

## TIGER: Linguistic Interpretation of a German Corpus

SABINE BRANTS<sup>1</sup>, STEFANIE DIPPER<sup>2</sup>, PETER EISENBERG<sup>3</sup>,  
SILVIA HANSEN-SCHIRRA<sup>1</sup>, ESTHER KÖNIG<sup>2</sup>, WOLFGANG  
LEZIUS<sup>2</sup>, CHRISTIAN ROHRER<sup>2</sup>, GEORGE SMITH<sup>3</sup> and HANS  
USZKOREIT<sup>1</sup>

<sup>1</sup>*Computational Linguistics, Saarland University (E-mail: Sabineb@9mail.com);*

<sup>2</sup>*Institute of Natural Language Processing (IMS), Stuttgart University;* <sup>3</sup>*Institut für Germanistik, Potsdam University*

**Abstract.** This paper reports on the TIGER Treebank, a corpus of currently 40,000 syntactically annotated German newspaper sentences. We describe what kind of information is encoded in the treebank and introduce the different representation formats that are used for the annotation and exploitation of the treebank. We explain the different methods used for the annotation: interactive annotation, using the tool ANNOTATE, and LFG parsing. Furthermore, we give an account of the annotation scheme used for the TIGER treebank. This scheme is an extended and improved version of the NEGRA annotation scheme and we illustrate in detail the linguistic extensions that were made concerning the annotation in the TIGER project. The main differences are concerned with coordination, verb-subcategorization, expletives as well as proper nouns. In addition, the paper also presents the query tool TIGERSearch that was developed in the project to exploit the treebank in an adequate way. We describe the query language which was designed to facilitate a simple formulation of complex queries; furthermore, we shortly introduce TIGER in, a graphical user interface for query input. The paper concludes with a summary and some directions for future work.

**Key words:** annotation, German, treebank

### 1. Introduction

Corpus-based methods play an important role in empirical linguistics as well as in machine learning methods in NLP. In these two areas of research, large natural language corpora, enriched with syntactic information, are needed. Thus, in recent years, there has been an increasing interest in the construction of these syntactically annotated corpora, commonly called *treebanks* (Lezius, 2001).

For German, the first initiative in the field of treebanks was the NEGRA Corpus (cf. (Skut et al., 1998; Brants et al., 1999a)), which contains

syntactically interpreted newspaper texts. Furthermore, there is the Verbmobil Corpus (Wahlster, 2000), which covers the area of spoken language.

This paper reports on the TIGER Treebank project, which aims at building the largest and most exhaustively annotated treebank for German. The annotation format and scheme are based on the NEGRA corpus; however, the TIGER Treebank exceeds the NEGRA corpus in size as well as in detail of annotation. Since the NEGRA Corpus is rather restricted in its size (20,000 syntactically annotated sentences) and the Verbmobil Corpus in its domains (i.e. spontaneous speech for the appointment negotiation domain), the construction of the TIGER Treebank as a comprehensive resource for the German language was a necessary step to overcome these drawbacks.

This paper is structured in the following way: section 2 describes the annotation format and provides general information on the annotation scheme. Furthermore, it contains a short overview of treebank initiatives for languages other than German. In section 3, the different methods used for the annotation of the treebank are presented. The linguistic extensions that were made in the TIGER project concerning the annotation scheme are covered in section 4. Section 5 gives an overview of the query language and query tool that were developed in the project for the exploitation of the treebank. Finally, section 6 summarizes the paper and sketches some ideas for future work.

## 2. The TIGER Corpus

The basis of the TIGER Treebank are texts from the German newspaper *Frankfurter Rundschau*. Only complete articles were used, which were taken from all kinds of domains<sup>1</sup> so as to cover a broader range of language variation. At the current stage of the first release in July 2003, the corpus contains 40,000 syntactically annotated sentences (i.e. about 800,000 words). For the end of the project, this amount is to be extended to approximately 80,000 sentences (about 1,500,000 words).

### 2.1. LEVELS OF ANNOTATION

In the NEGRA as well as in the TIGER corpus, a hybrid framework is used which combines advantages of dependency grammar and phrase structure grammar. The syntactic structure is represented by a tree. The branches of a tree may cross, allowing the encoding of local and non-local dependencies and eliminating the need for traces. This approach has considerable advantages for free-word order languages such as German, which show a large variety of discontinuous constituency types (Skut et al., 1997).

The linguistic annotation of each sentence in the TIGER Treebank is represented on a number of different levels (see Figure 1): Part-of-speech information is encoded in terminal nodes (on the word level). Non-terminal nodes are labelled with phrase categories. The edges of a tree represent syntactic functions. Furthermore, a supplementary annotation on the word level facilitates the encoding of information on lemmata and morphology.<sup>2</sup> For part-of-speech tagging, the Stuttgart-Tübingen-Tagset (Schiller et al., 1999) is used in a slightly modified version (cf. (Kramp and Preis, 2000; Smith and Eisenberg, 2000)). Information on lemmata and morphology was not annotated in the NEGRA corpus; this is a new feature that was added to the annotation in the TIGER project.

Syntactic structures are rather flat and simple in order to reduce the potential for attachment ambiguities. The distinction between arguments and adjuncts, for instance, is not expressed in the constituent structure, but is instead encoded by means of syntactic functions.

Apart from the annotation of morphology and lemmata, another annotation level was added to the TIGER corpus: Secondary edges, i.e. labelled directed arcs between arbitrary nodes, are used to encode coordination information. Currently, these secondary edges are only employed for the annotation of coordinated sentences and verb phrases; another potential use might be the systematic annotation of attachment ambiguities.

## 2.2. ANNOTATION FORMATS

In the TIGER project, we use several annotation formats for corpus storage, export and querying. There exist scripts that enable the transformation from one format to another.

First of all, the annotated sentences are stored and maintained in a MySQL database; information about the annotation is contained in tables. An additional output format is used for the export of the sentences. The

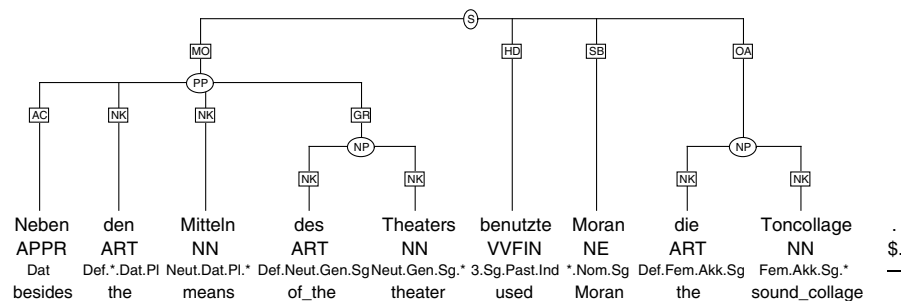


Figure 1. Different levels of annotation.

```

#BOS 37 3 863207489 1
%%word      tag      morph      edge      parent
Ausgerechnet ADJD    --         MO        502
Iggy        NE      Masc.Nom.Sg PNC       500
Pop         NE      *.Nom.Sg   PNC       500
verkoerpert VVFIN  3.Sg.Pres.Ind HD        503
gesanglich ADJD    Pos        MO        503
den         ART     Def.Masc.Akk.Sg NK        501
Staatsanwalt NN      Masc.Akk.Sg.* NK        501
.           $.      --         --        0
#500        MPN    --         NK        502
#501        NP     --         OA        503
#502        NP     --         SB        503
#503        S      --         --        0
#EOS 37

```

Figure 2. TIGER export format.

database entries (words, morphological tags, terminal nodes, non-terminal nodes and edges) can be exported to a table stored in a line-oriented and ASCII-based format (Brants, 1997) (see Figure 2). The major advantage of this export format is that it is easily readable for humans as well as easily processable for machines. Sentence boundaries are identified through sentence start and end tags. Furthermore, information on sentence origins, editors and tags used is stored at the beginning of each export file.

Based on this export format, the TIGER corpus can be transferred into a third format, namely TIGER-XML (Lezius et al., 2002a) (see Figure 3). A TIGER-XML file is typically split up into header and body. The corpus header contains meta-information on the corpus (such as corpus name, date, author, etc.) and a declaration of the tags that are used for morphology, part-of-speech, non-terminal nodes and edges. In the corpus body, directed acyclic graphs are used as the underlying data model to encode the linguistic annotation. Words, part-of-speech tags, morphological tags and lemmata occur as attributes of the element ‘terminal’, whereas non-terminals are represented in an additional element called ‘non-terminal’ referring to the corresponding terminal ID. Secondary edges are encoded explicitly as well. By using an XML format, the TIGER Treebank is exchangeable and usable with a large range of tools. The XML format is also the basis for the use of the corpus query tool TIGERSearch (Lezius and König, 2000).

### 2.3. COMPARABLE TREEBANK INITIATIVES

One of the first and best known treebanks is the Penn Treebank for the English language (Marcus et al., 1994), which consists of about 1 million

```

<s id="s37">
<graph root="s37_503">
<terminals>
  <t id="s37_1" word="Ausgerechnet" pos="ADJD"
    morph="--" />
  <t id="s37_2" word="Iggy" pos="NE"
    morph="Masc.Nom.Sg" />
  <t id="s37_3" word="Pop" pos="NE"
    morph="*.Nom.Sg" />
  <t id="s37_4" word="verkoerpert" pos="VVFIN"
    morph="3.Sg.Pres.Ind" />
  <t id="s37_5" word="gesanglich" pos="ADJD"
    morph="Pos" />
  ...
</terminals>
<nonterminals>
  <nt id="s37_500" cat="MPN">
    <edge label="PNC" idref="s37_2"/>
    <edge label="PNC" idref="s37_3"/>
  </nt>
  <nt id="s37_501" cat="NP">
    <edge label="NK" idref="s37_6"/>
    <edge label="NK" idref="s37_7"/>
  </nt>
  ...
</nonterminals>
</graph>
</s>

```

Figure 3. TIGER-XML format.

words of newspaper text. It contains part-of-speech tagging and rough syntactic and semantic annotation. A bracketing format is used to encode predicate-argument structure and trace-filler mechanisms are used to represent discontinuous phenomena. Other comparable treebanks for English are, for instance, the Susanne Corpus (Sampson, 1995) (containing detailed part-of-speech information and phrase structure annotation), the Lancaster Parsed Corpus (Leech, 1992) (representing phrase structure annotation by means of labelled bracketing) and the British part of the International Corpus of English (Greenbaum, 1996) (about 1 million words of British English that were tagged, parsed and checked afterwards).

For languages other than English, a fairly well-known treebank is the Prague Dependency Treebank for Czech (Hajic, 1999). It contains about 450,000 tokens and is annotated on three levels: on the morphological level (tags, lemmata, word forms), on the syntactic level (using dependency

syntax) and on the tectogrammatical level (encoding functions such as Actor, Patient, etc.). Recently, treebank projects for other languages have come to life as well, e.g. for French (Abeillé et al., 2000b), Italian (Bosco et al., 2000), Spanish (Moreno et al., 2000), Dutch (Schuurman et al., 2003), Turkish (Oflazer et al., 1999), Russian (Boguslavsky et al., 2000) and Bulgarian (Simov et al., 2002). More initiatives for linguistically interpreted corpora can be found in Uszkoreit et al. (1999), Abeillé et al. (2000b) and Abeillé et al. (2003).

### 3. Annotation Methods

We use two different methods for the syntactic annotation of the TIGER corpus: Interactive annotation and LFG parsing. The first is a combination of probabilistic parsing and human intervention (section 3.1). After the parsing is completed, morphological annotation is performed semi-automatically, using the given syntactic annotation for disambiguation. For the second method, a symbolic LFG grammar is used to parse large parts of the corpus; the output is disambiguated by a human annotator (section 3.2).

#### 3.1. INTERACTIVE TAGGING AND PARSING

Interactive annotation is an efficient combination of automatic parsing and human annotation. Instead of having an automatic parser as preprocessor and a human annotator as postprocessor, the two steps are interwoven in our approach. The parser generates a small part of the annotation, which is immediately presented visually to the human annotator, who can either accept, correct or reject it. Based on the annotator's decision, the parser proposes the next part of the annotation, which is again submitted to the annotator's judgement. This process is repeated until the annotation of the sentence is complete.

The advantage of this interactive method is that the human decisions can be used by the automatic parser. Thus, errors made by the automatic parser at lower levels are corrected instantly and do not 'shine through' on higher levels. The chances grow that the automatic parser proposes correct analyses on higher levels.

The interactive annotation works on several layers. The lowest one is the part-of-speech layer. Higher layers are defined by the depth of the syntactic structure. Each layer is represented by a different Markov Model, hence the name *Cascaded Markov Models* (Brants, 1999). The first step in the annotation process is the generation of part-of-speech tags. This step is performed using the statistical tagger TnT (Brants, 2000a). In addition to the tags, TnT also generates probabilities that help to decide on the reliability

of a proposed tag. The lower the probability of alternative tags, the higher the reliability of the best tag for a word (Brants and Skut, 1998). Approximately 84% of all tag assignments are classified as reliable by the tagger. The remaining 16% need to be proof-read by human annotators.

Once the part-of-speech tagging is done, Markov Models for higher layers start processing. Hypothetical phrases are generated, and the one with the highest probability is displayed to the annotator. The structure can be accepted, rejected or manually corrected by the annotator. Intervention by the human annotator immediately changes the set of hypotheses used by the parser. The syntactic structure is built phrase by phrase, bottom up. About 71% of the phrases suggested by the parser are correct, 17% need minor intervention (i.e., at most one non-terminal node needs to be added or deleted). The remaining 12% require major intervention by the human annotator.

Both tagger and parser are entirely trained on previously (manually) annotated data. No manual grammar or lexicon development are necessary. The annotation scheme is learnt automatically by tagger and parser. In case of changes in the annotation scheme, only a small amount of data needs to be changed manually. The tagger and parser are then trained on the changed data and are immediately ready for annotation with the new scheme.

The annotation is performed with the help of the tool ANNOTATE (Figure 4), a graphical user interface with a comprehensive set of tree manipulation functions and database access (Plaehn and Brants, 2000). ANNOTATE runs the TnT tagger and the Cascaded Markov Models in the background.

In order to achieve a high level of consistency and to avoid mistakes, we use a very thorough approach to the annotation: First, each sentence is annotated independently by two annotators. With the support of scripts, they then compare their annotations and correct obvious mistakes. Remaining differences are submitted to a discussion between the annotators. Although this process is rather time-consuming, it has proven to be highly beneficial for the accuracy of the annotation Brants (2000b). Furthermore, it also supports the continuous improvement of the annotation scheme: It is in the discussion between the annotators that discrepancies between the annotation scheme and the data become obvious. If this happens to be the case, new rules and better tests for operationalization are added to the annotation scheme. Thus, there is a cross-fertilization between the corpus and the annotation scheme.

For the analysis of lemmata and morphological tags, we use a tool called TIGERMorph which was developed by Berthold Crysmann. This morphological analyser is interleaved with ANNOTATE. TIGERMorph disambiguates the output of a morphological analyser on the basis of the

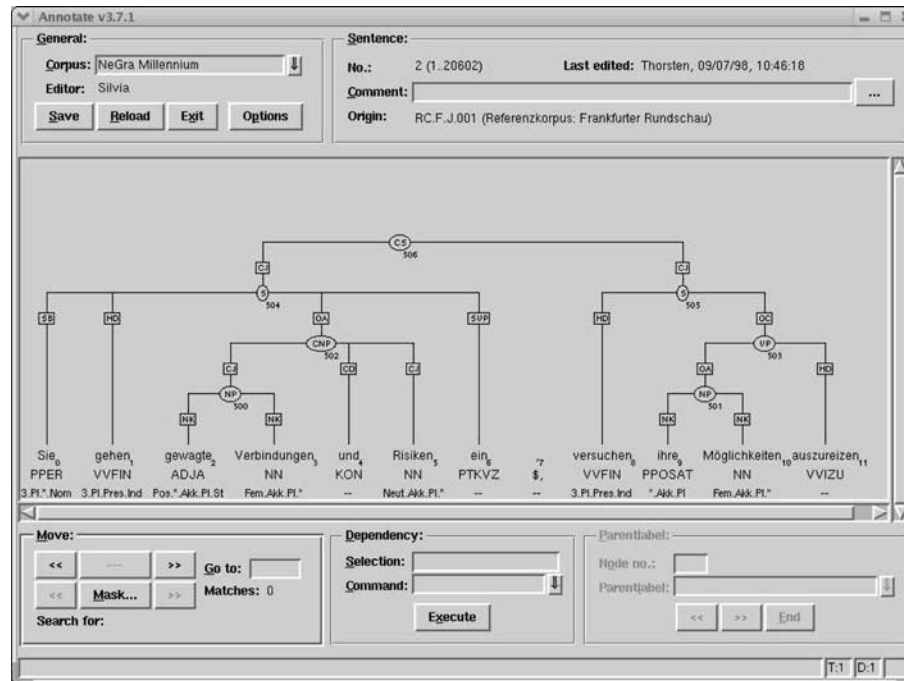


Figure 4. The annotation tool ANNOTATE.

already existing syntactic structure. It proposes lemmata and morphological tags for the words of a sentence, proceeding from left to right. The annotation of morphology and lemmata resembles the interactive annotation of the syntactic structure described above. Ambiguous tags are presented to the annotator, who then decides which one is the appropriate alternative. This information is returned to TIGERMorph and used for the disambiguation of the morphological analysis for the remaining words.

### 3.2. ANNOTATION BY LFG PARSING

As an alternative to interactive tagging and parsing, a broad coverage symbolic LFG grammar (Lexical Functional Grammar (Bresnan, 1982)) is used to parse parts of the corpus (Dipper, 2000). Usually, the LFG grammar outputs several analyses for a corpus sentence. The output is first filtered by a grammar internal ranking mechanism and then disambiguated by a human annotator. A transfer component converts the selected analysis into the TIGER format (Zinsmeister et al., 2002).

One advantage of this approach is the accuracy of the grammar's output. An LFG analysis is always syntactically consistent. It does not contain inconsistencies such as, e.g., missing subject-verb agreement, in case of



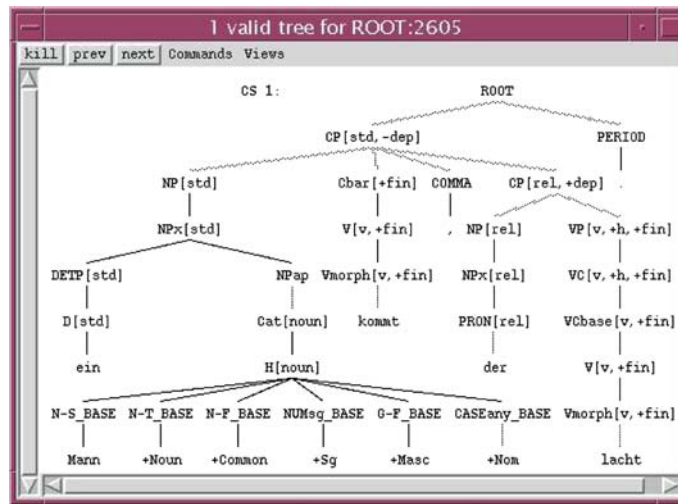


Figure 5. LFG c-structure.

which the parse would have failed. On the other hand, the grammar certainly is not error-free. But those errors which do occur are systematic and hence easier to correct than errors that occur with manual annotation.

### 3.2.1. Parsing

The LFG grammar applied in parsing has been developed in the ParGram project at the University of Stuttgart, using the Xerox Linguistic Environment (XLE) (ParGram, 2002; Dipper, 2003). The output of an LFG grammar basically consists of two representations, the constituent structure (c-structure) of the sentence being parsed, and its functional structure (f-structure). C-structure encodes information about morphology, constituency, and linear ordering. F-structure represents information about predicate argument structure, about modification, and about tense, mood etc. Examples of c- and f-structures are given in Figures 5 and 6 for the sentence *Ein Mann kommt, der lacht* ('a man is coming who laughs').

### 3.2.2. Disambiguation

Almost every sentence of a newspaper corpus is syntactically ambiguous. Hence the output of a purely symbolic grammar has to be disambiguated, i.e. a human annotator has to select the appropriate analysis. This task is supported by XLE which allows for 'packing' the different readings into one complex f-structure representation.<sup>3</sup>

On average, however, a sentence of the TIGER corpus receives several thousands of LFG analyses. Clearly it is impossible to disambiguate those

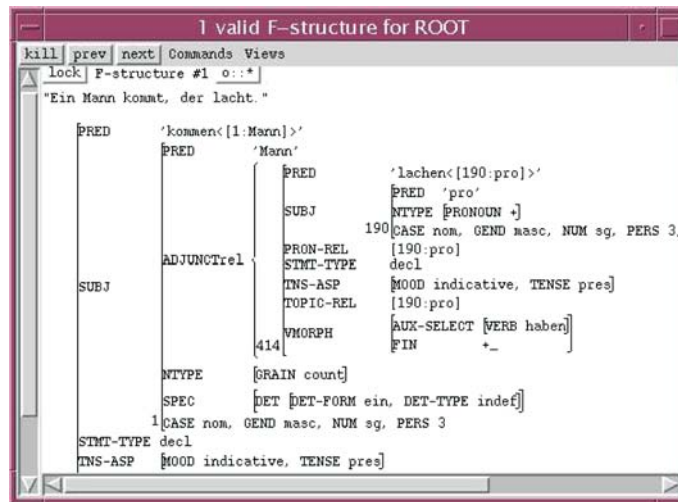


Figure 6. LFG f-structure.

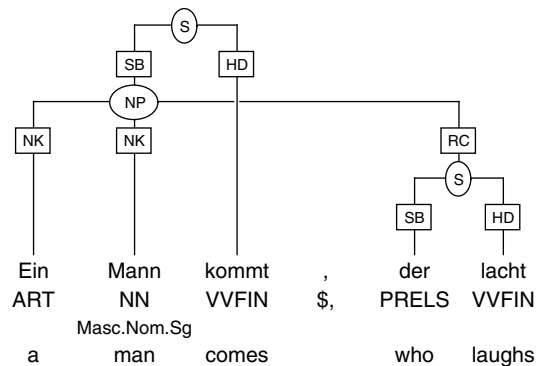


Figure 7. TIGER representation.

analyses manually. Therefore XLE provides a (non-statistical) mechanism for suppressing certain ambiguities automatically (Frank et al., 1998). By means of this mechanism, the average number of solutions drops down to 17, the median being 2.

### 3.2.3. Conversion into TIGER Format

All information that is required by the TIGER annotation scheme is contained in c- and f-structure representations of LFG. Compare the LFG representation (Figures 5 and 6) with the TIGER representation (Figure 7) of the sentence *Ein Mann kommt, der lacht*.

(i) LFG c-structure contains categorial information (e.g., NP, CP), lemmata (Mann), part-of-speech tags (+ Noun + Common), and morphological

tags (+Sg +Masc +Nom); in Figure 5, lemma and tags are only shown for the terminal *Mann*. In the TIGER scheme, this information is encoded in a slightly different terminology: the nodes and tags mentioned above correspond to TIGER nodes NP, S, the part-of-speech tag NN, and the morphological tag Masc.Nom.Sg, respectively.

(ii) LFG f-structure represents dependency relations such as the head argument relation SUBJ and the head modifier relation ADJUNCTrel ('relative clause'). Note that in c-structure, NP and CP (the relative clause) do not form a constituent; however, their f-structures, SUBJ and ADJUNCTrel, are linked. This linking information is encoded by a crossing branch in TIGER, cf. Figure 7.

Often, there is a one-to-one correspondence between LFG and TIGER representations. In these cases the transfer component simply converts one format into another, e.g. +Sg +Masc +Nom is mapped to Masc.Nom.Sg, CP to S, SUBJ to SB, etc. However, in other cases the transfer has to combine information both from c- and f-structure, as in the case of the extraposed relative clause in Figures 5–7. Here the transfer component makes use of the f-structure link between SUBJ and ADJUNCTrel to form a (discontinuous) constituent. The categorial label (NP) is derived from c-structure.

#### 3.2.4. Results and Outlook

When parsed with the current grammar version, 50% of the corpus sentences receive at least one analysis; approx. 70% of the parsed sentences receive the correct analysis (possibly among others).<sup>4</sup> About 2,000 sentences of the TIGER corpus have been annotated this way.

To enlarge coverage, XLE allows for partial parses, providing, e.g., N or P chunks. In first experiments, N chunks were found with a precision of 89% and a recall of 67%; for P chunks, the precision was 96%, and the recall was 79% Schrader (2001). Furthermore, to minimize manual effort, a statistical disambiguation tool can be integrated (Riezler et al., 2002).

## 4. Extensions in the TIGER Annotation Scheme

The annotation in TIGER is based on the annotation scheme that was used for the NEGRA corpus (Brants et al., 1999b). This annotation scheme covered a broad variety of phenomena. However, there was still room for improvement in its linguistic adequacy. A vital part of the work in the TIGER project is the linguistic extension of this annotation scheme. In the following, the major changes that were made in the TIGER project are presented. A more detailed account of these changes and an evaluation of the improved annotation scheme can be found in (Brants and Hansen, 2002).

## 4.1. COORDINATION

An essential linguistic extension in the TIGER annotation scheme was made concerning the annotation of coordinated sentences and verb phrases. In the NEGRA corpus, arguments that are shared by both verb conjuncts of a coordination, but that are only mentioned once, were structurally linked only to the nearest part of the coordination. Thus, the NEGRA annotation is in many cases not suitable for the extraction of subcategorization information. In the TIGER treebank, these shared arguments are provided with secondary edges in order to represent their syntactic relation to the more distant verb conjuncts.

Figures 8 and 9 illustrate the difference between the NEGRA and the TIGER treatment of these cases. In the example sentence *Der Mann liest und zerknüllt die Zeitung* ('the man reads and rumples the newspaper'), the common subject of both verbs is the NP *der Mann*, the common object is the NP *die Zeitung*. However, in the NEGRA annotation scheme, shared arguments are linked only to the nearest verb (cf. Figure 8). The structure of the tree would be exactly the same if the first verb were intransitive and did not have *die Zeitung* as its object (e.g. *Der Mann lacht und zerknüllt die Zeitung* ('the man laughs and rumples the newspaper')). In contrast, Figure 9 shows the annotation of the sentence according to the TIGER annotation scheme, making use of the secondary edges that were introduced. Thus, the TIGER scheme allows the differentiation between transitive and intransitive verbs in coordination.

## 4.2. VERB-SUBCATEGORIZATION

The NEGRA corpus provides no distinctions between prepositional phrases with respect to their syntactic functions; all PPs occurring in sentences

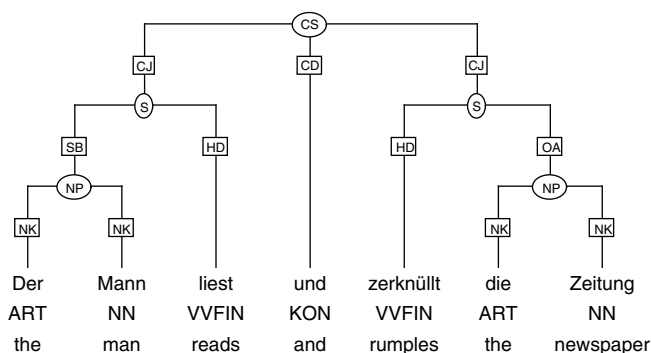


Figure 8. Coordination with shared arguments in NEGRA.

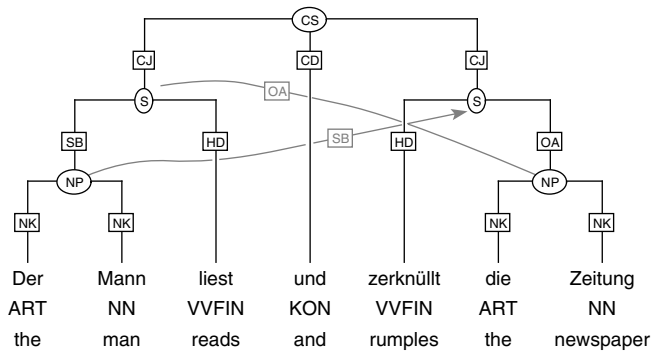


Figure 9. Coordination with shared arguments in TIGER.

or verb phrases are unexceptionally marked with the label MO (modifier). In the TIGER project, two additional function labels for PPs were introduced: prepositional objects (OP) and collocational verb constructions (CVC). The label OP is applied to constructions like *auf jemanden warten* (‘to wait for somebody’). These phenomena are marked by the fact that the preposition *auf* (‘on’) has lost its lexical meaning.

Figure 10 exemplifies the fact that NEGRA did not allow the distinction between complements and adjuncts on the level of edge labels.<sup>5</sup> In contrast, the TIGER annotation (Figure 11) mirrors the functional difference

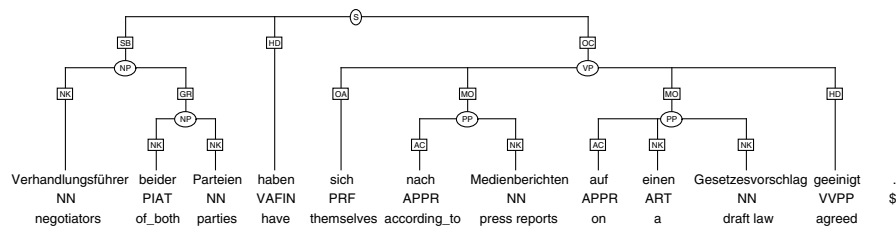


Figure 10. Annotation of PPs in NEGRA: MO.

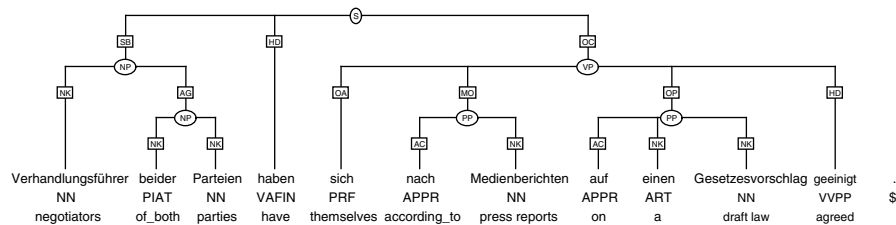


Figure 11. Annotation of PPs in TIGER: MO and OP.

between the first PP and the second PP in the use of different edge labels: the first PP is functionally independent of the verb and serves as an adverbial; it still receives the label MO. The second PP represents a typical example of a prepositional object (OP) in German: the preposition *auf* ('on') has completely lost its lexical meaning and is purely functional.

The other newly introduced edge label for prepositional phrases, CVC (collocational verb construction), serves to label verb + PP constructions in which the main semantic information is contained in the noun of the PP, not in the verb. This label can only be used with a very limited class of verbs (usually a semantically weak verb with an originally directional or local meaning, e.g. *stellen*, *kommen*, etc. ('to put', 'to come')) that occur in connection with an equally limited class of prepositions (mostly *zu* and *in* ('to' and 'in')). A typical example of this is the German collocational expression *in Kraft treten* (literal translation: 'to step into force', meaning: 'to take effect').

#### 4.3. EXPLETIVES

In the TIGER corpus, finer distinctions with regard to the usage of *es*, the German expletive, have been introduced. In the NEGRA annotation scheme, only one label (PH, meaning place-holder) was used for the non-semantic usage of *es*; in the TIGER scheme, we distinguish three types:

- *Vorfeld es*: This type of *es* is used to fill the first position of a sentence, called the *Vorfeld* slot. It is marked with the label PH (place-holder). Example: *Es naht ein Gewitter* (literal translation: 'it approaches a thunderstorm'; meaning: 'a thunderstorm is approaching'). As soon as another component occupies the *Vorfeld* slot, the *es* disappears: *Ein Gewitter naht*.
- *Correlative es*: This second type of *es* is always correlated to some propositional argument in the sentence. It is usually optional. Example: *Mich freut es, dass ...* ('it makes me happy that ...'). This type is also labeled as PH but can be easily distinguished from the first type because it always occurs in connection with a propositional sister node functioning as RE (repeated element) (cf. Figure 12).
- *Expletive es*: The last type of *es* functions as a non-thematic argument, e.g. in connection with weather verbs: *Heute regnet es* ('today, it is raining'). It receives the label EP (expletive).

#### 4.4. PROPER NOUNS

In the NEGRA as well as in the TIGER corpus, the parent label PN is used to mark ordinary proper nouns, such as *Gerhard Schröder*. The

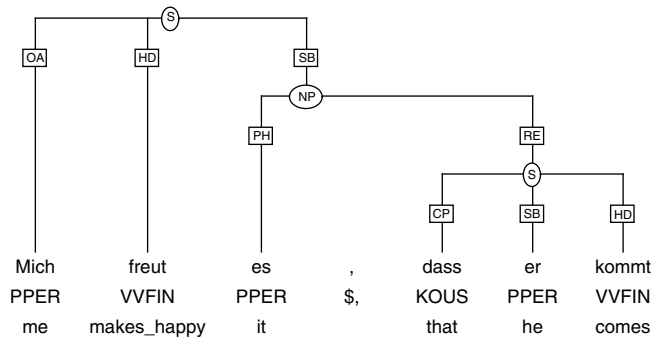


Figure 12. Correlative *es*.

single components receive the edge label PNC (proper noun component). Furthermore, the label PN is also used for multi-token company names, newspaper names (e.g. *The San Francisco Chronicle*) etc. In the TIGER annotation scheme, the usage of this label was extended to cover titles of films, books, exhibitions etc. that have a complex, sometimes sentence-like structure. Occurrences of these phenomena are first annotated structurally and then receive an additional unary parent label PN. The examples in Figures 13 and 14 illustrate the different annotations in NEGRA and TIGER.

Thus, the TIGER annotation permits the identification of structures that function as names, but do not feature the corresponding part-of-speech tag NE (proper noun) in one of their terminal nodes.

### 5. TIGER Search

Syntactically annotated corpora such as the TIGER treebank provide a wealth of information which can only be exploited with an adequate query tool and query language. Thus, a powerful search engine for treebanks has

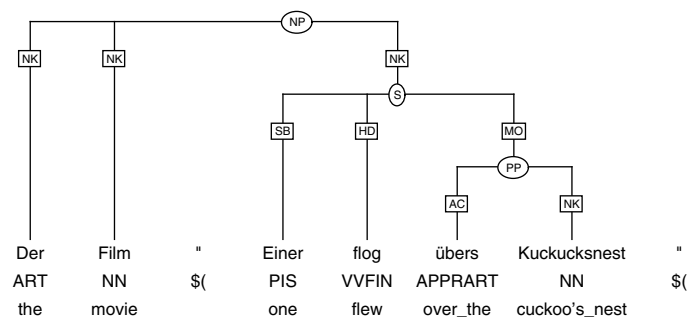


Figure 13. Treatment of structured proper nouns in NEGRA.

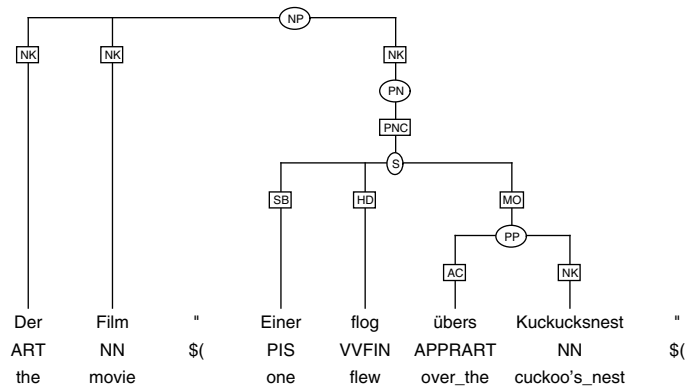


Figure 14. Treatment of structured proper nouns in TIGER.

been developed within the TIGER project (König and Lezius, 2000b; Lezius, 2002).

The search engine is freely available for research purposes and can be downloaded from the TIGER web site. To support all popular platforms, the tool has been implemented in Java. Treebanks have to be encoded according to the TIGER-XML format which is the default import format of TIGERSearch (Mengel and Lezius, 2000; Lezius et al., 2002a). However, input filters (i.e. converters to TIGER-XML) are provided for many popular treebank formats. A complete list of supported formats is given in (Lezius et al., 2002b). The search engine is thus a tool which can be used by the entire treebank community.

### 5.1. QUERY LANGUAGE

The query language (König and Lezius, 2000b; König and Lezius, 2000c) has been designed to fulfil two conflicting requirements: On the one hand, it is close to grammar formalisms, thus easy to learn. It allows modular, understandable code, even for complex queries. A user can pose queries intuitively, mapping linguistic descriptions directly into the query language. On the other hand, its expressiveness has been constrained to guarantee efficient query processing.

The query language consists of three levels. On the *node* level, nodes can be described by Boolean expressions over feature-value pairs. The following query is matched by the terminal node *lacht* ('laughs') in Figure 7:

```
[word="lacht" & pos="VVFIN"]
```

On the *node relation* level, descriptions of two or more nodes are combined by relations. Since graphs are two-dimensional objects, we need one basic relation for each dimension. These are immediate precedence (".")



for the horizontal dimension and immediate dominance (“>”) for the vertical dimension.<sup>6</sup> There are also derived node relations such as underspecified dominance or siblings, e.g.:

- >\* dominance (minimum path length 1)
- >n dominance in  $n$  steps ( $n > 0$ )
- > $m, n$  dominance in  $d$  steps ( $m \leq d \leq n$ )
- >L labelled dominance (edge label L)
- >@1 leftmost terminal successor (‘left corner’)
- . \* precedence (min. number of intervals: 1)
- § siblings

For example, the following query is matched by a subgraph of Figure 7 (the NP node):

```
[cat="NP"] > RC [cat="S"]
```

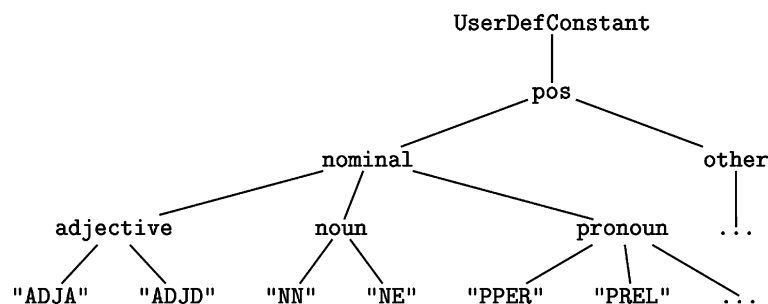
Finally, on the *graph description* level restricted Boolean expressions over node relations can be used (without negation). For example, a subgraph of Figure 7 (the second S node) satisfies the following query:

```
( [cat="S"] > [pos="PRELS"] ) &  
( [cat="S"] > [pos="VFIN"] )
```

Variables can be used to express coreference of nodes or feature values. For example, the two node descriptions [cat="S"] in the above query could refer to different nodes. A reformulation of the query using variables prevents this:

```
(#n: [cat="S"] > [pos="PRELS"] ) &  
(#n > [pos="VFIN"] )
```

In addition, the user can define type hierarchies. Subtypes may also be constants, e.g. in the case of part-of-speech symbols. Here is a part of a type hierarchy for the STTS tag set:



This hierarchy can be used to formulate queries more concisely:

```
[pos=nominal] .* [pos="VVFIN"]
```

There are also several useful predicates such as `discontinuous(#n)`, `continuous(#n)` (phrase does/does not contain crossing branches) or `arity(#n,num)` (phrase comprises *num* children). The following example query determines extraposed relative clauses in the TIGER treebank (cf. Figure 7):

```
(#n:[cat="NP"] >RC [cat="S"]) &
discontinuous(#n)
```

In order to define large collections of queries in a modular way, one can make use of a template notation (König and Lezius, 2000b; König and Lezius, 2000c). The following template describes a clause that comprises a personal pronoun and a finite verb (cf. Figure 9 and Figure 15 for a matching corpus graph):

```
MyClause(#s) <-
  #s:[cat="S"] &
  #t1:[pos="PPER"] &
  #t2:[pos="VVFIN"] &
  (#s > #t1) & (#s > #t2)
```

Now this template can be used in other queries. In the following example, the context of the clause is further specified (cf. example in Figure 15):

```
MyClause(#s) & [cat="CS"] > #s
```

## 5.2. QUERY TOOL

To ensure efficient query processing we have chosen an indexed-based approach. In a preprocessing step a corpus is imported and indexed. Many partial searches are performed during indexing in order to save processing time during query processing. The indexing of a corpus is realized in a tool called TIGERRegistry, the corpus query tool is called TIGERSearch. To increase performance we have also implemented query optimization strategies and search space filters. The query processing strategy is described in detail in (Lezius, 2002).

The TIGERSearch GUI comprises a graph viewer to view the matching sentences of a query. Figure 15 illustrates the visualization of a corpus graph that matches the example template above. Users can browse through the matching sentences using a navigation bar and export their favourite matches. Matches can be exported in the TIGER-XML format, but also as an interactive SVG image. Thus, match forests can be viewed in a format that does not depend on the TIGERSearch software suite.

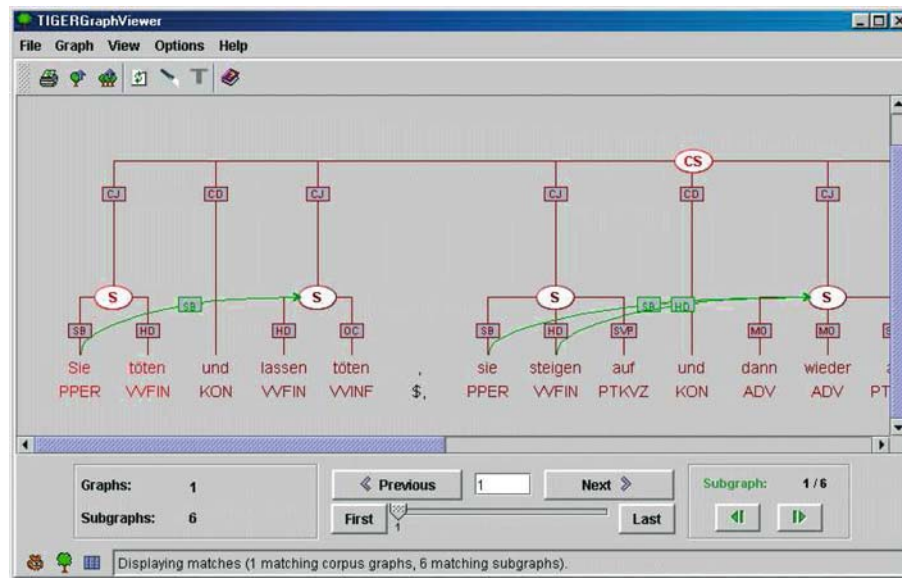


Figure 15. Visualization of query results.

We have also developed a graphical query input front-end which enables users to ‘draw’ queries in a very intuitive way (Voorman, 2002). Queries are expressed by combining nodes and node relations. Figure 16 illustrates how the example query above can be expressed using the graphical query editor.

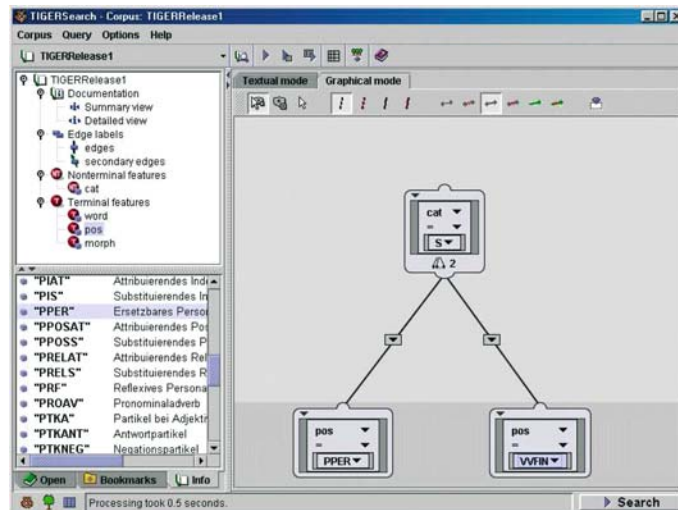


Figure 16. Graphical query input.

## 6. Summary and Outlook

In this paper, we presented the TIGER treebank, the largest and most comprehensive treebank for the German language. We explained the different levels of annotation: part-of-speech tags, phrase categories and syntactic functions. Furthermore, information about lemmata and morphology is also encoded in the corpus. The methods of annotation – interactive annotation with ANNOTATE and LFG parsing – as well as the different representation formats used for the TIGER treebank were demonstrated. We also gave a short overview of related work in comparable treebank projects.

Moreover, we also described the TIGER annotation scheme, which is based on the annotation scheme used for the NEGRA corpus. The paper also outlined the most important extensions in the TIGER annotation scheme, which concern the use of secondary edges in coordination, verb-subcategorization, finer distinctions concerning the German expletive *es* and a different treatment of structured proper nouns.

The last section presented TIGERSearch, a query tool that was developed in the project and which can be used to exploit the TIGER treebank and several other treebank formats. We explained the query language, which was designed to pose intuitive queries to the treebank. We also shortly introduced the graphical query input TIGERin.

Future work will be concerned with additional improvements in the annotation scheme. For instance, we envision the introduction of further distinctions concerning verbal arguments in order to facilitate the identification of thematic roles (Smith, 2000). After the first release of 40,000 sentences in July 2003, a final release is planned for the end of 2003 which will contain about 80,000 sentences completely annotated (part-of-speech tagging, syntactic structure, morphology, lemmata) according to the new TIGER annotation scheme and thoroughly checked for consistency.

All the tools developed in the TIGER project and the first release of the treebank are freely available for research purposes. For further information on the corpus, the corpus tools and how to obtain them, please refer to the project web page: <http://www.coli.uni-sb.de/cl/projects/tiger>

## Notes

<sup>1</sup> We only excluded regional news and sports news because experience from the past showed that these texts often contain tables, enumerations, etc. instead of complete sentences.

<sup>2</sup> The example in Figure 1, however, contains no lemmata annotation, but a literal translation instead.

<sup>3</sup> Furthermore XLE provides various browsing tools applying to c-structure as well as to f-structure which can be used for manual disambiguation (cf. (King et al., To appear) where these tools are described in detail).

<sup>4</sup> About 10% of the sentences failed because of gaps in the morphological analyser; 6% failed because of storage overflow or timeouts (with limits set to 100 MB storage and 100 s. parsing time). But 10% of the parsed sentences were not evaluated wrt. the correct analysis because they received more than 20 analyses.

<sup>5</sup> A correct English translation of the example sentence is the following: 'According to press reports, negotiators from both parties have agreed on a draft law'.

<sup>6</sup> The precedence of two inner nodes is defined as the precedence of their leftmost terminal successors (König and Lezius, 2000a).

## References

- Abeillé A., Brants T., Uszkoreit H., (eds.), (2000a.) *Proceedings of the COLING-2000 Post-Conference Workshop on Linguistically Interpreted Corpora LINC-2000*, Luxembourg.
- Abeillé A., Clement L., Kinyon A. (eds.), (2000b). Building a Treebank for French. *Proceedings of the Second International Conference on Language Resources and Evaluation LREC-2000* pp. 87–94. Athens, Greece.
- Abeillé A., Hansen-Schirra S., Uszkoreit H. (eds.), (2003). *Proceedings of 4th International Workshop on Linguistically Interpreted Corpora (LINC-03)*, Budapest.
- Boguslavsky I., Grigorieva S., Grigoriev N., Kreidlin L., Frid N. (2000). Dependency Treebank for Russian: Concept, Tools, Types of Information. *Proceedings of the 18th International Conference on Computational Linguistics COLING-2000*, Saarbrücken, Germany.
- Bosco C., Lombardo V., Vassallo D., Lesmo L. 2000. Building a Treebank for Italian: A datadriven annotation schema. *Proceedings of the Second International Conference on Language Resources and Evaluation LREC-2000*, pp. 99–106, Athens, Greece.
- Brants T. 1997. *The NEGRA Export Format*. CLAUS Report No. 98. Saarbrücken, Germany: Dept. of Computational Linguistics, Saarland University.
- Brants T. 1999. *Tagging and Parsing with Cascaded Markov Models – automation of corpus annotation*. Saarbrücken, Germany: German Research Center for Artificial Intelligence and Saarland University: Saarbrücken Dissertations in Computational Linguistics and Language Technology, Vol. 6.
- Brants T. (2000a). TnT A Statistical Part-of-Speech Tagger. *Proceedings of the Sixth Conference on Applied Natural Language Processing ANLP-2000*. Seattle, WA.
- Brants T. (2000b). Inter-annotator agreement for a German newspaper corpus. *Proceedings of Second International Conference on Language Resources and Evaluation LREC-2000*. Athens, Greece.
- Brants S., Hansen S. 2002. Developments in the TIGER Annotation Scheme and their Realization in the Corpus. *Proceedings of the Third Conference on Language Resources and Evaluation LREC-2002*, pp. 1643–1649, Las Palmas de Gran Canaria, Spain.
- Brants T., Hendriks R., Kramp S., Krenn B., Preis C., Skut W., Uszkoreit H. 1999b. *Das NEGRA-Annotationsschema*. (Tech. Rep.). Saarbrücken, Germany: Dept. of Computational Linguistics, Saarland University.
- Brants T., Skut W. 1998. Automation of treebank annotation. *Proceedings of New Methods in Language Processing NeMLaP-98*. Sydney, Australia.
- Brants T., Skut W., Uszkoreit H. 1999a. Syntactic Annotation of a German Newspaper Corpus. *Proceedings of the ATALA Treebank Workshop*, pp. 69–76, Paris, France.
- Bresnan J. (ed.), 1982. *The Mental Representation of Grammatical Relations*. MIT Press.

- Dipper S. 2000. Grammar-based Corpus Annotation. *Proceedings of the COLING-2000 Post-Conference Workshop on Linguistically Interpreted Corpora LINC-2000*, pp. 56–64, Luxembourg.
- Dipper S. 2003. *Implementing and Documenting Large-Scale Grammars – German LFG*. Doctoral Dissertation, University of Stuttgart. Arbeitspapiere des Instituts für Maschinelle Sprachverarbeitung (AIMS), Vol. 9(1).
- Frank A., King T. H., Kuhn J., Maxwell J. 1998. Optimality Theory Style Constraint Ranking in Large-scale LFG Grammars. *Proceedings of the LFG-98 Conference*. Brisbane, Australia: CSLI Online Publications, <http://www-csli.stanford.edu/publications>.
- Greenbaum S. (ed.), 1996. *Comparing English worldwide: The International Corpus of English*. Clarendon Press, Oxford, UK.
- Hajic J. 1999. Building a Syntactically Annotated corpus: The Prague Dependency Treebank. Hajicova E., (ed.), *Issues of Valency and meaning. Studies in honour of Jarmila Panevova*. Charles University Press, Prague, Czech Republic.
- King T. H., Dipper S., Frank A., Kuhn J., Maxwell J. To Appear. Ambiguity Management in Grammar Writing. *Journal of Language and Computation*.
- König E., Lezius W. 2000a. A description Language for syntactically annotated corpora. *Proceedings of the 18th International Conference on Computational Linguistics COLING-2000* pp. 1056–1060, Saarbrücken, Germany.
- König, E., Lezius W., Voormann H. 2003. *TIGERSearch User's Manual*. IMS, University of Stuttgart. (<http://www.tigersearch.de>).
- König E., Lezius W. 2003. *The TIGER language – A Description Language for Syntax Graphs, Formal Definition*. Tech. Rep. IMS, University of Stuttgart. (<http://www.ims.uni-stuttgart.de/projekte/complex/paper/lezius/tigerLangForm.ps.gz>).
- Kramp S., Preis C. 2000. *Konventionen für die Verwendung des STTS im NEGRA-Korpus*. Tech. Rep. Saarbrücken, Germany: Department of Computational Linguistics, Saarland University.
- Leech G. 1992. The Lancaster Parsed Corpus. *ICAME Journal*, 16(124).
- Lezius W. 2001. Baumbanken. K.-U. Carstensen, Ebert C., Endriss C., Jekat S., Klabunde R., Langer H. (eds.), *Computerlinguistik und Sprachtechnologie – eine Einführung*, pp. 377–385, Spektrum Akademischer Verlag, Heidelberg, Germany.
- Lezius W. 2002. *Ein Werkzeug zur Suche auf syntaktisch annotierten Textkorpora*. PhD Thesis. IMS, University of Stuttgart.
- Lezius W., Biesinger H., Gerstenberger C. 2002a. *TIGER-XML Quick Reference Guide*. Tech. Rep. IMS, University of Stuttgart.
- Lezius W., Biesinger H., Gerstenberger C. 2002b. *TIGERRegistry Manual*. Tech. Rep. IMS, University of Stuttgart.
- Lezius W., König E. 2000. Towards a Search Engine for Syntactically Annotated corpora. *Proceedings of the Fifth KONVENS Conference*, Ilmenau, Germany.
- Marcus M., Kim G., Marcinkiewicz M., MacIntyre R., Bies A., Gerguson M., Katz K., Schasberger B. 1994. The Penn Treebank: Annotating predicate Argument structure. *Proceedings of the ARPA Human Language Technology Workshop*, Morgan Kaufman, San Francisco, CA.
- Mengel A., Lezius W. 2000. An XML-based encoding format for syntactically annotated corpora. *Proceedings of the Second International Conference on Language Resources and Evaluation LREC-2000*, pp. 121–126, Athens, Greece.
- Moreno A., Grishman R., López S., Sánchez F., Sekine S. 2000. A Treebank of Spanish and its Application to Parsing. *Proceedings of the Second International Conference on Language Resources and Evaluation LREC-2000*, pp. 107–112, Athens, Greece.

- Oflazer K., Hakkani-Tür D., Tür G. 1999. Design for a Turkish treebank. *Proceedings of the Workshop on Linguistically Interpreted Corpora LINC-99*, Bergen, Norway.
- ParGram. 2002. *The ParGram Project*. (URL: <http://www2.parc.com/istl/groups/nltt/pargram/>).
- Plaehn O., Brants T. 2000. ANNOTATE – An Efficient Interactive Annotation tool. *Proceedings of the Sixth Conference on Applied Natural Language Processing ANLP-2000*, Seattle, WA.
- Riezler S., King T. H., Kaplan R., Crouch R., Maxwell J., Johnson M. 2002. Parsing the Wall Street Journal using a Lexical-Functional Grammar and Discriminative Estimation Techniques. *Proceedings of the ACL-02*, Philadelphia, PA.
- Sampson G. 1995. *English for the Computer. The SUSANNE Corpus and Analytic Scheme*. Clarendon Press, Oxford, UK.
- Schiller A., Teufel S., Stöckert, C. 1999. *Guidelines für das Tagging deutscher Textcorpora mit STTS*. Tech. Rep. University of Stuttgart, University of Tübingen.
- Schrader B. 2001. *Modifikation einer deutschen LFG-Grammatik für Partial Parsing*. University of Stuttgart, Studienarbeit.
- Schuurman I., Schoupe M., Hoekstra H., van der Wouden T. 2003. CGN, An Annotated Corpus of Spoken Dutch. *Proceedings of 4th International Workshop on Linguistically Interpreted Corpora (LINC-03)*. Budapest.
- Simov K., Osenova P., Slavcheva M., Kolkovska S., Balabanova E., Doikoff D., Ivanova K., Simov A., Kouylekov M. 2002. Building a Linguistically Interpreted Corpus of Bulgarian: the BulTreeBank. *Proceedings of Third International Conference on Language Resources and Evaluation LREC-2002*, pp.1729–1736. Las Palmas de Gran Canaria, Spain.
- Skut W., Brants T., Krenn B., Uszkoreit H. 1998. A Linguistically Interpreted Corpus of German Newspaper Text. *Proceedings of the Conference on Language Resources and Evaluation LREC-98*, pp. 705–711. Granada, Spain.
- Skut W., Krenn B., Brants T., Uszkoreit H. 1997. An Annotation Scheme for Free Word Order Languages. *Proceedings of the Conference on Applied Natural Language Processing ANLP-97*, Washington, DC.
- Smith G. 2000. Encoding thematic roles via syntactic categories in a German treebank. *Proceedings of the Workshop on Syntactic Annotation of Electronic Corpora*. Tübingen, Germany.
- Smith G., Eisenberg P. 2000. *Kommentare zur Verwendung des STTS im NEGRA-Korpus*. Tech. Rep. University of Potsdam.
- Uszkoreit H., Brants T., Krenn B. (eds.), 1999. *Proceedings of the Workshop on Linguistically Interpreted Corpora LINC-99*. Bergen, Norway.
- Voorman H. 2002. *TIGERin – Graphische Eingabe von Suchanfragen in TIGERSearch*. Diploma Thesis. IMS, University of Stuttgart.
- Wahlster W. (ed.), 2000. *Verbmobil: Foundations of Speech-to-Speech Translation*. Heidelberg, Germany: Springer.
- Zinsmeister H., Kuhn J., Dipper S. 2002. Utilizing LFG Parses for Treebank Annotation. *Proceedings of the LFG-02 Conference*, Athens, Greece.