

L. W. TOSH

**STRATIFICATIONAL GRAMMAR AND INTERLINGUAL MAPPING
FOR AUTOMATIC TRANSLATION**

Since 1959, the University of Texas Linguistics Research Center has been engaged in investigation into the possibilities of machine translation of languages. This has been a long-term project involving a number of funding sources. The effort is currently supported by contracts from the U.S. Air Force.

During this time, we have concentrated on the problems of syntactic description and the development of generalized parsing algorithms as tools for verifying descriptive research. Here, I propose to outline some aspects of the recent linguistic research. I will not deal with the technical aspects of the computer programming involved.

From the linguistic point of view there are two major areas of endeavor : one is pilot research of problematic areas; the second is comprehensive description deriving from the results of pilot research.

The principal feature which both types of effort have in common is the determination of the parameters of interlingual matching — parameters which become the foundation for subsequent structural description within a language.

The first step in the descriptive process is to determine the corresponding parts to be matched between pairs of sentences. For pilot descriptions we compare real pairs of sentences from actual texts. This channels the descriptive effort in the most productive direction toward our purpose : it focuses attention on the constructions characteristic of the particular area of scientific discourse with which we are dealing. For later comprehensive description, abstract classes of constructions on all levels are considered as well. Both traditional and more recent descriptive literature is consulted in all phases of research.

The initial determination of corresponding parts is done by submitting a pair of sentences to the judgment of bilingual informants. They specify the corresponding elements, providing us with the interlingual matching specifications. We then prepare a surface structure description within the framework set by these interlingual matching specifications.

Let me give you an example of our present and projected activity, using a pair of sentences from Russian and English. These demonstrate

some of the problems of matching various kinds of structures. The sentences are given in item (0) of the annex. The matching of lexical items is obvious so that we can write the necessary rules, such as, items (1), (2), and (3). The symbols N_k , N_l , etc. denote form classes of nouns. The classes are determined mainly by paradigmatic features, but extraparadigmatic, that is, syntagmatic, features are also considered. The symbols T_a , T_b , ..., T_z denote interlingual transfer classes. The transfer classes establish the correspondence between entries in pairs of monolingual grammars.

Morphological relationships must also be accounted for. They occasion such entries in the rule inventory as in item (4). Symbols NO_{FN} and $NO_{sing_{AN}}$ denote classes of inflected nouns. The linking symbol T , denotes the interlingual transfer class. It contains all inflecting entries for the singular. For a comprehensive description of a language, it is necessary to write the entire set of such entries into the inventory of rules. This would have to cover all instances of morphological inflection — for all adjectives, nouns, verbs, and any other inflecting classes.

Moreover, some constructions have greater complexity. In these sample sentences, the noun phrases exemplified by *этой окрестности* or *(of) this area*, for instance. Their surface structures are to be accounted for by items (5), (6), and (7). Here for the Russian, the symbol NP_{FG} denotes the noun phrase construction feminine and genitive. Similarly, the symbol DT denotes determiners or articles and the symbol NO , the classification of inflected nouns. Here, *этой* is classified as a determiner (DT). Full subscripts, such as $FGDPI$, thus indicate the gender and the range of cases which characterize the form or construction so classified. In this instance, the gender is feminine and the case range would be genitive, dative, prepositional, and instrumental.

Then, in a similar fashion, corresponding structures are established for the English counterpart. For this, there may be such entries as items (8) and (9). Here, the symbol NP_{sing} classifies the noun phrase construction as singular. Symbols like $NO_{sing.A}$ denote a further classification of inflected nouns. They are used to include some extraparadigmatic features as cooccurrence with the indefinite article *a* or *an*.

If we look at a tree model of the description outlined here, we will see an analysis like item (10) for the Russian noun phrase and item (11) for the English noun phrase. You can see that the two parsings do not match element for element.

Nevertheless, the nodes marked NP at the top correspond for the phrases which are subsumed by these nodes. Likewise, the nodes DT_{FGDPI} and DT correspond for the expressions *этой* and *this*. Finally, there is correspondence between the symbols NO_{FGDP} and $NO_{sing_{AN}}$ for *окрестности* and *area*, respectively.

What we now need is to establish interlingual transfer correspondence between such differing but equivalent nodes in a pair of parsing structures.

It is in such areas that recent efforts at the Linguistics Research Center have been focused. We have developed parsing algorithms of a very wide generality. Therefore, we are able to write transfer entries over these and more complex trees.

Such transfer entries would present the descriptions to the parsing algorithms in such a guise as item (12). In this example, the * denotes that the entries shown do not exist as explicit phrase structure rules in the respective grammars. Rather, they are implicitly defined in the interlingual transfer system. Such a mechanism defines the class of bi-directional transformations on the surface structures.

To further illustrate this feature, let us consider the problem of correlating discontinuous elements in a string. In the pair of example sentences there are two occurrences of discontinuities. One is in the verb phrase illustrated in item (13). The other is in the correlative conjunction structure in item (14). The surface structure entries devised for part of the structure over the two verb phrases are given in items (16) through (23). Interlingual transfer is shown by the classes T_k , T_b , and T_n items (15), (16) and (18). These transfer mappings are straightforward and they will generate the matrix structure into which the rest of the structures for the lexical discontinuity will fit. Transfer is shown later for items (17) and (19) through (23). In diagram form, the matrix structure for Russian would appear as in item (24). A corresponding structure would be generated in the output for English. The more complex problem remains, that of accounting for the discontinuities in the lexical string.

For this, the notational form (in our system of approach to such discontinuities) is similar to that given earlier for less complex problems. The list of entries pertaining to discontinuous elements is given in items (25) through (28). The symbol *ZAVIS.OT* denotes the class specifying the discontinuous concatenation of the rule given in item (17) with the class *OT*. If rule (28) were applied to rule (25) we would generate the structure in item (29).

The parsing algorithms which manipulate the descriptive rules may be regarded as composed of two main systems : primary and secondary. The primary algorithm manipulates entries of the kind shown in solid lines, for instance the structure in item (20). The secondary algorithm manipulates entries of the kind shown in broken lines, for instance the structure in item (31). In items (25) through (27) the symbol *d* denotes a function similar to + in phrase "structure rules. However, unlike +, *d* does not denote an immediate concatenation between symbols. Instead, *d* denotes that the two elements concerned are more loosely associated, that is, in a discontinuous manner. See for instance the association of solid structures in item (29).

We may generate the corresponding English structure similarly from the class *DEPEND.ON* as in item (32). By applying all of items (25) through (28) we may generate the more complex structure with the correlative conjunctions as in item (33). Careful inspection of item (33) will

reveal that each of the uppermost symbols over solid lines fits into the matrix structure in item (24). I have not illustrated the English counterpart, since it is structured quite similarly.

As I have mentioned, the relationship between primary and secondary algorithms has been kept purposely general. In this way, we can accommodate mappings or translation problems which may involve both primary and secondary algorithms in one language and only the primary algorithm in the other. An example of this type of problem is to be found in the translation of the noun phrase *железной дороги* into *railroad*. Here, the Russian input is regarded as discontinuous with respect to translation, since the lexical stems are interrupted by inflectional endings. Item (34) shows a parsing of the Russian string. By writing the set of transfer rules shown in items (35) and (36), we may map onto the English structure in item (37). The new notation $\langle \alpha, \alpha \rangle$ and $\langle \alpha \rangle$ represents the interlingual mapping mechanism. It records the mapping equivalences of class symbols, between grammars. Thus, the notation in item (35) denotes that the lexical contents of the positions α and α in Russian are to be mapped into the single position α in English. It is also possible to perform the reverse operation: under the condition that the matrix structure has been defined, one content position may be mapped into many. Mapping notation was not used earlier because the relationships were straightforward. In actual practice, however, the relationship must be defined for each set of equivalent interlingual transfer entries.

There is one final problem area of interest in the examples, and that is the generation of articles in English. This was briefly touched upon earlier. I do not intend to suggest that the descriptive problem has been solved, but only that the parsing algorithms are geared to manage the descriptive data. In translating the phrase *построения железной дороги* we may generate several alternatives such as item (38). For the alternative *the building of a railroad*, items (39) and (40) list the transfer entries necessary for mapping the Russian and English. Item (39) will generate the Russian structure shown in item (41) and the English structure shown in item (42). When item (40) is applied it will fill in the missing inflectional step in both languages. The symbol *VO* in the secondary parsing denotes the class of constructions which will come to account for transformation between verb plus object and deverbalized noun plus object.

The symbol $N_{DE-VERB}$ denotes this class of deverbalized nouns. In English, such must be accompanied by some matrix of expression like *the ... of a ...* whenever the translation alternative with *of* is chosen.

For the present, the descriptive solutions offered here are but projections, since research into description utilizing secondary algorithms has only begun.

However, we have verified the performance of the primary algorithms with a reasonably large data base, one which is very comprehensive for English, German and Russian noun phrases (exclusive of relative clauses).

This data base contains dictionary entries numbering some 190,000, 40,000 and 140,000 items respectively. Entries describing syntactic constructions number several thousand in each language. Small, pilot descriptions of a few hundred entries have been tested in French, Spanish, Hebrew, Japanese, and Chinese. With the development of the secondary algorithms, we are beginning research into the description of problems involving discontinuity, transformations, and semantic features. In conclusion, one might say, "successful description of these areas will depend both on the developing of the secondary algorithm and on good fortune". I am indebted to my colleagues, Mrs. H-J. Hewitt and Mr. S. Whelan, for their aid in matters of presentation.

REFERENCES

- R. W. JONAS, *Generalized Translation of Programming Languages*, AFIPS Conference Proceedings, 1967 (Fall joint Computer Conference, Vol. 29 (in-press)).
- R. W. JONAS, *System Design for Computational Linguistics* (Austin, 1967).
- W. P. LEHMANN, *An Experiment in Machine Translation*, "The Graduate Journal", (1965) 7.111—131.
- W. P. LEHMANN, *Research in German-English Mechanical Translation*, Technical Report RADC-TR-67-98 (Griffiss Air Force Base, 1967).
- W. P. LEHMANN, *Research on Syntactic and Semantic Analyses for Mechanical Translation*: (Austin, 1967).
- L. W. TOSH, *Content Recognition and the Production of Synonymous Expressions*, Proceedings of the ninth International Congress of Linguists (The Hague, 1964), p. 722—729.
- L. W. TOSH, *Development of Automatic Grammars*, "Linguistics" (1965), 12.49—60.
- L. W. TOSH, *Initial Results of Syntactic Translation at the Linguistics Research Center*, "Linguistics" (in press).
- L. W. TOSH, *Syntactic Translation* (The Hague, 1965).
- L. W. TOSH, *Translation Model with Semantic Capability*, "Linguistics" (in press).

ANNEX

(0) **Индустриализация этой окрестности зависела и от построения железной дороги и от электрификации.**

Industrialization of this area depended both on the building of a railroad and on electrification.

(1) $[N_x \rightarrow \text{индустриализация}] \leftarrow T_x \rightarrow [N_x \rightarrow \text{industrialization}]$

(2) $[N_x \rightarrow \text{окрестность}] \leftarrow T_x \rightarrow [N_x \rightarrow \text{area}]$

(3) $[N_x \rightarrow \text{ящур}] \leftarrow T_x \rightarrow [N_x \rightarrow \text{hoof-and-mouth disease}]$

$$(4) \left[\begin{array}{l} \text{NO}_{FN} \rightarrow \{N_k + \text{Я}\} \\ \text{NO}_{FG} \rightarrow \{N_y + \text{А}\} \\ \text{NO}_{FG} \rightarrow N_x + \text{И} \\ \vdots \end{array} \right] \leftarrow T_1 \rightarrow \left[\begin{array}{l} \text{NO}_{\text{sing}\cdot\text{AN}} \rightarrow \{N_a\} \\ \text{NO}_{\text{sing}\cdot\text{AN}} \rightarrow \{N_b\} \\ \text{NO}_{\text{sing}\cdot\text{AN}} \rightarrow \{N_c\} \\ \vdots \end{array} \right] + \emptyset$$

$$(5) \text{NP}_{FG} \rightarrow \text{DT}_{FG} \rightarrow \text{NO}_{FG}$$

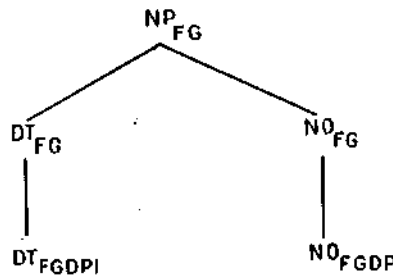
$$(6) \text{DT}_{FG} \rightarrow \text{DT}_{FGDPI}$$

$$(7) \text{NO}_{FG} \rightarrow \left\{ \begin{array}{l} \text{NO}_{FGDP} \\ \text{NO}_{FGDPI} \\ \dots \end{array} \right\}$$

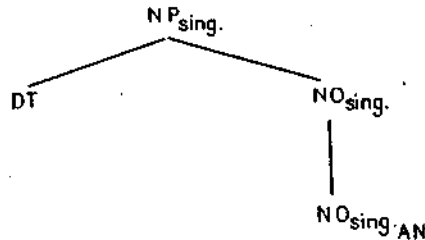
$$(8) \text{NP}_{\text{sing}\cdot} \rightarrow \left\{ \begin{array}{l} \text{DT} \\ \text{DT}_{\text{sing}\cdot} \end{array} \right\} + \text{NO}_{\text{sing}\cdot}$$

$$(9) \text{NO}_{\text{sing}\cdot} \rightarrow \left\{ \begin{array}{l} \text{NO}_{\text{sing}\cdot\text{A}} \\ \text{NO}_{\text{sing}\cdot\text{AN}} \end{array} \right\}$$

(10)



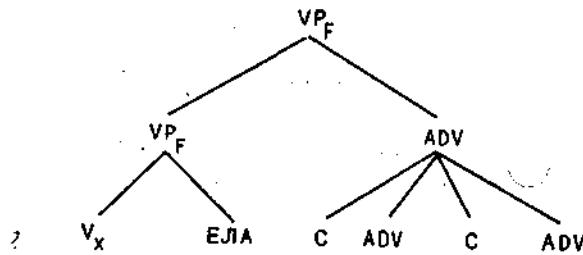
(11)



$$(12)^* [\text{NP}_{FG} \rightarrow \text{DT}_{FGDPI} + \text{NO}_{FGDP}] \leftarrow T_1 \rightarrow *[\text{NP}_{\text{sing}\cdot} \rightarrow \text{DT} + \text{NO}_{\text{sing}\cdot\text{AN}}]$$

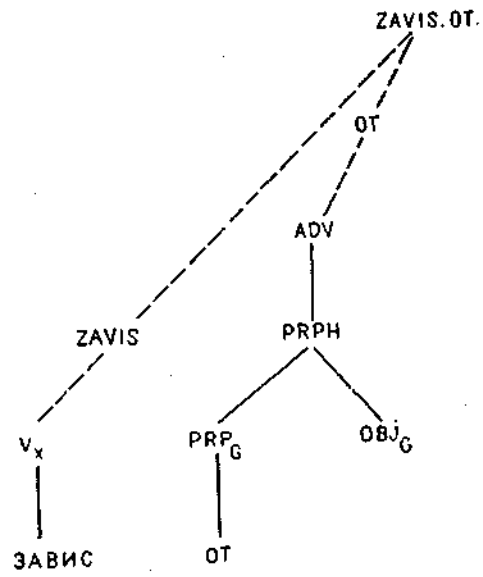
- (13) зависела ... от ... от
depended ... on ... on
- (14) и ... и
both ... and
- (15) $[VP_F \rightarrow VP_F \rightarrow ADV] \leftarrow T_x \rightarrow [VP \rightarrow VP \rightarrow ADV]$
- (16) $[VP_F \rightarrow V_x + ела] \leftarrow T_1 \rightarrow [VP \rightarrow V_x + ed]$
- (17) $V_x \rightarrow$ завис $V_x \rightarrow$ depend
- (18) $[ADV \rightarrow C + ADV + C + ADV] \rightarrow T_m \leftarrow [ADV \rightarrow C + ADV + C + ADV]$
- (19) $C \rightarrow$ и $C \rightarrow$ both
- (20) $ADV \rightarrow$ PRPH $ADV \rightarrow$ PRPH
- (21) $PRPH \rightarrow$ PRP_G + OBJ_G $PRPH \rightarrow$ PRP + OBJ
- (22) $PRP_G \rightarrow$ OT $PRP \rightarrow$ on
- (23) $C \rightarrow$ and

(24)

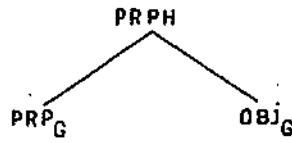


- (25) $*[ZAVIS. OT \rightarrow [V_x \rightarrow$ завис $] \text{ d } OT] \leftarrow T_n \rightarrow [DEPEND. ON \rightarrow [V_y \rightarrow$ depend $] \text{ d } ON]$
- (26) $[OT \rightarrow CC \text{ d } OT \text{ d } OT] \leftarrow T_o \rightarrow [ON \rightarrow CC \text{ d } ON \text{ d } ON]$
- (27) $[CC \rightarrow [C \rightarrow$ и $] \text{ d } [C \rightarrow$ и $] \leftarrow T_p \rightarrow [CC \rightarrow [C \rightarrow$ both $] \text{ d } [C \rightarrow$ and $]]$
- (28) $[OT \rightarrow *[ADV \rightarrow OT + OBJ]] \leftarrow T_q \rightarrow [ON \rightarrow *[ADV \rightarrow on + + OBJ]$

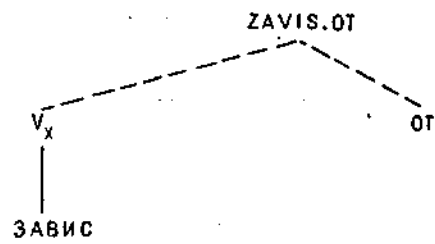
(29)



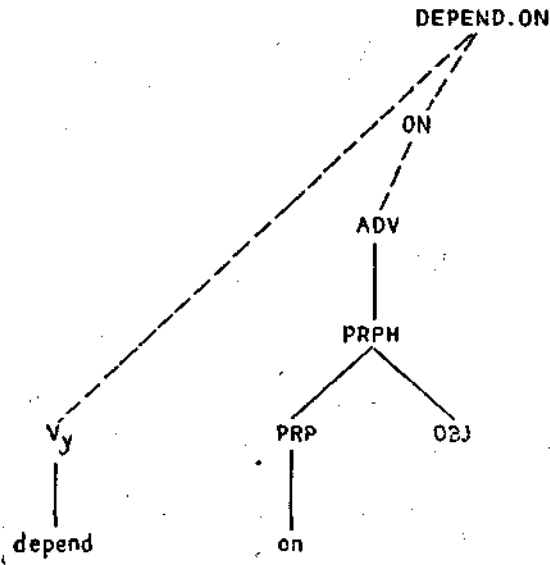
(30)



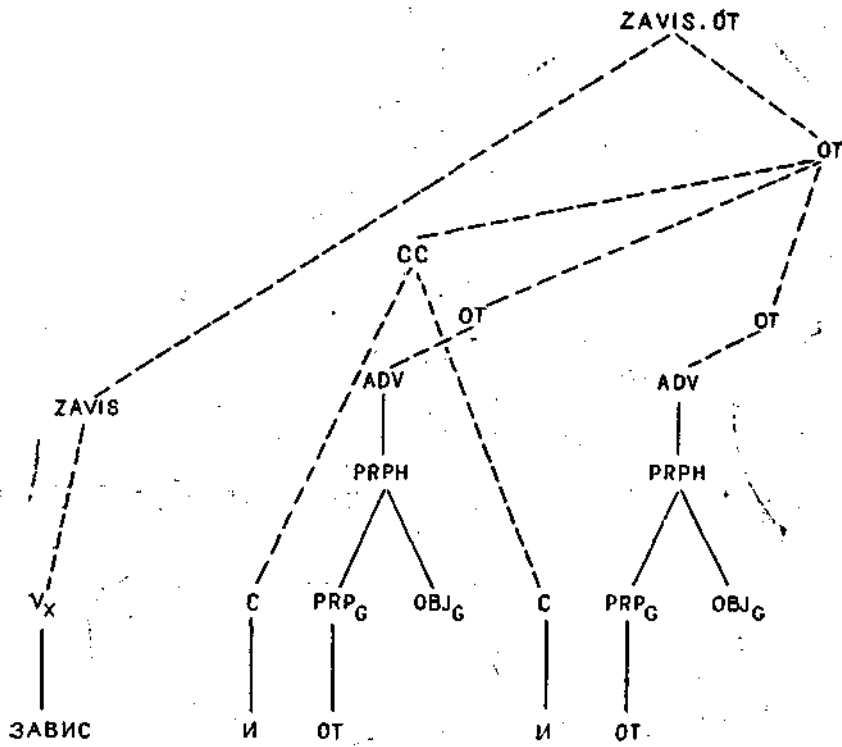
(31)

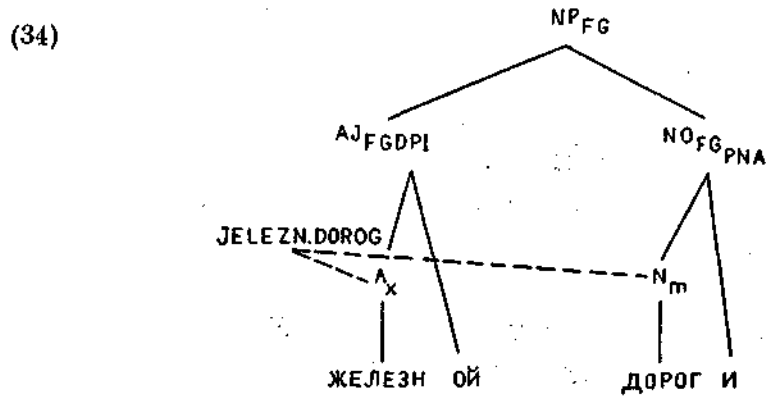


(32)



(33)





(35) $*[NP_{FG} \rightarrow A_x + ой + N_m + и] \leftarrow T_r \rightarrow *[NO_{sing.} \rightarrow N_b]$

(36) $[JELEZN.DOROG \rightarrow [A_x \rightarrow железн] \text{ d } [N_m \rightarrow дорог]] \leftarrow T_s \rightarrow [N_b \rightarrow railroad]$

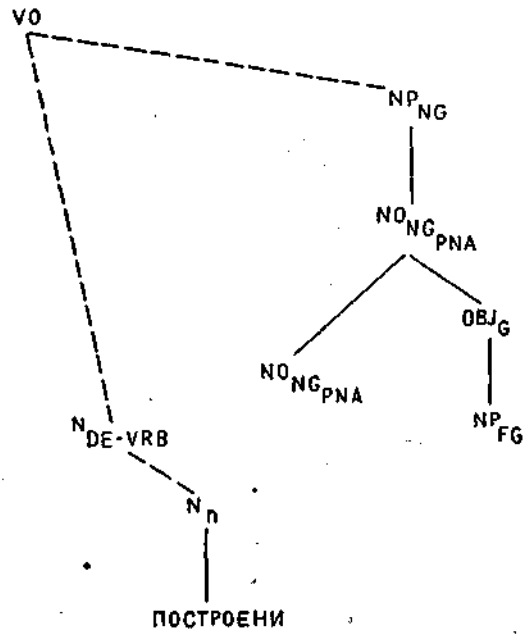


(38) building $\left(\left\{ \begin{matrix} a \\ the \end{matrix} \right\} \right)$ railroad (s)
the building of $\left(\left\{ \begin{matrix} a \\ the \end{matrix} \right\} \right)$ railroad(s)

(39) $[VO \rightarrow * [NP_{NG} \rightarrow NO_{NG}PNA + NP_{FG}] \text{ d } [N_n \rightarrow построена]]$
 $[VO \rightarrow * [NP_{sing.} \rightarrow the + NO_{sing.} + of \text{ a } + NO_{sing.}]]$
d $[N_b \rightarrow building]$

(40) $[NO_{NG}PNA \rightarrow N_n + я] \rightarrow T_i \rightarrow * [NO_{sing.} \rightarrow N_b]$

(41)



(42)

