# EFFICIENT SENTENCE DISAMBIGUATION BY PREFERRED CONSTITUENT ORDER

*S. Kronenberg\*, K. Skuplik+*

University of Bielefeld

\*Technische Fakultät, +Fakultät für Linguistik und Literaturwissenschaft

P.O. Box 100131, 33501 Bielefeld, Germany,

{susanne,kristina}@techfak.uni-bielefeld.de

## ABSTRACT

*A major problem with (partially) free constituent order is to manifest preferences among structurally distinct parses of ambiguous sentences. In order to obtain scoring criteria a preferred constituent order can considerably support a best-first strategy. This work presents an experimentally evaluated model of preferred German constituent order in the middle field and its application for the implementation of a robust and efficient parsing strategy for spontaneous speech. This constituent order is used to guide the parallel LR(1)-parser to derive the preferred interpretation of ambiguous sentences.*

## 1  INTRODUCTION

Although German is a partially free word order language native speakers show a definite preference for a certain constituent order. An experimental evaluation of the influence of the grammatical function on the constituent order established a close relation between both. For the considered constituents a most frequently used order basing on their grammatical functions is obtained. This preferred constituent order is used to manifest preferences among structurally distinct parses of ambiguous sentences. These preferences are used to guide a LR(1)-parser to derive the preferred parse for ambiguous sentences. The model for the disambiguation of conflicting parses is based on the assumption that the parse closest to the preferred constituent order is assigned the preferred interpretation. Therefore, the action of the LR(1)-parser which leads to a parse tree most similar to the preferred constituent order will be judged best. All parses will be parallel derived from the LR(1)-parser but the one which is judged best will be given preference.

## 2  THE EVALUATED MODEL

### 2.1  The Model

Assuming that there are several factors controlling the linearization of constituents we focus on the influence of the grammatical function on the constituent order. Our model - based on Heidolph, Flaemig & Motsch's [1] work - aims at definite, non-pronominal and non-modified constituents of the middle field of isolated verb-second sentences. We call a non-finite part of a sentence a constituent if it can be topicalized and if there is no further segmentation of this part so that each segment can be topicalized for itself without causing a semantic difference [7]. For such constituents having one of the following grammatical functions our model predicts as the most frequently realized order:

subject < temporal adverbial < locational adverbial < modal or instrumental adverbial < dative object < accusative object < absolute directional adverbial < relative directional adverbial < non-finite part of the predicate (separable prefix, infinite or participle), where '<' stands for 'is realized before'.

### 2.2  Experimental Evaluation

Our model was tested in a sentence generation experiment with a stimuli-set of 41 sentences with the constant length of five constituents [7]. Each of the 49 subjects had to alternately verbalise a verb-second sentence and - as a distracting item - an arithmetical equation. In both cases the constituents were visually presented simultaneously. For the statistical analysis of the sentence realizations two meassuremants are mainly determined. First there is the binary measurement whether subjects follow the predicted order with 1/6 'right' of 24 possible orderings - assuming each non-finite constituent can be topicalized. Then there is the Levenshtein-distance L indicating the degree of deviation from the predicted order. The results of a statistical analysis containing both measurements support our model. More than 71% of the sentences were realized following the preferred order. Most of the subjects were even realized as the first constituent (81%, that is 58% of all sentences). A fraction from the remaining 29% does not have to be analysed as a deviation because in these cases locational adverbials were not used as a constituent but as a modification of a constituent. In case of non-preferred ordering the degree of deviation is rather small (L<0,5 on a scale from 0 to 3). Additionally, we found a second preferred order not predicted by

the model where the subject is in sentence initial position and where the first two constituents of the middle field occur in switched order (14%). Results of post hoc calculated standard residues of tests against the hypothesis that each of the constituent orders is equally probable strongly support the preference of the predicted constituent order with the subject in sentence initial position.

The encouraging empirical results go together with results from two differently but nonetheless comparably designed sentence generation experiments from Pechmann, Uszkoreit, Engelkamp & Zerbst [3] looking at constituent order. Their method correlates errors and articulation times with the grammaticality and acceptability judgments of given constituent orders. Contrary our experiment dealt with spontaneously chosen orderings and it supplementary covers an expanded set of constituents.

## 3 THE PARSING MODEL

### 3.1 The Parser

The processing model for the automatic determination of the constituent order is based on a LR(1)-parser and a unification grammar. The advantage of this approach is that unification grammar puts most of the syntactic information that is standardly captured in context free phrase structure rules into the lexicon. Every word in the lexicon is represented by a feature structure which specifies the values for various attributes [5]. The parser is augmented with the facility for handling feature structures [2]. This implies that every reduce action of the parser includes a unification of the corresponding feature structures. Ambiguous hypotheses or ambiguous syntactic structures are derived by a parallel processing strategy of the LR(1)-parser. Therefore, the action table of the LR(1)-parser contains conflicting entries for ambiguous grammar rules and multiple parse trees are built for ambiguous sentences. Competitive parses are necessary because the kind of the constituent under consideration cannot be determined before it is completely derived. The preferred constituent order is used to guide the parallel LR(1)-parser to manifest preferences among structurally distinct parses of ambiguous sentences. Accordingly, the parse tree with a constituent order most similar to the preferred constituent order will be judged best.

### 3.2 The Determination of the Grammatical Function of Constituents

The information for the automatic determination of the constituent order which is implied by the grammatical function of the constituents is given by the corresponding feature structures. This determination of the grammatical function is archieved by the following schemata:

1. The entries of the subcategorization frame of a verb determine the kind of oblique constituents. By unifying a noun phrase with a verb the head value CASE of the noun determines the kind of the noun phrase.

2. For the determination of the kind of prepositional phrase the head value PFORM (prepositional form) of the preposition and the head value MAJOR of the noun or adverbial phrase is used. Additionally, the CASE value distinguishes certain kinds of prepositional phrases formed by a preposition and a noun phrase.

3. In order to distinguish temporal prepositional phrases from other ones, which cannot be done using only the information given by the CASE value, the head value NFORM (noun form) is extended by a value t-noun, which distinguishes temporal nouns from other ones.

4. Adverbials are determined by the head values MAJOR and ADVFORM (adverb form), which is extanded by the values temporal, locational, directional and modal.

By the unification of the described values most prepositional phrases can be uniquely determined.

### 3.3 The Processing Model

As described in section (4.3) the decision which parse will obtain preference is delayed until a prepositional phrase is reduced. By this reduction the kind of prepositional phrase will be determined as described in section (3.2). The constituent order is implied by the grammatical functions of the constituents. Each grammatical function is symbolized by a number, called the ordering number. The order of these numbers is defined to be increasing. So, the constituent order is the following: subject (1) < temporal adverbial (2) < locational adverbial (3) < modal or instrumental adverbial (4) < dative object(5) < accusative object (6) < absolute directional adverbial (7) < relative directional adverbial (8) < non-finite part of the predicate (separable prefix, infinite or participle) (9). In case of a conflicting configuration of the LR(1)-parser two different orders of the constituents are derived. The parse with decreasing order will be judged worse.

## 4 DISAMBIGUATION RULES

As described in section (2.1) a preferred constituent order is confirmed by the experimental evaluation. As described in [6] and [4] conflicting configurations, i.e. conflicting entries in the action table of the LR(1)-parser, can be solved by preferring one possible action. Our processing model proceeds on the assump-

tion that the preferred action in an ambiguous configuration is predicted by the preferred, i.e. increasing, order.

## 4.1 Lexical Disambiguation

Before presenting the processing model all cases where the attachment is unambiguous are listed. These cases must be excluded because no disambiguation is necessary and therefore our model shall not make any prediction in these cases.

**Example 1:**

- Du schraubst die Leiste auf den Würfel fest.
- *You screw the bar on the cube tight.*
- You tighten the bar on the cube.

In example (1) the noun phrase *'the cube'* of the prepositional phrase *'on the cube'* has accusative case. Therefore, the prepositional phrase is a direct directional adverbial phrase which can only modify the verb. To exclude the attachment to the noun *bar* prepositional phrases which can only modify a noun or a verb respectively will assigned categories which can only reduce with one of them. This choice of category is delayed until the determination of the grammatical function of a constituent [6] because no disambiguation is possible before.

**Example 2:**

- Du schickst einen Brief aus Amerika.
- You send a letter from America.

**Example 3:**

- Du gibst Peter einen Brief aus Amerika.
- *You give Peter a letter from America.*
- You give a letter to Peter from America.

In example (2) the prepositional phrase *'from America'* can be a modifier of the verb *'send'*. In example (3) the verb *'give'* cannot be modified by this prepositional phrase. Both verbs will be assigned different categories which can or cannot reduce with the prepositional phrase.

## 4.2 The Phenomena to be Modeled

The remaining ambiguous cases will be derived by the use of the preferred constituent order if possible.

**Example 4:**

- Der Mann (1) schickt dem Freund (5) das Buch (6) aus Amerika (7).
- *The man sends the friend the book from America.*
- The man sends the book to the friend from America.

**Example 5:**

- Der Mann (1) schickt dem Freund (5) aus Amerika (7) das Buch (6).

- *The man sends the friend from America the book.*
- The man sends the book to the friend from America.

In example (4) the order of the constituents in the middle field follows the preferred constituent order. Therefore, no prediction of the attachment of the prepositional phrase *from America* can be made. In example (5) the preferred constituent order in the middle field is violated if the prepositional phrase is considered to be a constituent. By assuming that the prepositional phrase modifies the noun *'friend'* the modified noun phrase *'the friend from America'* will inherit the ordering number (5) from the noun phrase. This implies that the order is still increasing. Therefore, this parse will be preferred to the parse where the prepositional phrase is attached to the verb.

**Example 6:**

- Der Mann (1) schraubt die Leiste (6) auf dem Tisch (3) fest (9).
- *The man screwss the bar on the table tight.*
- The man tighten the bar on the table.

In example (6) the preferred order of the constituents in the middle field is violated if the prepositional phrase is considered to be a constituent. Therefore, the prepositional phrase *on the table* will be considered as a modifier of the noun *bar* which leads to an increasing order. The parse where the prepositional phrase is attached to the verb will be judged worse.

## 4.3 The Processing Model

Summarizing the disambiguation model follows the underlying ideas:

- If the constituents in the middle field follow the preferred constituent order no prediction of the attachment of prepositional phrases can be done.
- If the preferred order of the constituents in the middle field is violated by a prepositional phrase the parse where the prepositional phrase is considered as a modifier of the preceding noun - if admissible (compare section (4.1)) - is preferred.

This will be modeled by the LR(1)-parser by the following disambiguation rules:

- If a consituent will be reduced with the verb phrase the ordering number is obtained from the right constituent of the corresponding grammar rule.[1]

---

[1] This exludes automatically constituents in the front field, which are not considered by our model.

- The judgment of a reduce action will be decreased if two constituents are reduced where the ordering numbers of both are not in an increasing order.
- Otherwise the judgment remains unchanged.

The disambiguation model is illustrated on example (6). For the processing model the finite part of the verb is included in the constituent order and it is assigned the ordering number (0). The prepositional phrase *'on the table'* can modify the verb or the noun *'bar'*. In the first case the noun phrase *'the bar'* is reduced with the verb. Since the noun phrase is an accusative object it obtains the ordering number (6). After this reduction the verb phrase is also assigned the ordering number (6). In the first case the prepositional phrase is a locational adverbial phrase and obtains the ordering number (3). By reducing the verb phrase with the prepositional phrase the increasing constituent order is violated and the judgment will be reduced. In the second case the noun *'bar'* will be first reduced with the prepositional phrase and no ordering number can be determined due to schema 1 in section (3.2). The modified noun phrase *'the bar on the table'* will be reduced with the verb. By unifying the verb and this modified noun phrase the noun phrase is assigned the ordering number (6). Now the order is increasing and the judgment will be maintained. This implies that the final judgment of the second case is better than the one of the first case and so this interpretation will be preferred.

In several cases modal or instrumental adverbials can occur at a different position in a sentence as defined in our model. This implies that the judgment of the attachment of modal or instrumental adverbials is rather insecure. Because our parser is designed to derive utterances of spoken dialogues in a special assembly scenario we can eliminate this kind of uncertainties by including visual information [8] into the analysis process. The visual analysis component determines which kind of objects are in a scene. Basing on this judgment the corresponding interpretation will be chosen although it might be that an other interpretation will be judged better by our model. Modal or instrumental adverbials and local adverbials with the preposition 'auf' (*on*) cannot be distinguished by our model. Both possibilities will be derived by our model. Visual information is used as far as possible to distinguish both cases.

## 5 CONCLUSION

The model presented here describes an attachment preference resolution based on a preferred constituent order. Although, German is a partially free word order language native speakers show a definite and mostly consistent preference for a certain constituent order. This order is taken to guide a par-

allel LR(1)-parser to derive the preferred reading of structurally distinct parses of ambiguous sentences. Conflicting configurations will be solved based on this preferred order. Multiple parse trees are built for ambiguous sentences but preference will be given to the one most similar to the preferred constituent order, i.e. the parse with an increasing order of the corresponding ordering numbers.

## 6 ACKNOWLEDGEMENT

## References

[1] K. E. Heidolph, W. Flaemig, and W. Motsch. *Grundzüge einer deutschen Grammatik*. Berlin: Akademie-Verlag, 1984.

[2] S. Kronenberg and F. Kummert. Soft unification: Towards robust parsing of spontaneous speech. In *IASTED International Conference on Artificial Intelligence and Soft Computing, 9.- 12. August 1999*, Honolulu,USA, 1999. to appear.

[3] T. Pechmann, H. Uszkoreit, J. Engelkamp, and D. Zerbst. Wortstellung im Deutschen Mittelfeld. Linguistische Theorie und psycholinguistische Evidenz. In C. Habel, editor, *Perspektiven der kognitiven linguistischen Modelle und Methoden*, pages 257–299. Opladen: Westdeutscher Verlag, 1996.

[4] Fernando Pereira. A new characterization of attachment preference. In D. Dowty, L. Kartunnen, and A. Zwicky, editors, *Natural Language Parsing: Psychological, Computational and Theoreticla Perspectives*, ACL Studies in Natural Language Processing, pages 307–319. Cambridge University Press, 1985.

[5] Carl Pollard and Ivan A. Sag. *Information-Based Syntax and Semantics*, volume 1: Fundamentals of *CSLI Lecture Notes no. 13*. Stanford: Center for the Study of Language and Information, 1987.

[6] Stuart Sieber. Sentence disambiguation by a shift-reduce parsing technique. In *Proceedings of the 21st Annual Meeting of the Association for Computational Linguistics*, pages 113–118, 1983.

[7] K. Skuplik and R. Flach. Präferierte Satzgliedfolge im Mittelfeld: Modell und experimentelle Evaluation. Technical report, 98-9, Universität Bielefeld, SFB 360, 1998.

[8] G. Socher, G. Sagerer, and P. Perona. Bayesian reasoning on qualitative descriptions from images and speech. In H. Buxton and A. Mukerjee, editors, *ICCV'98 Workshop on Conceptual Description of Images*, Bombay, India, 1998.