

CO-ORDINATIVE ELLIPSIS IN RUSSIAN TEXTS:
PROBLEMS OF DESCRIPTION AND RESTORATION

Igor A. BOLSHAKOV

VINITI, Academy of Sciences of USSR
Moscow, 125219, USSR

ABSTRACT. Russian elliptic constructions are examined from the point of view of syntactic analysis. Reciprocal elements in a co-ordinative elliptic sentence are exposed and possible types of their similarity are explored. Linear formulae of ellipsis for most textual cases are constructed and statistics of their use is discussed. As a result the main steps of ellipsis restoration algorithm are outlined.

INTRODUCTION. The investigations of ellipsis (gapping) in natural language sentences with structural methods have been carried on for more than 20 years, but algorithms of automatic restoration of omitted words either in coherent Russian texts or in sequences of Man-Machine interaction replicas have not yet been proposed. The problem is still topical. Indeed, at an average each 7th entry in Great Soviet Encyclopedia and 25th abstract in a common Soviet abstract journal contain at least one sentence of a following kind:

В первой базе данных содержится информация по всем отраслям знаний, во второй - только по гуманитарным наукам.

/The first database contains information on all knowledge fields, the second one - only on the Humanities./

Conversion of such a sentence to a formalized language, e.g. for automatic updating of factographic databases, without the ellipsis restoration is impossible. Meanwhile ellipsis in Russian sci-tech texts is very diversified and covers any part of a sentence and most frequently the predicate with adjacent words.

Early Soviet works /Leontieva, 1965/, /Paducheva et al., 1973/, /Korelskaya et al., 1973/ had examined the phenomenon from the point of view of synthesis. But when synthetic rules transform a non-empty entity to an empty one, corresponding analytical rules are not their trivial reversion. We examine elliptic constructions in co-ordinative sentences with orientation to analysis, i.e. to parsing algorithms including restoration of omitted words.

Several important issues should be forced for our purpose: 1) introducing the notion of reciprocal elements in non-omitted parts of elliptic sentences; 2) defining new types of syntagmata for restoration of semi-destroyed links between words within the reciprocals; 3) exploring possible types of reciprocal similarity; 4) describing co-ordinative sentences with minimal number of elliptic transformation formulae; 5) collecting statistics of the formulae use, which implies a search order of a specific formula for a given sentence. Thus a base for outlining the main steps of an ellipsis restoration algorithm is formed.

RECIPROCAL ELEMENTS. There are several kinds of word omission in natural languages. Among them co-ordinative reduction is carried out according to the formula $XQ \& X^*Q$

$\Rightarrow X \& X^*Q$ or $QX \& QX^* \Rightarrow QX \& X^*$, where X and X^* are different words or word groups, Q is a recurrent group of words, and $\&$ stands for a co-ordinative conjunction or just a comma: (He took)_Q (the bread)_X (and)_& (he took)_Q (the milk)_{X*} \Rightarrow (He took)_Q (the bread)_X (and)_& (the milk)_{X*}.

Co-ordinative ellipsis includes co-ordinative reduction as a subset. Two or more phrase segments co-ordinated in a single sentence have some identical parts not necessarily standing at the borders of the segments. For economy's sake the natural language omits, wholly or partially, the recurrent part of a segment, as a rule, in the second or subsequent one. The omission may be carried out according to the formula $XQ_0 Y \& X^*Q_0 Y^* \Rightarrow XQ_0 Y \& X^* - Y^*$, where X and X^* are a pair of somewhat similar, but compared and even semantically opposed elements referred to hereinafter as reciprocal; Y and Y^* are another such pair; the co-ordinative segments are $XQ_0 Y$ and $X^* - Y^*$. In Russian writing ellipsis is often accompanied with a dash.

The common feature of elliptic omissions is that a connected dependence tree for a sentence cannot be built without their restoration. The parsing algorithm should take into account and properly distinguish the specific phenomena occurring in a given sentence, e.g. ellipsis and zero copula going together.

BASIC ASSUMPTIONS. We examined a co-ordinative ellipsis assuming the following:

- A context confined to a standalone elliptic sentence is sufficient for restoration.

- An ellipsis in a sentence is explicated with a dash (along with some implied indications), but omission may occur not where the dash is or not only there. Comparative and other constructions without a dash have not been examined, but most of our statements cover them, too.

Omitted parts, after restoration and, if necessary, morphological re-agreement, exactly restore the original meaning of a given sentence as a whole. Therefore recurrent parts in different segments have had exactly the same length, linear order, and lexical content before omission.

- A dependency tree was selected as a tool for description of links between word-forms in a sentence, and the nature of the links corresponds to those suggested in "Meaning - Text" model /Mel'čul, 1973/.

- Co-ordinatively reduced words should be restored, along with the ellipsis proper, where this provides connectedness of a parse tree and saturation of obligatory valencies.

- There exists an algorithm for syntactic analysis of Russian texts which can cope with parsing any ellipsis-free sentence. It can also parse, consistently and unambiguously, the complete segments of an elliptic

sentence and tie links, even ambiguous, within word groups in incomplete segments, preparing them as disjointed "bushes" for a final parse tree.

STATISTIC OBSERVATIONS. As many as 600 elliptic sentences were extracted from large (more than 1000 pages) corpus of texts in various fields and genres (abstracts, articles, popular science books, and brochures). The material was arranged in a minimal number of elliptic formulae. A total of 24 formulae was found, but only four of them exceeded the 4% threshold (see Table).

In the Table the arrows show the direction of synthesis; Q_0 stands for an omitted part replaced with a dash; Q_1, Q_2 are additional omissions possible to the left and to the right of the dash; P and P_1 are optional segments (modifiers and the like), not involved in the ellipsis.

Available statistics permits to conclude:

- The most widely used formula (single omission between two reciprocals) has exceeded the use total of all other formulae, and four leading formulae together account for more than 80 per cent.

- Number N_1 of reciprocal pairs and number N_0 of omissions in a sentence usually satisfy the formula $|N_1 - N_0| \leq 1$, since reciprocals and omissions are commonly interleaving each other, but it is possible to construct counter-examples.

- Cases with $N_0 > 2$ and/or $N_1 > 2$ amount for less than two per cent.

- Approximately in 7% of cases nonsymmetric ("chlastic") formulae occurred. For example, formula $P QXY \& X^* QY^* \Rightarrow P QXY \& X^* * Y^*$ may have the following realization:

(Совместная работа)_P (требует)_Q (от инженера)_X (расширения математического кругозора)_Y (и)_& (от математика)_{X*} - (владения некоторыми математическими знаниями)_{Y*} / (A joint work)_P (requires)_Q (from engineer)_X (broadening of his mathematical outlook)_Y (and)_& (from mathematician)_{X*} - (the possession of some technical background)_{Y*} /

- Approximately four per cent of all cases did not fit into our formulae. Either a possible formula was too complex to be practical, or the shape of the sentence was dubious and even incorrect from the point of view of a human editor.

ENCLOSING SYNTAGMATA. Before searching reciprocals it is necessary to establish syntactic links within word groups in an incomplete segment. The convenient tools for describing these links, the so-called syntagmata, are not always sufficient here, see, e.g. the word combinations: **обобщенные экономические/отдельные производственные (показатели)** /generalized economic/some productional (indications)/; **17 старых/пятью новыми (предприятиями)** /by 17 old/by five new (enterprises)/; **из капиталистических/из социалистических (стран)** /from capitalist/from socialist (countries)/; **число имеющихся/число вновь приобретенных (книг)** /the number of available/the number of newly bought (books)/; **80 тыс. публицистических/200 тыс. научных (публикаций)** /80,000 publicistic/200,000 scientific (publications)/. Reciprocals in the

pairs above are divided with slashes, and connected nouns (derived from complete segments) are given in parentheses.

We propose dependencies of a new type, named **enclosing syntagmata (ES)**. These are established either between co-subordinated attributes, linking brother nodes from right to left (**обобщенные - экономические**) or between a number and co-subordinated attribute (**17 - старых**) or between a proposition or a noun with predicative properties and subordinate noun's attribute (**из - капиталистических, число - имеющихся, тыс. - публицистических**).

Several ES constitute the links between a dominating word and its indirectly subordinated one. Via ES the nodes most important for reciprocal matching turned to be within easy distance from their bushes' roots (may be, at the very roots). Similar links can be established in complete segments, too. They transform a relevant subtree to an acyclic graph, which facilitates matching nodes and arcs.

The set of ES is subject to updating so far. Caution should be used however about updating. Indeed, attempts to directly link words arbitrarily distant within a convenient-dependency tree, though eliminating the very notion of ellipsis, lead to superfluous complexity of ES and of a global model of natural language, too.

SIMILARITY FEATURES. Manual segmentation of all available sentences has made clear that antagonists in reciprocal pairs are not in general case mutually isomorphic, i.e. their subtree do not quite coincide. Therefore, the labels at the matched nodes should be forcedly involved and these are of the following types:

Lexical. Lexemes at the roots and/or their direct subordinates in 22 per cent were strictly coincide.

Morphological. In most cases lexemes at the matched nodes belonged to the same part of speech, and their wordforming characteristics were in agreement: nouns and numbers expressed by words - in case; adjectives (pronominal included) and participles - in gender, number, and case; personal verb-forms - in gender, number, and person. Only in 30 per cent of cases the agreement has not been revealed (numbers in digits, abbreviations, etc.).

Syntactic. Some indicators treated in the "Meaning - Text" model as syntactic, might coincide, e.g. interrogativity of lexemes **КУДА** and **СКОЛЬКО** in the sentence **Я спросил, куда идти, а он - сколько времени.** /I asked where I should go, and he asked what time it was./

Semantic. If matching labels of the three kinds mentioned failed or at once several nodes in a complete segment were similar to the node in an incomplete one, then coincidence of even one semantic indicator is important. Taxonomy of such indicators is not established yet. Several facets (classification aspects) with admissible intersection of their scopes are expected to suit well, but simple thesaurical hierarchies are not excluded, either. We have specifically observed: quantitative and cardinal words (28%), synonymy and antonymy (10%), hyponymy and hyperonymy (i.e. genus/species, part/whole) (3%), and other kinds of such similarity (2%).

When semantic similarity within pairs X/X^* and Y/Y^* failed, semantic proportion $X/Y = X^*/Y^*$ has been sometimes observed,

but algorithmic verification of the proposition is difficult.

DRAFT ALGORITHM. An algorithm we propose for restoration of co-ordinative ellipsis (CE) consists of the following main steps:

1) A basic ("ellipsis-free") algorithm in carrying on syntactic parsing, sentence by sentence. CE-flag is simultaneously derived from the three other flags signalling connectedness of the parse tree, co-ordinativity of the sentence structure, and presence of a zero copula.

2) If the CE-flag is set for the sentence then algorithm of ellipsis restoration starts up. Given borders of the incomplete segment and a dash position, X^0 and Y^0 elements are extracted. They should be inwardly linked via connectors and, if necessary, enclosing syntagmata.

3) The X^0 and Y^0 found, their reciprocals X and Y in the complete segment are searched starting from its borders inwardly. The node-by-node matching for all nodes starts from their roots, uses all node labels, and should be continued at lower tree level, if function words are encountered at the roots.

4) The X and Y found, dependency links in the complete segment are traced, mainly in ascending order. The number and location of borders between reciprocals and additions are thus revealed.

5) The specific CE-formula is diagnosed and omitted word groups are duplicated with

its help from the complete segment to the incomplete one. Re-agreement of the restored words is carried out, if necessary, e.g. in number, person, and/or gender for a predicate.

If there are more than two co-ordinative segments in the sentence, then steps 3 to 5 are repeated for every subsequent segment applying preceding (complete or just restored) one.

ACKNOWLEDGEMENTS. The author wishes to express his great gratitude to I.G.Bider, L.G.Lovdin, L.G.Mitushin, E.V.Paducheva, and especially to Yu.D.Apresian for their keen critique. He is also thankful to M.S. Pradkin and A.A.Raskina for editorial remarks.

REFERENCES

Isodicheva, E.N. (1965) 'Analysis and synthesis of Russian elliptic sentences.' Nauchno-tekhnicheskaya informatsiya, 11, Moscow: VINITI, pp. 44-46.
 Paducheva, E.V. & Idashchenko, T.K. (1973) 'Ellipsis as zero anaphoric sign.' Nauchno-tekhnicheskaya informatsiya, Ser. 2, 5, Moscow: VINITI, pp. 20-31.
 Korelskaya, T.D. & Paducheva, E.V. (1973) 'Transformation in symmetric constructions: co-ordination and ellipsis.' Nauchno-tekhnicheskaya informatsiya, Ser. 2, 9, Moscow: VINITI, pp. 29-38.
 Mel'čuk, I.A. Experience in the theory of linguistic "Meaning - Text" models (in Russian). Moscow: Nauka, 1973.

TABLE. The four most frequent ellipsis formulae

#	%	Formula	Russian example	Word-by-word translation
1	61.7	$[P_0]XQ_0Y \& X^0Q_0Y^0 \Rightarrow [P_0]XQ_0Y \& X^0 - Y^0$	(Без округления) P_0 (при сдвиге доопыляющих) X (кодов ошибки будет) Q_0 (отрицательной) Y (,а) Q_0 (при сдвиге обратных) X^0 (положительной) Y^0	(Without rounding of) P_0 (when shifting additional) X (codes an error will be) Q_0 (negative) Y (, and) Q_0 (when shifting reverse) X^0 - (positive) Y^0
2	3.0	$[P_0]XQ_0YQ_1 \& X^0Q_0Y^0Q_1 \Rightarrow [P_0]XQ_0YQ_1 \& X^0 - Y^0$	(2/3) X (микромИМ используют-ся) Q_0 (в коммерческой) Y (сфере) Q_1 (,а) Q_1 (1/3) X^0 (в научно-технической) Y^0	(2/3) X (of microcomputers are used) Q_0 (in commercial) Y (field) Q_1 (and) Q_1 (1/3) X^0 (in sci-tech) Y^0
3	6.7	$[P_0]Q_0X \& P_1Q_0X^0 \Rightarrow [P_0]Q_0X \& P_1 - X^0$	(Если температура одинакова, то) Q_0 (плазма называется) Q_0 (изотермической) X (,а) Q_0 (в противном случае) P_1 (неизотермической) X^0	(If the temperature is identical then) P_0 (plasma is named) Q_0 (isothermic) X (and) Q_0 (in opposite case) P_1 (nonisothermic) X^0
4	3.7	$[P_0]Q_1XQ_0Y \& Q_1X^0Q_0Y^0 \Rightarrow [P_0]Q_1XQ_0Y \& X^0 - Y^0$	(В 1982 г.) P_0 (фонд) Q_1 (университетских) X (библиотеч-ных составили) Q_0 (600 тыс. т) Q_0 (,а) Q_0 (специальных) X^0 (порядка 40 тыс. т) Y^0	(In 1982) P_0 (the resources) Q_1 (of university) X (libraries' account for) Q_0 (600,000 vol.) Y (and) Q_0 (of specialized) X^0 (approximately 40,000 vol.) Y^0