretically adequate method, one could think of

certain improvements which would reduce the required number of operations by many orders of magnitude. Unfortunately, however, the actual situation seems to be much worse. It is not only a matter of practicality, but it seems that the whole model is just not good enough. Already six years ago I was worried *by* sentences such as

John, unfortunately, slept soundly

which, so it appears at least, cannot be handled by a model incorporating the two above-mentioned basic features. Notice that there is no trouble with the slightly different and semantically, though perhaps not stylistically, equivalent sentence

Unfortunately, John slept soundly.

Assigning unfortunately to the category s/s, a wholly natural and intuitive assignment, we arrive at an adequate syntactical analysis. This assignment, however, clearly does not work for John, unfortunately, slept soundly, as the reader will easily verify for himself. It is of course possible that some other less natural category assignment to unfortunately, perhaps combined with some ingenious treatment of the commas (which so far have been completely disregarded in the immediate constituent model), would do the trick. It seems, however, unlikely that such an assignment could be made in a fashion which would not be almost entirely ad hoc. And this would not only be esthetically and methodologically repugnant but also, in all likelihood, have unpleasant repercussions inasmuch as word sequences which intuitively would not be regarded as grammatical sentences would have derivations with an exponent of s.

A similar situation, but even simpler since no commas are involved, arises with regard to the word sequence

He looked it up.

Regarding he and it as belonging to the categories n - leaving aside once more the clearly required refinements -, looked as belonging to the category n\s/n, as seems natural, it seems highly implausible that any category assignment of up which would not be woefully ad hoc would insure the sentencehood of the given word sequence. Assigning up, for instance, to the category s\s would obviously result in a derivation with an exponent s, but this unnatural saving of the phenomena would immediately retaliate with the unwanted imposition of sentencehood to such sequences as

He went home up.

(For further examples of the breakdown of the phrase structure model see Chomsky's Syntactic Structures,⁷⁾ to which I owe much of the present argument.)

Every English speaker, I presume, feels that in our sentence

He looked it up

looked and up belong somehow together. Indeed there is no trouble with such a sequence as

He looked up this argument,

as the reader will easily verify for himself, if only up is assigned in a completely intuitive fashion to the category n\s/n\n\s/n. This being so, assigning up to a different category, whatever it now may" be, in the sentence

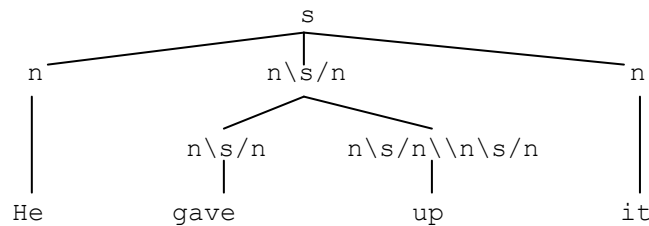
He looked it up

looks now even more artificial than before.

These simple facts indicate, though it cannot be said that they prove in the strong sense used in mathematics, that the immediate constituent model is not an adequate one as such, but has to be supplemented in one way or another.

Let me finish this discussion by presenting a very brief outline of one such supplementation method, referring the reader for a fuller discussion to Chomsky's mentioned book and other publications of his⁸⁾. The new model, called

the transformational model, assumes that sentences are generated not only by the procedure we called above expansion, but also in addition by so-called transformations. One such transformation, for instance, would transform the so-called terminal string of the following expansion



i.e. He gave up it, which is, of course, not an English sentence, into He gave it up by a certain obligatory transformation. This transformation rule which states in effect that in certain environments certain word sequences have to be turned around is clearly beyond the reach of an immediate constituent model. On the other hand, this way of looking at how the sentence He gave it up was generated has a rather natural appearance, and might well correspond, at least in spirit, to the way old-fashioned, traditional grammar has dealt with the situation.

Other transformations transform two terminal strings into one sentence. One of these, an optional one, would operate upon the sequence of the two terminal strings (which are in this special case sentences in their own right)

Paul thought it. John slept soundly.

and turn this sequence into the sentence

Paul thought that John slept soundly.

This very same transformation operates upon the sequence

Paul thought it. That John slept soundly.

and transforms it into

Paul thought that that John slept soundly.

Yet another transformation to the effect that under certain determined conditions that may be omitted would transform this last sentence into

Paul thought that John slept soundly.

This way of looking at the situation results now in a natural and adequate explanation of the constructional homonymy of the last sentence. We also realize, by the way, that transformations may operate upon the results of prior transformations .

Linguists, such as Harris, Chomsky, and their associates, who are at work at the development of this new kind of model⁹⁾ have already unveiled a large number of transformations amounting to many hundreds in English. It is, however, quite clear that the transformations introduced so far are not yet sufficient to account for all intuitively possible English sentences. It is at this state that the question mentioned at the beginning of this paper arises – whether there exists a decision procedure for structure in English, or in other natural languages for that matter, since it is unlikely that the natural languages should differ among themselves in this respect. Obviously the answer to our question will depend upon the exact nature of the transformations. Only when we will have a better and more extensive understanding of the kind of transformations at work, will we be in a position to fruitfully attack our problem. At this moment one could only speculate about this answer, and it is doubtful whether such speculations would be worthwhile. In any case, even the possibility that for a certain set of formation rules in English the notion of English sentence would not be a decidable (or general recursive) one seems exciting enough to warrant an increase in interest in our problem among mathematical logicians who by

training are in many respects in a better position to attack it than are linguists. Chomsky has already been able to show that there exist highly interesting connections between the theory of linguistic models and such theories as the theory of automata, recursive function theory (perhaps especially conspicuous in the form of the theory of algorithms) and the theory of Post canonical systems. This multiple relationship indicates that we have in all probability in the theory of language models an interesting new field in which cross-fertilization of mathematical logic and structural linguistics should lead to important results.

NOTES

- 1) See, e.g., A. Church, Introduction to mathematical logic, I, Princeton, 1956, p. 51. There exist, however, less demanding conceptions.
- 2) Ibid., p. 53.
- 3) In "Some linguistic obstacles to machine translation", forthcoming in the Proceedings of the Second International Congress of Cybernetics, held in Namur, September 1958.
- 4) See N. Chomsky, "Three models for the description of language", IRE Transactions on Information Theory, Vol. IT-2, No. 3 (1956) and Syntactic structures, 's-Gravenhage, 1957.
- 5) "A quasi-arithmetical notation for syntactic description", Language 29:47-58 (1953).
- 6) See K. Ajdukiewicz, "Die syntaktische Konnexitaet", Studia Philosophica 1:1-27 (1935-36); cf. A.A. Fraenkel and Y. Bar-Hillel, Foundations of set theory, Amsterdam, 1958, pp. 169-170.
- 7) In the paper mentioned in note 5, I used a less convenient symbolism. The present symbolism is due to J. Lambek, "The mathematics of sentence structure", American Mathematical Monthly 65:154 (1958).
- 8) See above, note 4.
- 9) Viz., to those mentioned above in note 4, as well as, for instance, to a forthcoming paper, "A transformational approach to syntax".
- 10) In addition to Chomsky's publications, see Z.S. Harris, "Cooccurrence and transformations in linguistic structure, Language 33:283-340 (1957) and the excellent review of Chomsky's Syntactic structures by R.L. Lees in Language 33:375-408 (1957).