

TOWARDS BETTER UNDERSTANDING OF ANAPHORA

Barbara Dunin-Kępcicz
Institute of Informatics, Warsaw University
P.O. Box 1210
00-901 Warszawa, Poland

ABSTRACT

This paper presents a syntactical method of interpreting pronouns in Polish. Using the surface structure of the sentence as well as grammatical and inflexional information accessible during syntactic analysis, an area of reference is marked out for each personal and possessive pronoun. This area consists of a few internal areas inside the current sentence and an external area, i.e. the part of the text preceding it. In order to determine that area of reference several syntactic sentence-level restrictions on anaphora interpretation are formulated.

Next, when looking at the area of pronoun's reference, all NPs which number-gender agree with the pronoun can be selected and this way the set of surface referents of each pronoun can be created. It can be used as data for further semantic analysis.

I INTRODUCTION

Reference is one of the central concepts of any linguistic theory. In recent research into anaphora the term "reference" has been used in three different senses (Szwedek, 1981):

- (a) as a relation between the name and the thing named (Hall Partee, 1978)
- (b) as an association between noun phrases and mental entities in the language user's (Nash-Webber, 1978)
- (c) as an association between the occurrence of phrases in the text (Reinhart, 1981)

However the reference is understood, in order to interpret correctly anaphora on the semantic level ((a) and (b)), first a stage (c) is necessary.

In this paper I have taken the point of view presented under (c). I shall discuss the problem of anaphora in Polish sentences. My attention is focused on personal and possessive pronouns explicitly occurring in

the text and moreover on zero pronouns, i.e. ellipsis of NP in the subject position, specific for Slavonic languages.

My purpose in the description of regularities of the reference in the Polish language, I shall express them by defining the area of pronoun's references, i.e. those regions of the text where its antecedents should be found. These surface referents will be selected from among NPs occurring in the sentence.

The research on anaphora made for English has led to the formulation of some structural rules using such relations as command, c-command and precede-and-command (Reinhart, 1981).

I have been searching for analogous rules for Polish. But two essential differences have to be considered:

- (i) grammatical and morphological properties of Polish and English;
- (ii) different grammatical traditions.

For English the rules concerning the coreference of entities were formulated on the basis of generative-transformational grammar. For Polish the first precise description of Polish syntax was formulated only recently by Szpakowicz, who based his work on the framework created by Saloni (Saloni, 1976; Saloni and Swidzinski, 1981). It is a kind of immediate-constituent grammar; the grammatical categories (case, gender, etc) are applied not only to single words, but also to compound phrases. In my present work I have limited my attention to the subset of Polish described by Szpakowicz (Szpakowicz, 1983).

Polish is a highly inflexional language and this fact has many and varied consequences. Surface referents of the pronoun will be selected from among those NPs which number-gender agree with the pronoun. Strictly speaking, the grammatical categories of the pronoun should be compatible with the categories of the NP, but in cases of neutralization they cannot be fully determined.

My method of determining the areas of pronoun's reference is a syntactic one, because it is based on morphological and syntactical properties of the Polish language. I assume

the availability of the surface structure of the sentence as well as grammatical and inflexional information accessible during a syntactic analysis. I deliberately do not make use of any semantic information, trying to get the most out of grammar. The feature I intend to provide is a complete definition of the area of pronoun's reference.

II AREA OF REFERENCE

A. Internal and external areas of reference

In the process of determining the surface referents of the pronoun, first the area of its reference should be marked out. This area, i.e. those regions of the text, where its antecedents should be found, is usually made up of several internal reference areas, i.e. the appropriate bits of the current sentence, and an external area, the part of the text preceding the current sentence. The list of internal areas depends on the syntactic position of the pronoun in the sentence. To determine these areas it is necessary to formulate sentence-level anaphora restrictions for Polish. These rules will determine the conditions of both obligatory coreference and obligatory non-coreference of entities. Thus we have two situations to consider:

- (i) in the case of obligatory coreference one internal area of reference containing the appropriate referent should be marked out;
- (ii) in the case of obligatory non-coreference the elements which are forbidden as surface referents of the pronoun should be excluded from the internal area.

The coreference of entities which is qualified on the basis of some other premises will be called admissible coreference.

At our disposal we have a multileveled, hierarchic surface structure of the sentence. Generally, it seems that internal areas can be identified with the constituents on the highest level: subject, objects, modifiers, regardless of their syntactic realization. Strictly speaking, noun as well as NP or any sentential structures can be instances of internal areas of reference.

The partitioning of sentence (1) illustrates it:

- (1) "(Ewa i Piotr) poszli (do niego) (z dziewczyną, którą właśnie spotkali)".

"Ewa and Peter went to him with a girl which just met".

B. Rules concerning coreference of entities in Polish

1. The basic criterion of excluding coreference

The following rules of excluding the coreference of entities concern a level

deeper than that on the surface, because they refer to syntactical functions of phrases in the sentence. The first rule presents the problem of coreference of the subject and other nominal groups, i.e. objects and nominal modifiers, in short called objects. It concerns reflexive pronouns, so it should be noted first that they differ from those in English, eg.:

- possessive pronoun "swoj" may have one of the following meanings: his, her, its.
- reflexive pronoun "siebie" can mean: himself, herself, itself, myself, ourself, yourself, themselves.

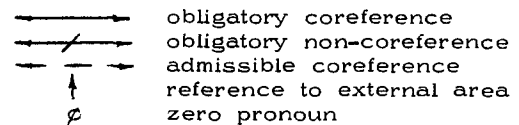
The basic criterion of excluding coreference I have formulated from the analytical point of view:

- (R 1) If the object is expressed by means of a reflexive pronoun, then it is coreferential with the subject; in other cases the referential identity of the subject and object is excluded.

This criterion is applied both to look for coreferents of objects - blocking the subject, and in testing the possible antecedents of the subject - blocking the objects.

Let us consider some examples:

Meaning of symbols:



- (2) "Ewa zapytała ją o to"
"Eva asked her about it"
- (3) "∅ zapytała ją o to"
"Asked her about it"
- (4) "Ona zapytała ją o to"
"She asked her about it"
- (5) "On zapytał Jana o Piotra"
"He asked John about Peter"
- (6) "Piotr nalał sobie piwa"
"Peter poured himself beer"

Rule R 1 holds for possessive pronouns:

- (7) "Ewa uwielbia swoją przyjaciółkę"
"Eva adores her friend"

Now let us have a look at the case of the preposed PPs so difficult to interpret in English. The basic criterion of excluding coreference covers these phrases too:

- (8) "Nagle, obok Jana, ∅ zobaczył węża"
"Suddenly, near John, saw a snake"

- (9) "Nagle, obok niego, \emptyset zobaczył węża"
 "Suddenly, near him, saw a snake"
 (10) "Nagle, obok siebie, \emptyset zobaczył węża"
 "Suddenly, near himself, saw a snake"
 (11) "Nagle, obok siebie, on zobaczył węża"
 "Suddenly, near himself, he saw a snake"

In examples (10) and (11) the reflexive pronoun has appeared. These are the only two cases in which the coreference with the subject of the main sentence is permitted and even obligatory. Such an interpretation is correct irrespective of the position of FP in the sentence, i.e. it does not depend on whether this phrase precedes or follows the subject.

The basic criterion of excluding coreference works as follows:

- (i) it is valid only for a simple clause, without blocking coreference between the elements of the main sentence and the constituents of embedded clauses;
- (ii) it is obligatory on every level of the sentence, i.e. it concerns all the sentence constructions irrespective of their position in the structure of the whole sentence.

Examples (12) to (14) illustrate this:

- (12) "Piotr nie wiedział, czy \emptyset pójdzie do kina"
 "Peter did not know, whether would go to the movies"
 (13) "Jan zapomniał, o co Piotr go pytał"
 "John forgot, what Peter asked him about"
 (14) "Jan spotkał chłopca, który go dawno nie odwiedził"
 "John met a boy, who didn't visit him for long"

The interpretation of reflexive pronouns is not so easy as the criterion R 1 suggests. These pronouns can be involved in various compound phrases which often are ambiguous. Especially infinitive phrases are hard to interpret. In order to do this correctly, an implicit agent which will be called further the deep subject, should be obtained. It often needs a few hypotheses to be formulated. Let us consider an example. The sentence:

- (15) "Jan kazał służącemu umyć się"

can be translated in two ways which exactly give the sense of possible Polish interpretations:

- (15.1) "John told (the servant) (to wash him)"
 (15.2) "John told (the servant) (to wash himself)"

In the infinitive phrase "umyć się" ("to wash him" or "to wash himself") which is standing in the object position, the reflexive pronoun "się" is coreferential with the deep subject of this phrase. Thus its interpretation has to be determined. Here we have two possibilities:

- (i) the previous object - "servant" - interpretation (15.1)
- (ii) the subject of the main sentence - "John" - interpretation (15.2)

One of them is the referent of the deep subject. And so we come to the next rule:

- (R 2) In order to interpret the infinitive phrase, the deep subject of the phrase has to be selected from among the previous object (if any) and the subject of the main sentence.

2. Excluding the coreference between objects

The next sentence-level restriction of anaphora interpretation regulates the problem of coreference of NPs other than a subject, i.e. objects, between them.

- (R 3) The coreference of particular objects is excluded. This is an obligatory non-coreference.

- (16) "Jan zapytał go o Piotra"
 "John asked him about Peter"
 (17) "Jan zapytał go o niego"
 "John asked him about him"
 (18) "Jan zapytał Piotra o niego"
 "John asked Peter about him"

This rule does not hold for possessive pronouns which in Polish do not create NPs by themselves. If these pronouns occur in objects, they may be coreferential with objects preceding them (admissible coreference).

- (19) "Jan zapytał Piotra o jego brata"
 "John asked Peter about his brother"

Rule R 2 is only valid for a simple clause, but it concerns all the sentence constructions irrespective of their position in the whole sentence.

3. Rules of interpreting compound sentences

The next group of problems concerns the coreference of entities in a compound sentence, including the question of the subject. In a Polish sentence it needs not be explicit. Ellipsis of the NP in the subject position, often called "the elided subject", is a natural way of expressing "thematic continuity" and exemplifies an unaccented position in the sentence. On the other hand, the pronoun as the subject stands in syntactic opposition to the elided subject (zero pronoun) and exemplifies an accented position in the sentence.

While determining the antecedent of the subject of a simple sentence or a main clause in a compound sentence (explicit or implicit) we reach out to the external area of references. However, the basic criterion of excluding coreference is still valid.

- (20) "On zapytał go o Piotra"
 "He asked him about Peter"

The interpretation of compound sentences is difficult and sometimes leads to ambiguous results. The following rules concern mainly the coreference (or non-coreference) of elided subjects in co-ordinate and subordinate clauses. In the case of co-ordinate clauses two rules can be formulated:

- (R 4) For each two clauses in a sequence, if the elided subject is in the second clause, then the subject of the first clause should be extrapolated there (obligatory coreference).

- (21) "Piotr wstał od stołu i podszedł do okna"
 "Peter left the table and approached the window"

- (R 5) For each two clauses in a sequence, the pronoun or zero pronoun subject in the first clause cannot be coreferential with the non-pronoun subject of the second clause (obligatory non-coreference).

- (22) "Ø wstał od stołu, a Piotr podszedł do okna"
 "He left the table and Peter approached the window"

Interpreting subordinate clauses depends on the relative position of the main and the embedded clause.

- (R 6) If the embedded clause precedes the main clause and if both have elided subjects, these have to be coreferential (obligatory coreference).

- (23) "Zanim Ø wyszedł, Ø zgasił światło"
 "Before left_{masc}, turned off_{masc} the light"

- (24) "Ponieważ Ø zapomniał, Ø zapytał o to"
 "Because forgot_{masc}, asked_{masc} about it"

- (R 7) The elided subject in the embedded clause is a natural way of indicating the nearest candidate - the previous object (if it is there) or the subject of the main sentence (admissible coreference).

- (25) "Jan zapewnił Piotra, że Ø pójdzie do kina"
 "John promised Peter, that will go to the movies"

- (R 8) The pronoun or zero pronoun subject in the main sentence can be coreferential with the non-pronoun subject of the embedded clause which precedes the main sentence (admissible coreference), but cannot be coreferential with the non-pronoun subject of the embedded clause following the main sentence (obligatory non-coreference).

- (26) "Zanim Jan wyszedł, Ø zgasił światło"
 "Before John left, turned_{masc} off the light"

- (27) "Ø zgasił światło, zanim Jan wyszedł"
 "Turned_{masc} off the light, before John left"

- (28) "On nie wiedział, czy Piotr pójdzie do kina"
 "He didn't know, whether Peter will go to the movies"

4. Interpretation of relative clauses

Relative clauses are quite easy to interpret in Polish. Either their subject or object is replaced with pronoun "which" or "what" or their equivalents (only such types of relative clauses are described in the Szpakowicz grammar). These pronouns always indicate the NP next to which they stand and inherit gender, number and person from it. Thus the obligatory coreference of relative pronoun and this NP is determined. Let us have a look at some examples:

- (29) "Ewa zaprosiła Anie, która Ø znana od dawna"
 "Eva invited Ann, which (object) had known fem for long"

(30) "Ewa zaprosiła Anię, która znała ją od dawną"

"Eva invited Ann, which had known her for long"
(subject)

III CONCLUSION

The above syntactic method of interpreting pronouns yields only partial results - the list of internal areas of reference or the external area, both with certain restrictions on coreference, are determined. Next, more detailed results can be obtained. When looking at the internal areas, all NPs which number-gender agree with the pronoun should be selected and a list of surface referents of pronoun together with a list of elements blocked as the referents can be drawn up. If no internal areas are marked out, the external area with the list of blocked elements is the result of the method presented here. Similarly, while only admissible coreference is determined, the external area is marked out too and the list of blocked elements remains valid. On the other hand the obligatory coreference makes it possible to define the appropriate antecedent of the pronoun. The list of surface referents may be ordered by assuming the specific method of traversing the parsing tree. I expect, that as for English, recency understood as a physical distance between the pronoun and its antecedent can be the first approximation of the probability.

As expected the results of the method applied here need semantic verification. But at the same time they are a reasonable data for further semantic analysis. Data arrived at in this way make this process much easier.

It seems that a similar procedure can be carried out for other languages. Full grammatical information should be used wherever it can simplify such complex process as the semantic analysis.

IV REFERENCES

- HIRST, Graeme (1979). Anaphora in Natural Language Understanding: A Survey. Dept. of Compute Science, University of British Columbia.
- HOEBS, Jerry R (1976). Computational Approach to Discourse Analysis. Artificial Intelligence Center, SRI International.
- HOEBS, Jerry R (1978). Coherence and Coreference. Technical note 168. Artificial Intelligence Center, SRI International.

NASH-WEBBER, Bonnie Lynn (1978). A Formal Approach to Discourse Anaphora. PhD thesis, Harvard University

PARTEE, Barbara Hall (1978). Bound Variables and Other Anaphors in: Waltz 1978, 79-85.

REINHART, Tanya (1981). Definite NP Anaphora and C-Command Domains. in: Linguistic Inquiry, Vol 12, No 4, Fall 1981.

SALONI, Zygmunt (1976). Cechy składniowe polskiego czasownika (Syntax Properties of Polish Verb). Ossolineum, Prace językoznawcze, 1976.

SALONI Zygmunt, SWIDZINSKI Marek (1981). Składnia współczesnego języka polskiego (Syntax of Contemporary Polish Language). Wydawnictwa Uniwersytetu Warszawskiego, 1981.

SZPAKOWICZ, Stanisław (1983). Formalny opis składniowy zdań polskich. (Formal Syntactic Description of Polish sentences). Wydawnictwa Uniwersytetu Warszawskiego, 1983.

SZWEDEK, Aleksander (1981). Word Order, Sentence, Stress and Reference in English and Polish. WSP Bydgoszcz, 1981.

V ACKNOWLEDGEMENTS

I would like to acknowledge Janusz Bień and Stanisław Szpakowicz for their helpful comments on this paper.