

Stylebook for the Tübingen Treebank of Written German (TüBa-D/Z)

Heike Telljohann, Erhard W. Hinrichs, Sandra Kübler,
Heike Zinsmeister, Kathrin Beck

Universität Tübingen
Seminar für Sprachwissenschaft
Wilhelmstr. 19
D-72074 Tübingen

August 2015

Abstract

This stylebook is an updated version of Telljohann et al. (2012). It describes the design principles and the annotation scheme for the German treebank TüBa-D/Z developed by the Division of Computational Linguistics (Lehrstuhl Prof. Hinrichs) at the Department of Linguistics (*Seminar für Sprachwissenschaft – Sfs*) of the Eberhard Karls Universität Tübingen, Germany. The guidelines focus on the syntactic annotation of written language data taken from the German newspaper 'die tageszeitung' (taz). The unannotated taz newspaper material was taken from the Science CD (*Wissenschafts-CD*) of 'die tageszeitung' (taz) that can be licensed from contrapress media GmbH (http://shop.taz.de/index.php?cat=c18_taz-Archiv.html).

At present, the treebank comprises 3,644 articles (95,595 sentences) selected from the taz editions between 1989 and 1999. The average sentence length is 18.7 words and the total number of tokens currently amounts to 1,787,801. The TüBa-D/Z treebank is still under development. Thus, the number of annotated sentences will increase over time. Periodic data updates and accompanying updates of this stylebook will be made available at:

<http://www.sfs.uni-tuebingen.de/en/ascl/resources/corpora/tueba-dz.html>

Please consult this website in order to ensure that you are using the most recent and most complete version of the treebank.

The annotation scheme for the TüBa-D/Z treebank is derived from the VERB-MOBIL treebank for spoken German, developed earlier (1997–2000) by the Division of Computational Linguistics of the Sfs (Hinrichs et al. 2000). The TüBa-D/Z annotation scheme has been extended along various dimensions to accommodate the characteristics of written texts. In order to ensure the reusability of the data, a surface-oriented annotation scheme has been adopted that is inspired by the notion of topological fields and is enriched by a level of predicate-argument structure. The linguistic inventory used in the treebank annotation is based on a minimal set of assumptions that are uncontroversial among major syntactic theories. In this sense it is an attempt at theory-neutrality.

Acknowledgements

Funding for the TüBa-D/Z has come from a variety of sources:

- the *Competence Center for Text- and Information Technology* (Kompetenzzentrum für Text- und Informationstechnologie – KIT)) grant by the Ministry of Science, Research and the Arts Baden-Württemberg (funding since 2000);
- the collaborative research center (Sonderforschungsbereich) grant *SFB 441 – Linguistic Data Structures*, project A1 – Representation and Automatic Acquisition of Linguistic Data funded by the German Research Council (Deutsche Forschungsgemeinschaft – DFG);
- the collaborative research center (Sonderforschungsbereich) grant *SFB 833 – The construction of meaning - the dynamics and adaptivity of linguistic structures*, project A3 – Disambiguating Discourse Connectives using Corpus-induced Semantic Relations funded by the German Research Council (Deutsche Forschungsgemeinschaft – DFG);
- the ESFRI research infrastructure project grants *D-SPIN* and *CLARIN-D* funded by the Federal Ministry of Education and Research (BMBF) (funding since 2008).

A project of this scale would not be possible without the generous support from many contributors:

Our special thanks go to 'die tageszeitung' (taz) who kindly granted permission to process the newspaper data and to release the treebank.

We would like to acknowledge Rosmary Stegmann for her many contributions to the treebank of spoken German in VERBMOBIL. Her research laid the foundations for the annotation scheme of that treebank, which has been summarized in the 'Stylebook for the German Treebank in VERBMOBIL' (Stegmann et al. 2000).

We would like to thank Manfred Sailer and Frank Richter for their helpful comments and support in form of encouragement and critical discussions from which we could strongly benefit for the challenging task of developing a data-oriented syntactic annotation scheme for spoken as well as for written German.

Furthermore, we are indebted to Tylman Ule for his assistance with part-of-speech tagging of the data and with data conversion.

We would also like to acknowledge the support of Martina Liepert and Jorn Veenstra, who initiated and developed the integration of named entities into the annotation scheme.

Moreover, we would like to thank Julia Trushkina (Trushkina 2004) and Yannick Versley (Versley et al. 2010) who provided the tools for morphological preprocessing.

Furthermore, Yannick Versley (Versley et al. 2010) supported the project by developing a tool for lemma disambiguation and for the automatic integration of semantic classes of named entities.

The quality of the treebank has been considerably improved by feature oriented consistency checks developed by Ventsislav Zhechev. Further consistency tests were contributed by Tylman Ule and Frank H. Müller in the course of their research work in the SFB 441. They deserve special mention for their support.

We would like to thank Marie Hinrichs for managing the complete tool chain and carrying out the many steps of data pre-processing, integration, and post-processing required to support the life cycle of a TüBa-D/Z release.

We would like to thank Vera Möller and Karin Naumann (2007) for annotating anaphora and coreference relations and also for doing an excellent job in documenting the concepts.

Yannick Versley and Holger Wunsch supported the project in various aspects. In the course of their Ph.D. projects in the SFB 441 they enhanced the conceptual aspects of the anaphora resolution as annotated in the treebank. They also wrote mapping and conversion tools for integrating the anaphora annotation in the Export-XML format.

For their diligence and dedication to the arduous task of linguistic annotation and of post-editing we thank our research assistants Janne Berlacher, Anne Brock, Armin Buch, Nadine Cetin, Heike da Silva Cardoso, Marisa Delz, Silke Dutz, Katrin Eichler, Emilia Ellsiepen, Steffen Froemel, Holger Gauza, Simone Hartung, Daniel Hüttl, Heike Johannsen, Miriam Käshammer, Laura Kassner, Sarah Klug, Julia Koch, Janina Kopp, Anuschka Kranz, Christian Kreß, Rebecca Kreß, Michael Kossack, Anne Lohse, Wolfgang Maier, Nicole Maruschka, Kai Metzger, Vera Möller, Simone Müller, Till Pachalli, Maja Pietsch, Brigitta Rist, Andreas Rudin, Maria Schmidt, Marie Schreier, Insa Starr, Melanie Störzer, Isabel Trott, and Dominikus Wetzler. They also improved the linguistic quality of the annotation by dedicated discussions on problematic and interesting examples.

The development of the TüBa-D/Z treebank was notably facilitated by a number of former VERBMOBIL partners whose contributions went well beyond the call of duty. Hans Uszkoreit and his colleagues at the Saarland University kindly provided us with the graphical annotation tool *Annotate* (Plaehn 1998) which was developed as part of the research project (*Teilprojekt C3*; Principal investigators: Uszkoreit/Smolka) *Nebenläufige grammatische Verarbeitung* (NEGRA) in the collaborative research center (Sonderforschungsbereich) 378. The *Annotate* tool provides human annotators with a graphical, user-friendly interface for annotating and editing trees and also offers database support for maintaining large treebanks. We would like to express our special gratitude to Thorsten Brants, who has kindly and generously provided us with software support and user assistance for the *Annotate* tool from the very beginning of the Tübingen treebank project.

Contents

List of Tables	8
1 Introduction	9
2 Major Challenges and Design Decisions	11
3 The Theoretical Basis of the Annotation Scheme	14
3.1 Topological Fields	14
3.1.1 The Concept of Topological Fields	14
3.2 Constituent Analysis and Topological Fields	17
3.3 General Annotation Principles	18
3.3.1 Flat Clustering Principle	18
3.3.2 Longest Match Principle	18
3.3.3 High Attachment Principle	18
3.4 The Structure of an Annotated Tree	19
3.4.1 The Levels of Annotation	19
3.4.2 The Inventory of Labels	19
3.4.3 What Is a Syntactic Unit?	22
3.4.4 Printing and Spelling Errors	28
3.4.5 Isolated Phrases	29
3.4.6 Long-Distance Dependencies	31
3.4.7 Empty Categories	32
3.5 Lemma Information	33
3.5.1 Lemmatization Rules for POS-Tags	33
3.5.2 Lemmatization Rules for Specific Linguistic Phenomena	37
4 The Annotation of the Internal Structure of Phrases	40
4.1 Premodification and Postmodification in Phrases	40
4.2 Noun Phrases	40
4.2.1 Noun Phrases without Modifiers	40
4.2.2 Prenominal Modification	41
4.2.3 Postnominal Modification	46
4.2.4 Appositional Constructions	49
4.2.5 Foreign Language Material	53
4.2.6 Named Entity Annotation	56
4.2.7 Ordinal Numbers	64

4.2.8	Cardinal Numbers	64
4.2.9	Letters and Non-Words	66
4.2.10	Expletive and Other Uses of <i>es</i>	67
4.3	Determiner Phrases	70
4.4	Prepositional Phrases	71
4.4.1	Prepositions	71
4.4.2	Circumpositions and Postpositions	74
4.5	Adjectival Phrases	74
4.6	Adverbial Phrases	79
4.7	Verb Phrases	81
4.7.1	Head of a Sentence and Verb Complex	81
4.7.2	Verb Complexes in Verb-second and Verb-final Clauses	81
4.7.3	Ersatzinfinitiv Constructions	83
4.7.4	Infinitives with <i>zu</i>	85
4.7.5	Coherency and Incoherency of Verbal Constructions	87
4.7.6	AcI Constructions	88
4.7.7	Imperatives	89
4.7.8	Particle Verbs	90
4.7.9	Verbs with Predicate	91
4.7.10	Modal Verbs	94
5	Attachment Principles for Phrases	96
5.1	Attachment to Fields	96
5.2	Attachment of Ambiguous Complements	96
5.3	Modifier Attachment	97
5.3.1	Modifier Attachment in the Initial Field	99
5.3.2	Attachment across Punctuation Marks	99
5.3.3	Ambiguous Modifiers in Isolated Phrases	100
6	The Annotation of Sentences	102
6.1	Sentence Initial Fields	102
6.1.1	The C-Field in Verb-Final Clauses	102
6.1.2	The KOORD-Field in all Clause Types	104
6.1.3	The PARORD-Field in Verb-Second Clauses	105
6.1.4	Resumptive Constructions: The LV-Field	105
6.2	Questions	106
6.2.1	W-Questions	106
6.2.2	Yes - No Questions	107
6.3	Clauses of Comparison	108
6.4	Relative Clauses	109
6.4.1	Event-modifying Relative Clauses	111
6.4.2	Independent Relative Clauses	111
6.5	Coordination	112
6.5.1	Coordination of Phrases	113
6.5.2	Asymmetric Coordination	114
6.5.3	Coordinations with Complex Conjunctions	115

6.5.4	Coordinations with Truncated Words	116
6.5.5	Attachment Principles of Coordination within Phrases	118
6.5.6	Coordination of Topological Fields	119
6.5.7	Attachment of Ambiguous Modifiers in Coordination	120
6.5.8	Coordination of Sentences	122
6.5.9	Paratactic Constructions	124
6.5.10	Conjunctions Occurring with Isolated Phrases	124
6.5.11	Split Coordinations	126
6.6	Elliptical Constructions	127
7	The Annotation of Specific Syntactic Phenomena	130
7.1	Superlative and Comparative Forms	130
7.1.1	Superlative Forms	130
7.1.2	The Comparative Particles <i>wie</i> and <i>als</i>	130
7.2	Verbal and Adjectival Use of Participles	133
7.3	Topicalization	134
7.4	Headlines	135
7.5	Discourse Markers	137
7.6	Parentheses	139
8	Criteria for the Distinction of Grammatical Functions	141
8.1	Subcategorization of Verbs	141
8.2	Subcategorization of PREDs	141
8.3	Distinction of FOPP, OPP, and V-MOD	142
8.4	Distinction of MOD, MOD-MOD, and V-MOD	143
8.5	Distinction of ON, PRED, ON-MOD, and PRED-MOD	143
9	The TüBa-D/Z Data Formats	146
9.1	The NEGRA Export Format	146
9.2	The Penn Treebank Format	150
9.2.1	The Penn Treebank Format Version 1	150
9.2.2	The Penn Treebank Format Version 2	153
9.3	The Export-XML Format	155
9.4	The CoNLL Format (2006, 2010, 2011/2012)	157
9.4.1	The CoNLL 2006 Format	157
9.4.2	The CoNLL 2010 Format	158
9.4.3	The CoNLL 2011/2012 Format	159
	References	160
	Index	163

List of Tables

3.1	Three clause types according to Höhle (1986)	15
3.2	Topological fields	16
3.3	Levels of annotation	19
3.4	The STTS tag set	21
3.5	Morphological feature combinations for lexical elements	23
3.6	Values of morphological features	24
3.7	Node labels	25
3.8	Edge labels	26
3.9	Syntactic-Semantic Node Labels for Named Entities	27
3.10	Lemmatization rules for POS-tags	33
3.11	Lemmatization rules for specific linguistic phenomena	37
4.1	Semantic Classes and Subclasses for Named Entities	57
4.2	Types of <i>es</i>	69

Chapter 1

Introduction

The purpose of this report is to describe the design principles and annotation scheme for the TüBa-D/Z treebank of German. It is intended as a guide for the treebank annotators in Tübingen and for theoretical and computational linguists who want to use annotated treebank data for their own research. In addition, we hope that this report may be of some use for researchers who want to construct their own treebank for German or for some other language. We would like to emphasize that the annotation scheme is language-specific, and we advise against adopting this scheme without modification for some other language. However, we do believe that the type of design decisions that are reported here for German will arise for other languages as well. And it is in this sense that the current report could provide an useful point of reference.

The TüBa-D/Z treebank was developed by the Division of Computational Linguistics (Lehrstuhl Prof. Hinrichs) at the Department of Linguistics (*Seminar für Sprachwissenschaft – SfS*) of the Eberhard Karls Universität Tübingen, Germany. The guidelines focus on the syntactic annotation of written language data taken from the German newspaper 'die tageszeitung' (taz). The unannotated taz newspaper material was taken from the Science CD (*Wissenschafts-CD*) of 'die tageszeitung' (taz) that can be licensed from contrapress media GmbH (http://shop.taz.de/index.php?cat=c18_taz-Archiv.html).

At present, the treebank comprises 95,595 sentences. The newspaper material is taken from the taz editions from

1989 628 articles from 251 days over 12 months, 32,267 sentences.

1989 632 articles from 4 days over 1 month, 12,245 sentences.

1995 1,107 articles from 6 days over 1 month, 21,391 sentences

1997 238 articles from 154 days over 12 months, 7,497 sentences

1999 1,039 articles from 6 days over 2 months, 22,195 sentences

Total 3,644 articles from 421 days over 28 months from 5 years, 95,595 sentences

The average sentence length is 18.7 words and the total number of tokens currently amounts to 1,787,801. The TüBa-D/Z treebank is still under development. Thus, the number of annotated sentences will increase over time. Periodic data updates and accompanying updates of this stylebook will be made available at:

<http://www.sfs.uni-tuebingen.de/en/ascl/resources/corpora/tueba-dz.html>

Please consult this website in order to ensure that you are using the most recent and most complete version of the treebank.

The annotation scheme for the TüBa-D/Z treebank is derived from the VERBMOBIL treebank for spoken German, developed earlier (1997–2000) by the Division of Computational Linguistics of the Sfs (Hinrichs et al. 2000). The annotation scheme for the VERBMOBIL treebank has been summarized in the ‘Stylebook for the German Treebank in VERBMOBIL’ (Stegmann et al. 2000). The TüBa-D/Z annotation scheme has been extended along various dimensions to accommodate the characteristics of written texts. In order to ensure the reusability of the data, the linguistic inventory used in the treebank annotation is based on a minimal set of assumptions that are uncontroversial among major syntactic theories. In this sense it is an attempt at theory-neutrality.

The TüBa-D/Z treebank is released in four different data formats : the Negra Export format, the Export-XML format, the Penn treebank format (version 1 and 2), and the CoNLL format (2006, 2010, 2011/2012). More information about each data format is given in chapter 9.

To the best of our knowledge, the VERBMOBIL treebank for spoken German is still the only treebank based on non-genre-specific German speech data. It is released as TüBa-D/S treebank (<http://www.sfs.uni-tuebingen.de/en/ascl/resources/corpora/tueba-ds.html>). For written texts, TüBa-D/Z is not the only treebank available for German. Two other (semi-)manually annotated treebanks are currently available, each with their own annotation scheme: the Negra treebank (<http://www.coli.uni-saarland.de/projects/sfb378/negra-corpus/>) and the TIGER treebank (<http://www.ims.uni-stuttgart.de/forschung/ressourcen/korpora/tiger.html>).

The Tübingen Partially Parsed Corpus of Written German (TüPP-D/Z; <http://www.sfs.uni-tuebingen.de/en/ascl/resources/corpora/tuepp-dz.html>) is a project closely related to the TüBa-D/Z treebank. It consists of 200 million word tokens of the Science CD (*Wissenschafts-CD*) of ‘die tageszeitung’ (taz), including the sentences which are annotated in the TüBa-D/Z treebank. The texts were automatically annotated with clause structure, topological fields, and chunks, in addition to more low level annotation including parts of speech and morphological ambiguity classes. The first release of TüBa-D/Z (12/2003) functioned as training corpus.

Chapter 2

Major Challenges and Design Decisions

Most syntactic theories consider individual sentences as the primary domain of linguistic theorizing and of syntactic annotation. For written language, the segmentation into sentences is largely unproblematic and coincides with the domain of syntactic analysis.

However, newspaper texts exhibit a number of phenomena that do not lend themselves easily to a purely sentence-based annotation. These phenomena include: headlines, titles, parentheses, discourse markers, and sentence conjunction by a colon. These cases are described in more detail in sections 3.4.3 to 3.4.5 of this stylebook.

The second main question, which needed to be addressed at the outset of the project was the inventory of syntactic categories and grammatical functions to be used for syntactic annotation and specification of predicate-argument structure. Here our choices were guided by two main considerations:

1. Linguistic adequacy and theory-neutrality: For the purposes of reusability of the treebank data, the annotation scheme should not reflect a commitment to a particular syntactic theory. Rather, the inventory of categories should be a reflection of common assumptions that syntacticians share across different frameworks concerning questions of constituenthood, phrase attachment, and grammatical functions. On this note, the annotations should be theory-neutral and minimal. This desideratum is of utmost importance so as to ensure the reusability of the annotated data.

At the same time, the annotation scheme should reflect as much as possible those empirical generalizations that syntacticians, especially from a descriptive perspective, have identified as characteristic of the language in question.

2. Balancing the needs of potential users: Since the construction of a treebank is a labor-intensive and costly enterprise, ideally the TüBa-D/Z treebank should appeal to as many potential users as possible. Moreover, the treebank should be of interest to researchers of a wide range of different fields. Considering the renewed interest in the use of corpora for both theoretical and computational linguistics, choicepoints in the annotation scheme should be resolved in such a way that the needs of potential users are balanced as much as possible.

To support the use of the TüBa-D/Z treebank in computational linguistics, the annotation scheme should be sensitive to processing considerations, as long as linguistic adequacy of the choice of annotations is not compromised. *Ceteris paribus*, processing considerations favor annotation schemes that pay close attention to properties of syntactic surface structure, particularly to word order regularities and distributional properties of words and phrases. At the same time, the use of empty categories and data structures with crossing dependencies among phrases are to be avoided if the annotations are to be used for parsers that rely on the context-freeness of the underlying grammar.

In order to satisfy the above aims, the annotation scheme is surface-oriented and context-free. The theoretical assumptions underlying the levels of annotation and the choice of labels themselves are as much as possible based on a rich tradition of theoretical and empirical research on German syntax.

For the treatment of word regularities of German, which is a language with relatively free word order, an inventory of topological fields is incorporated into the annotation scheme. Topological fields in the sense of Herling (1821), Erdmann (1886), Drach (1937), and Höhle (1986) are widely used in descriptive studies of German syntax. Such fields constitute an intermediate layer of analysis above the level of individual phrases and below the clause level. The concept of topological fields favors tree-based annotations, i.e. bracketings that do not rely on crossing or discontinuous dependencies. Instead, such non-linear dependencies are to be expressed at the level of predicate-argument structure which constitutes a second level of annotation with its own descriptive inventory of grammatical functions.

The framework of topological fields is widely used in empirical and theoretical accounts of German syntax. Thus, it is in the linguistics literature. This greatly facilitates thorough training of human annotators, since they can rely on the pre-existing body of literature. One purpose of this stylebook is to add to these reference materials.

Currently, a total of 25 syntactic node labels for the encoding of constituent structures are being used. These include labels for topological fields as well as labels for phrases and their constituent parts.

In order to capture grammatical functions of individual phrases and syntactic dependencies between phrases, constituent structure trees are enriched by a set of edge labels between constituent structure nodes. The current inventory of edge labels comprises 42 distinct categories. In addition to these primary edge labels, four secondary edge labels are used. These labels indicate phrase-internal government of elements in the verb complex, express phrase-internal modification of noun phrases, resolve long-distance dependencies among modifiers, or relate the phrasal complements of so-called *third-construction* control verbs.

For certain computational applications, robust identification of *named entities*, e.g. person names, names of companies and institutions, names of geographical locations, is a major concern. Therefore, such named entities are identified by a special node label, and their internal structure is sometimes identified by an additional secondary edge label that is used exclusively for named entities.

At the word level, part-of-speech labels are assigned according to the Stuttgart-Tübingen tag set, which is widely accepted for part-of-speech tagging for German and which provides an inventory of 54 distinct part-of-speech labels. In addition, information on inflectional morphology is given.

Detailed information about the complete inventory of node labels, edge labels, part-of-speech labels and inflectional feature clusters is given in section 3.4.2 of this stylebook.

The remainder of this stylebook is organized as follows: chapter 3 offers an overview of the theoretical foundations of the annotation scheme, focusing on the concept of topological fields (3.1) and its relation to constituent structure (3.2), on general annotation principles (3.3), as well as an overview of the annotation levels and of the inventory of the annotation labels for each level (3.4). Chapter 4 concerns the annotation of the internal structure of phrases, broken down into major word classes and their phrasal projections. Chapter 5 addresses the principles for relating individual phrases to each other, particularly for modifier and complement attachment. Chapter 6 discusses the annotation of entire sentences, focusing on the relationship between sentence types and topological fields, coordination (including phrasal conjunction) and elliptical constructions. Chapter 7 is devoted to the annotation of miscellaneous syntactic constructions such as comparatives, verbal and adjectival participles, topicalization, newspaper headlines, discourse markers, and parentheses, which each pose special challenges for the annotation tasks. Chapter 8 describes the criteria used for distinguishing different grammatical functions. Chapter 9 describes the five different data formats in which the TüBa-D/Z treebank is distributed. The stylebook concludes with a bibliography and a subject index.

We do not consider the annotation level of anaphora and coreference relations in this stylebook. Please consult (Naumann and Möller 2007) for a detailed description of these phenomena.

Chapter 3

The Theoretical Basis of the Annotation Scheme

3.1 Topological Fields

The annotation scheme for the TüBa-D/Z treebank has been developed with special regard to the characteristics of the German language: the interaction of configurational and non-configurational syntactic properties, which arise from the partially free word order. On the one hand, there exist three different clause types with respect to the fixed position of the finite verb (verb-second (V-2), verb-initial (V-1), and verb-final (V-end)). On the other hand, there is a high degree of variability of complements and adjuncts. In order to treat the relatively high degree of word order freedom in German, the treebank adopts the notion of topological fields as the primary clustering principle of a sentence.

The basic characteristics of the model of topological sequences within a German sentence were originally formulated by Herling (1821) and Erdmann (1886). Herling (1821) developed an adequate topological theory for complex sentences in which clauses form a topological carrying a syntactic function and he mentioned the special position of the finite verb in verb-second und verb-final clauses. Erdmann (1886) established the basics of a theory of topological fields and pointed out that the first position in a clause is not necessarily the subject position. The so called Herling/Erdmann scheme already covers a set of word order regularities which apply for all three clause types of German. Later Drach (1937) introduced the notion of *field*. Finally, Höhle (1986) developed topological schemes for the three clause types.

3.1.1 The Concept of Topological Fields

In a German clause, the finite verb can appear in three different positions: verb-second, verb-initial, and verb-final. Only in verb-final clauses the verb complex consisting of the finite verb and non-finite verbal elements forms a . The discontinuous positioning of the verbal elements in verb-first and verb-second clauses is the traditional reason for structuring German clauses into fields. The positions of the verbal elements form the *Satzklammer* (sentence bracket) which divides the sentence into a *Vorfeld* (initial field), a *Mittelfeld* (middle field), and a *Nachfeld* (final field). The *Vorfeld* and the *Mittelfeld*

are divided by the *linke Satzklammer* (left sentence bracket), which is the finite verb, the *rechte Satzklammer* (right sentence bracket) is the verb complex between the Mittelfeld and the Nachfeld. Thus, the theory of topological fields states the fundamental regularities of German word order. It is an important basis for the topological analysis of any German sentence, since subclauses and embedded clauses are treated within the bounds of fields. Identical word order regularities within a specific field can be realized in all three clause types. But the fields themselves differ in their possible elements and grammatical rules. Therefore, the theory is a descriptive rather than explanatory theory for a specific language.

Höhle (1986) denotes the three clause types as E-Sätze (verb-final clauses), F1-Sätze (verb-initial clauses), and F2-Sätze (verb-second clauses). The topological schemes of these types are listed in Table 3.1.

Table 3.1: Three clause types according to Höhle (1986)

E-Sätze	(KOORD) - (C) - X - VK - Y
F1-Sätze	(KOORD - (KL) - FINIT - X - VK - Y
F2-Sätze	(KOORD or PARORD) - (KL) - K - FINIT - X - VK - Y

Abbreviations and explanations used in Table 3.1:

VK: verb complex

FINIT: element denoting categories of finiteness

KOORD: coordinating particles (e.g. *und*, *oder*)

PARORD: non-coordinating particles (e.g. *denn*, *weil*)

X, Y: sequence of any number of constituents

C: complementizer

K: one constituent

KL: nominativus pendens, resumptive construction (*Linksversetzung*)

These schemes topologically analyse not only atomic sentences but also complex sentence constructions which contain embedded clauses. Such embedded clauses can occur in a *Linksversetzung* (resumptive construction), Vorfeld, Mittelfeld, or Nachfeld. Herling's theory of the coordination and embedding of sentences covers these phenomena in detail (Herling 1821).

According to Höhle (1986), we assume the existence of the following topological fields (cf. Table 3.2):

The following description of the topological fields does not claim completeness regarding all descriptive details but rather mentions their main characteristics.¹

¹In the following, the abbreviations for the fields listed in Table 3.2 are used.

Table 3.2: Topological fields

Field	Description
VF	Vorfeld (initial field)
LK	Linke (Satz-)Klammer (left sentence bracket)
MF	Mittelfeld (middle field)
VC	Verbkomplex (verb complex)
NF	Nachfeld (final field)
LV	Linksversetzungsfeld (field for resumptive constructions)
C	C-Feld (field for complementizers, left from MF)
KOORD	Koordinationsfeld (field for coordinating particles) left-most element, optionally in all clause types, (e.g. <i>und, oder</i>)
PARORD	Koordinationsfeld (field for non-coordinating particles) left-most element, optionally only in verb-second (e.g. <i>denn, weil</i>)

VF: The *Vorfeld* consists of only one constituent. Usually it is the subject². But because of the high degree of non-configurationality in German, the subject can also occur in the *Mittelfeld*, thus allowing almost every other constituent to occupy the *Vorfeld*.

LK: The *Linke Klammer* is the position of the finite verb in verb-second and verb-first clauses or a conjunction in verb-final clauses. It consists of exactly one element.

MF: Apart from those s which are optionally located in other fields, any non-verbal constituent may occur in the *Mittelfeld*. It consists of a sequence of any number of constituents. The linear order of the constituents depends on the specific word order principles for German and their interaction.

VC: The *Verbkomplex* is a sequence of verb forms. In verb-second and verb-first clauses it consists of one or more non-finite elements or - depending on the verb - of a separable prefix. In verb-final clauses it also contains the finite verb. The rule for the linear order in general is: right determines left. If there is a finite verb in the verb complex, it is usually the right-most element (exception: *Ersatzinfinitiv* constructions (*daß er sich ein neues Konzept wird überlegen müssen*) (cf. 4.7.3).

NF: For some clause types (e.g. *so daß*-Sätze), the *Nachfeld* is the obligatory position. Embedded complement clauses, relative clauses, and single constituents can optionally occur in the *Nachfeld*. In contrast to the *Vorfeld* it may be occupied by any number of constituents.

LV: The *Linksversetzungsfeld* is a field for the left-dislocated phrase of resumptive constructions. A *Linksversetzung* is a pendent constituent. It can be regarded as a

²In the fifth release, 52.5% of all *Vorfeld* fields host the subject.

syntactic anticipation of a part of a sentence (cf. 6.1.4). There are many restrictions which apply for this position.

C: The *C-Feld* only occurs in verb-final clauses. It is obligatorily occupied in finite verb-final clauses if there is no conjunction in the Linke Klammer. In non-finite verb-final clauses the C-position may be empty. This field can be occupied by conjunctions of sentential objects (e.g. *daß, ob*) or sentence initial conjunctions like *um, obwohl, wenn* and also by complex interrogative or relative phrases, e.g. ..., '*um wieviel Geld*' geht es dabei? / ..., '*an der*' *Max Daniel Professor für Klavier ist*. (cf. 6.1.1).

KOORD: The KOORD-field is the field for coordinating particles. In contrast to the PARORD-field, it can optionally occur as the left-most element of all clause types (cf. 6.1.2).

PARORD: The PARORD-field is the field for non-coordinating particles which optionally occur as the left-most element of a verb-second clause (cf. 6.1.3).

Concerning the distribution of constituents to topological fields see also the chapter *Deskriptive Generalisierungen* in Grewendorf (1991).

The combination of these fields in order to constitute verb-first, verb-second, or verb-final clauses is described in Höhle (1986).

The topological model, which is the basis of most traditional German grammars, only provides descriptive parameters concerning the sentence structure without making any statement about the regularities within the fields and the hierarchical constituent structure of the sentence. For more complicated phenomena, it offers only a catalogue of detailed descriptions.

3.2 Constituent Analysis and Topological Fields

The main weakness of the concept of topological fields is the above-mentioned fact that the hierarchical constituent structure of a sentence cannot be described. The aim is to find a form of representation which combines the topological model with a constituent analysis in order to describe the hierarchy of the linguistic s within the fields. In our annotation scheme, the integration of a constituent analysis was achieved by a second level of annotation strictly within the bounds of topological fields: a predicate-argument structure with its own descriptive inventory of syntactic categories and grammatical functions. The constituent structure is represented by phrase structure trees (phrase markers) whose node and edge labels carry this information.

In order to analyse syntactic constructions, it is necessary to define the number and types of constituents within the fields.

1. **Number of constituents within the fields:**

In general, C, LK, KOORD, PARORD, and VF contain only **one** constituent. More than one constituent is allowed within MF and NF.

2. **Types of constituents within the fields:**

Phrasal constituents occur in VF, MF, NF and C (interrogative or relative phrases). Embedded clauses either belong to NF, VF, LV, or in some cases to MF. Usually, outside the spoken language context, verb-final clauses do not occur isolated. They need to be attached if possible.

3.3 General Annotation Principles

Our annotation scheme tries to find a trade-off between pragmatic requirements on the one hand and linguistic reality on the other hand. The following three common annotation principles are adopted to group the constituents within a syntactic tree: the *flat clustering principle*, the *longest match principle*, and the *high attachment principle*.

3.3.1 Flat Clustering Principle

The *flat clustering principle* keeps the number of hierarchy levels in a syntactic structure as small as possible. As a consequence, any degree of branching is allowed. Constituents which cannot be assigned a grammatical function within a syntactic construction are structured as much as possible, but are not typically connected to surrounding constituents as a whole.

3.3.2 Longest Match Principle

The *longest match principle* demands that as many daughter nodes as possible are combined into a single mother node, provided that the resulting construction is syntactically as well as semantically well-formed.

3.3.3 High Attachment Principle

The *high attachment principle* prescribes that syntactically and semantically ambiguous modifiers are attached to the highest possible level in a tree structure. Premodifiers and postmodifiers are treated in a different way. First, both kinds of modifiers are projected to their phrase level. Since the modification scope of premodifiers is unambiguous, they are directly attached to the head of the phrase which they are modifying. By contrast, postmodifiers are always attached on a higher level to preserve ambiguity. This decision was taken to avoid the problematic distinction whether a postmodifier is a free adjunct or a complement of the modified phrase.

3.4 The Structure of an Annotated Tree

3.4.1 The Levels of Annotation

A syntactic tree consists of nodes and edges. Nodes represent constituents on different levels of annotation. Edges always link daughter nodes to a mother node. The root node of a tree is assumed as the sentence node of a construction. One level below the sentence node, the nodes of the topological fields are located. This is the reason why topological fields can be regarded as the top-level ordering principle for sentences in the treebank. The sequence of the fields in the three clause types never violates the topological schemes given by Höhle (1986). Within each sentence structure, in general at least two topological fields are occupied (exception: infinitive constructions, (cf. 4.7.4). Others may be left empty (elliptical constructions, cf. 6.6). Table 3.3 lists the four levels of annotation which we distinguish within the structure of an annotated syntactic tree³:

Table 3.3: Levels of annotation

Level	Inventory
clause level	root node labels for different types of clauses
field level	node labels for topological fields (including labels for conjuncts of fields)
phrase level	node labels for syntactic categories (including syntactic-semantic node labels for named entities) and edge labels for grammatical functions
lexical level	lexical entries tagged with the part-of-speech (POS-)tags taken from the STTS tag set (Schiller et al. 1995) and with morphological features (Trushkina 2004, Versley et al. 2010) and lemmata (Versley et al. 2010)

Node labels denote the syntactic category of a phrase or sentence, a topological field, or a grammatical property. Edge labels denote the grammatical function of lexical entries, phrases, topological fields, and clauses.

3.4.2 The Inventory of Labels

The **part-of-speech tags** used for the annotation are taken from the Stuttgart-Tübingen tag set (STTS) (Schiller et al. 1995).⁴ The STTS is a guideline for the annotation of German text corpora on the lexical level. Every single part-of-speech of a text is assigned one specific tag. The tag set consists of the tags listed in Table 3.4 (cf. (Schiller et al. 1995)). The tagging of the data was performed by the *tnt* tagger (Brants 1998) and manually corrected with the *Annotate* tool (Plaehn 1998).

³We do not consider the suprasentential annotation level of anaphora and coreference relations in this stylebook. Please consult (Naumann and Möller 2007) for a detailed description of these phenomena.

⁴PAV was changed into a new tag called PROP (pronominal form of a prepositional phrase) in order to justify PX as the syntactic category of its mother.

The **morphological tags** give information about inflectional morphology and include features such as *case*, *number*, *person*, etc. A specific combination of feature-value pairs is defined for each relevant part-of-speech category, see Table 3.5 for the list of part-of-speech categories that are annotated with morphological features and the corresponding feature combinations. The values are represented in a cluster by single character abbreviations, see Table 3.6 for the set of features and their values. Features can uniquely be identified by their position in the cluster.

Node labels indicate the syntactic category of a phrase or sentence, but they are also used to label topological fields and sequences of topological fields within coordinations or to indicate specific grammatical properties of constituents. Table 3.7 lists all node labels which are used in the treebank. (An additional node is introduced for named entities, see Table 3.9)

Edge labels indicate the grammatical function of lexical entries, phrases, topological fields, and clauses. Since case information is given and a distinction of different modifiers is made by these labels, the syntactic tree structures also contain semantic roles. The specific set of edge labels for the German treebank is listed in Table 3.8, including **secondary edge labels**. The latter ones are used to resolve ambiguities on a different level of description.

Two specific edge labels denote whether a constituent has the function of a head (HD), e.g. a phrase (NX, PX, ADJX, ADVX, VXFIN, VXINF), or a non-head (-), e.g. a determiner or a modifier attached to a phrase. On any annotation level, there is at most one head. Within phrases, these two labels indicate the internal dependency structure of the phrase. The head of a sentence structure (e.g. SIMPX) is always the finite verb. In coordinations, each conjunct depends on the head of the whole construction and is denoted with a specific edge label (KONJ) in order to distinguish them from conjunctions and modifying elements within a coordination (see 6.5.1 and 6.5.3). Edge labels below all root node labels carry only non-head labels (cf. (Kübler and Telljohann 2002)).

In an enhanced version of the TüBa-D/Z treebank, each named entity is assigned one of the following semantic classes: person (PER), organisation (ORG), location (LOC), geopolitical entity (GPE), or other (OTH). The semantic class OTH comprises all remaining named entities not fitting into PER, ORG, LOC, or GPE (cf. 4.2.6).

In order to annotate these semantic classes, **syntactic-semantic node labels** of the pattern *syntactic category = semantic class* are defined as the mother node of named entities (see Table 3.9). These syntactic-semantic nodes indicate that the structure below represents a (complex) named entity of a certain syntactic category belonging to one of the five semantic classes (e.g. *Ute Wedemeier* (NX=PER), *The Jim Wane Swingtett* (NX=ORG), *Sögestraße* (NX=LOC), *Auf die stürmische Art* (PX=OTH) (cf. 4.2.6).

The former node label 'EN-ADD' and the secondary edge label 'EN' are deleted.

The internal syntactic structure of named entities is governed by the general annotation rules. All parts below a syntactic-semantic node that do not belong to the named entity itself are marked as '-NE', e.g. *[[die (-NE)] AWO]* (NX=ORG), *[[Der (-NE)] zweite Weltkrieg]* (NX=OTH).

Table 3.4: The STTS tag set

POS =	description	examples
ADJA	attributive adjective	<i>[das] große [Haus]</i>
ADJD	adverbial or predicative adjective	<i>[er fährt] schnell, [er ist] schnell</i>
ADV	adverb	<i>schon, bald, doch</i>
APPR	preposition; left circumposition	<i>in [der Stadt], ohne [mich]</i>
APPRART	preposition + article	<i>im [Haus], zur [Sache]</i>
APPO	postposition	<i>[ihm] zufolge, [der Sache] wegen</i>
APZR	right circumposition	<i>[von jetzt] an</i>
ART	definite or indefinite article	<i>der, die, das, ein, eine</i>
CARD	cardinal number	<i>zwei [Männer], [im Jahre] 1994</i>
FM	foreign language material	<i>[Er hat das mit “] A big fish [” übersetzt]</i>
ITJ	interjection	<i>mhm, ach, tja</i>
KOUI	subordinating conjunction with <i>zu</i> + infinitive	<i>um [zu leben], anstatt [zu fragen]</i>
KOUS	subordinating conjunction with clause	<i>weil, daß, damit, wenn, ob</i>
KON	coordinative conjunction	<i>und, oder, aber</i>
KOKOM	particle of comparison, no clause	<i>als, wie</i>
NN	noun	<i>Tisch, Herr, [das] Reisen</i>
NE	proper noun	<i>Hans, Hamburg, HSV</i>
PDS	substituting demonstrative pronoun	<i>dieser, jener</i>
PDAT	attributive demonstrative pronoun	<i>jener [Mensch]</i>
PIS	substituting indefinite pronoun	<i>keiner, viele, man, niemand</i>
PIAT	attributive indefinite pronoun without determiner	<i>kein [Mensch], irgendein [Glas]</i>
PIDAT	attributive indefinite pronoun with determiner	<i>[ein] wenig [Wasser], [die] beiden [Brüder]</i>
PPER	irreflexive personal pronoun	<i>ich, er, ihm, mich, dir</i>
PPOSS	substituting possessive pronoun	<i>meins, deiner</i>
PPOSAT	attributive possessive pronoun	<i>mein [Buch], deine [Mutter]</i>
PRELS	substituting relative pronoun	<i>[der Hund,] der</i>
PRELAT	attributive relative pronoun	<i>[der Mann ,] dessen [Hund]</i>
PRF	reflexive personal pronoun	<i>sich, einander, dich, mir</i>
PWS	substituting interrogative pronoun	<i>wer, was</i>
PWAT	attributive interrogative pronoun	<i>welche [Farbe], wessen [Hut]</i>
PWAV	adverbial interrogative or relative pronoun	<i>warum, wo, wann, worüber, wobei</i>
PROP	pronominal adverb	<i>dafür, dabei, deswegen, trotzdem</i>

POS =	description	examples
PTKZU	<i>zu</i> + infinitive	<i>zu [gehen]</i>
PTKNEG	negation particle	<i>nicht</i>
PTKVZ	separated verb particle	<i>[er kommt] an, [er fährt] rad</i>
PTKANT	answer particle	<i>ja, nein, danke, bitte</i>
PTKA	particle with adjective or adverb	<i>am [schönsten], zu [schnell]</i>
TRUNC	truncated word - first part	<i>An- [und Abreise]</i>
VVFIN	finite main verb	<i>[du] gehst, [wir] kommen [an]</i>
VVIMP	imperative, main verb	<i>komm [!]</i>
VVINFINF	infinitive, main	<i>gehen, ankommen</i>
VVIZU	infinitive + <i>zu</i> , main	<i>anzukommen, loszulassen</i>
VVPP	past participle, main	<i>gegangen, angekommen</i>
VAFIN	finite verb, aux	<i>[du] bist, [wir] werden</i>
VAIMP	imperative, aux	<i>sei [ruhig !]</i>
VAINFINF	infinitive, aux	<i>werden, sein</i>
VAPP	past participle, aux	<i>gewesen</i>
VMFIN	finite verb, modal	<i>dürfen</i>
VMINFINF	infinitive, modal	<i>wollen</i>
VMPP	past participle, modal	<i>[er hat] gekonnt</i>
XY	non-word containing special characters	<i>D2XW3, letters</i>
\$,	comma	<i>,</i>
\$.	sentence-final punctuation	<i>. ? ! ; :</i>
\$(other sentence internal punctuation	<i>- [] ()</i>

The following POS categories do not contain any morphological information and are assigned the morphological label ”- -”: ADJD, ADV, APZR, CARD, FM, ITJ, KOUJ, KOUS, KON, KOKOM, PWAV, PROP, PTKZU, PTKNEG, PTKVZ, PTKANT, PTKA, TRUNC, VVIZU, VVPP, VAPP, VMPP, XY, \$, , \$. , \$(.

3.4.3 What Is a Syntactic Unit?

The newspaper articles of the taz have been defined as the primary segmentation domain of the data. They are preprocessed into syntactic units delimited by punctuation marks (. ? ! ; - ... /) for which specific rules demand or forbid segmentation. Each syntactic unit is assigned a specific code which identifies its origin in the newspaper data, e.g. T990507.123 (T (taz) 99 (year) 05 (month) 07 (day) 123 (article)).

A syntactic unit usually consists of one complete sentence structure with a root node (SIMPX, R-SIMPX, P-SIMPX). But it may also consist of one or more sentences and/or phrases, e.g. headlines, titles, sentences with parentheses, sentences with discourse markers, or sentence conjunction by a colon.

An annotated tree is a complete syntactically and semantically well-formed construction according to the *longest match principle*. The model of topological fields does not prescribe that all fields have to be occupied. The fact that fields can be left empty, also helps us to cope with elliptical constructions (cf. 6.6).

Table 3.5: Morphological feature combinations for lexical elements

POS	feature combination	comments
ADJA	case number gender	underspecified for gender if the plural noun is underspecified, i.e. the plural noun does not morphologically represent its gender, e.g. deadjectival nouns: <i>die/np* nordhessischen/np* Grünen/np*</i> invariant local description e.g. <i>Berliner/***</i> cardinal numbers as abbreviation: full morphology e.g. <i>im 4./dsn Jahrhundert/dsn</i>
APPR	case	without case if a preposition takes another PP as complement, e.g. <i>bis/_ zu/d einer/dsf Woche/dsf</i> and in the construction <i>was für ein(er/e/...)</i>
APPRART	case number gender	
APPO	case	
ART	case number gender	
NN	case number gender	can be underspecified for gender, e.g. deadjectival nouns (<i>Abgeordnete</i> (in plural)) or pluralia tantum (<i>Leute</i>)
NE	case number gender	
PDS	case number gender	
PDAT	case number gender	
PIS	case number gender	underspecified: <i>man/ns*</i> <i>nichts/***</i> (cf. <i>nix, sowas</i>) PIS or PIAT: <i>allerhand/***</i> (cf. <i>allerlei, allzuviel, dergleichen, derlei, etwas, genausoviel, genug, genügend, keinerlei, mehr, reichlich, soviel, viel, wenig, weniger, zuviel, zuwenig</i>) PIDAT or PIS: <i>sowas/***</i> (cf. <i>paar, bißchen</i>)

POS	feature combination	comments
PIAT	case number gender	plural is underspecified for gender, e.g. <i>lauter</i> /***, see also 'PIS or PIAT' below
PIDAT	case number gender	<i>solch</i> /*** (cf. <i>manch</i> , <i>welch</i> , <i>all</i>), see also 'PIS or PIDAT' below
PPER	case number gender person	
PPOSS	case number gender	
PPOSAT	case number gender	
PRELS	case number gender	plural is underspecified for gender
PRELAT	case number gender	
PRF	case number gender person	<i>sich</i> : underspecified for gender
PWS	case number gender	underspecified for gender: plural forms and <i>wer</i> , <i>wem</i> , <i>wen</i>
PWAT	case number gender	<i>wessen</i> /***
VAFIN	person number mood tense	
VAIMP	number	
VMFIN	person number mood tense	
VVFIN	person number mood tense	
VVIMP	number	German has only second person imperative forms

Table 3.6: Values of morphological features

Feature	Values
case	n (nominative), g (genitive), d (dative), a (accusative), * (underspecified)
gender	m (masculine), f (feminine), n (neuter), * (underspecified)
number	s (singular), p (plural), * (underspecified)
mood	i (indicative), k (subjunctive; German 'Konjunktiv')
person	1 (first), 2 (second), 3 (third), * (underspecified)
tense	s (present), t (past), * (underspecified)

Table 3.7: Node labels

Node Labels	Description
Phrase Node Labels	
ADJX	adjectival phrase
ADVX	adverbial phrase
DP	determiner phrase (e.g. <i>gar keine</i>)
FX	foreign language phrase
NX	noun phrase
PX	prepositional phrase
VXFIN	finite verb phrase
VXINF	non-finite verb phrase
Topological Field Node Labels	
LV	resumptive construction (Linksversetzung)
C	complementizer field (C-Feld)
FKOORD	coordination consisting of conjuncts of fields
KOORD	field for coordinating particles
LK	left sentence bracket (Linke (Satz-)Klammer)
MF	middle field (Mittelfeld)
MFE	middle field between VCE and VC
NF	final field (Nachfeld)
PARORD	field for non-coordinating particles
VC	verb complex (Verbkomplex)
VCE	verb complex with the split finite verb of <i>Ersatzinfinitiv</i> constructions
VF	initial field (Vorfeld)
FKONJ	conjunct consisting of more than one field
Root Node Labels	
DM	discourse marker
P-SIMPX	paratactic construction of simplex clauses
R-SIMPX	relative clause
SIMPX	simplex clause

Table 3.8: Edge labels

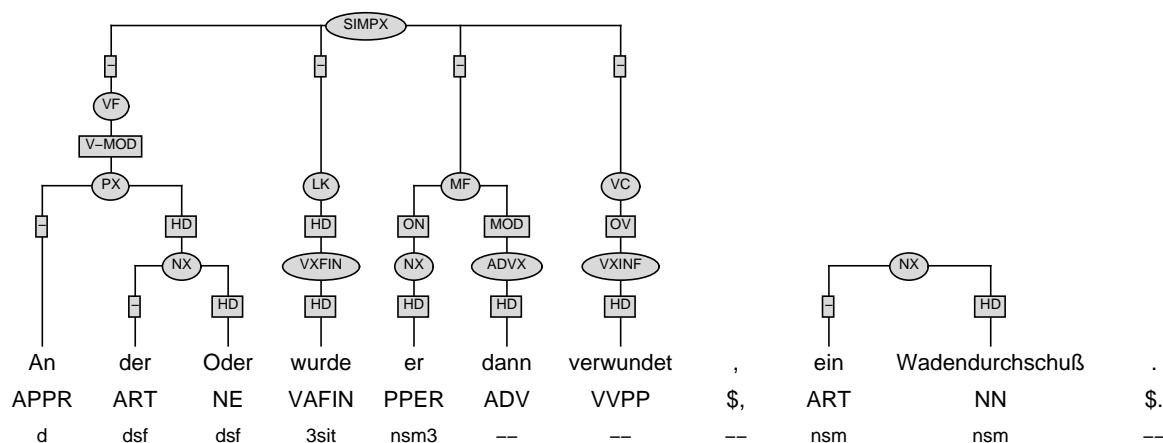
Edge Labels	Description
Edge Labels denoting Heads and Conjuncts	
HD	head
-	non-head
KONJ	conjunct
Complement Edge Labels	
ON	nominative object (i.e. subject; also clausal subjects)
OD	dative object
OA	accusative object
OG	genitive object
OS	sentential object
OPP	prepositional object
OADV	adverbial object
OADJP	adjectival object
PRED	predicate
OV	verbal object
FOPP	facultative (i.e. optional) prepositional object, passivized subject (<i>von</i> -phrase)
VPT	separable verb prefix
APP	apposition
Modifier Edge Labels	
MOD	ambiguous modifier
ON-MOD, OA-MOD, OD-MOD, OG-MOD, OS-MOD, OPP-MOD, FOPP-MOD, PRED-MOD, OADJP-MO, OADV-MO, V-MOD, MOD-MOD	modifiers modifying complements or modifiers, e.g. V-MOD = modifier of the verb
Edge Labels in Split Coordinations	
ONK, OAK, ODK, OGK, OPPK, FOPPK, OSK, OADVPK, OA-MODK, MODK, V-MODK	second conjunct (K) in PREDK, split coordinations e.g. ONK = second conjunct of a nominative object
Edge Label denoting Structural Expletive	
ES	Vorfeld- <i>es</i>
Secondary Edge Labels	
refvc	dependency relation between: two verbal objects in VC
refmod	two ambiguous modifiers
refint	a phrase internal part and its modifier
refcontr	control verb and its complement across clause boundaries

Table 3.9: Syntactic-Semantic Node Labels for Named Entities

Labels	Description
Syntactic-Semantic Node Labels	
ADJX=ORG	adjectival phrase, named entity of the semantic class “organisation”
ADJX=OTH	adjectival phrase, named entity of the semantic class “other”
ADVX=ORG	adverbial phrase, named entity of the semantic class “organisation”
ADVX=OTH	adverbial phrase, named entity of the semantic class “other”
DM=OTH	discourse marker, named entity of the semantic class “other”
FX=LOC	foreign language phrase, named entity of the semantic class “location”
FX=ORG	foreign language phrase, named entity of the semantic class “organisation”
FX=OTH	foreign language phrase, named entity of the semantic class “other”
FX=PER	foreign language phrase, named entity of the semantic class “person”
NX=GPE	noun phrase, named entity of the semantic class “geopolitical entity”
NX=LOC	noun phrase, named entity of the semantic class “location”
NX=ORG	noun phrase, named entity of the semantic class “organisation”
NX=OTH	noun phrase, named entity of the semantic class “other”
NX=PER	noun phrase, named entity of the semantic class “person”
PX=GPE	prepositional phrase, named entity of the semantic class “geopolitical entity”
PX=LOC	prepositional phrase, named entity of the semantic class “location”
PX=ORG	prepositional phrase, named entity of the semantic class “organisation”
PX=OTH	prepositional phrase, named entity of the semantic class “other”
PX=PER	prepositional phrase, named entity of the semantic class “person”
SIMPX=ORG	simplex clause, named entity of the semantic class “organisation”
SIMPX=OTH	simplex clause, named entity of the semantic class “other”
VXINF=ORG	non-finite verb phrase, named entity of the semantic class “organisation”
VXINF=OTH	non-finite verb phrase, named entity of the semantic class “other”
Edge Label	
-NE	non-head, the part below is not part of the named entity

Punctuation is not annotated, i.e., all punctuation marks are not attached to the tree structure. Exceptions are punctuation marks which carry a semantic meaning within a sentence, e.g. - (*bis, und*) in expressions like *15.30 - 17.30 Uhr*. They are tagged according to the part of speech that they represent in the text (cf. 4.4.1).

Constituents are not attached to a tree if they are not assigned a grammatical function within the specific syntactic construction. The following tree diagram shows two annotated trees in one syntactic unit:⁵



The leaves of the trees consist of pairs of non-terminal symbols and part-of-speech tags. Non-terminal symbols are represented by spherical nodes, whereas edge labels are depicted by rectangular nodes. The tree diagram consists of two trees, a SIMPX and an isolated phrase. In accordance with the four annotation levels shown in Table 3.3, the sentence is annotated top-down by the root node (SIMPX), the field nodes (VF, LK, MF, and VC), the phrase nodes (PX, VXFIN, NX, ADVX, and VXINF), and finally the tagged lexical entries. The edge labels between the field level and the phrase level indicate that the syntactic structure contains one unambiguous modifier (V-MOD), a subject (ON), one ambiguous modifier (MOD), a verbal object (OV), and the finite verb, which itself is the head (HD) of the entire syntactic construction. The noun phrase (*ein Wadendurchschuß*) is not attached to the sentence structure because otherwise the well-formedness of the construction would be violated. Thus, it has to be annotated as an isolated phrase lacking a verbal constituent.

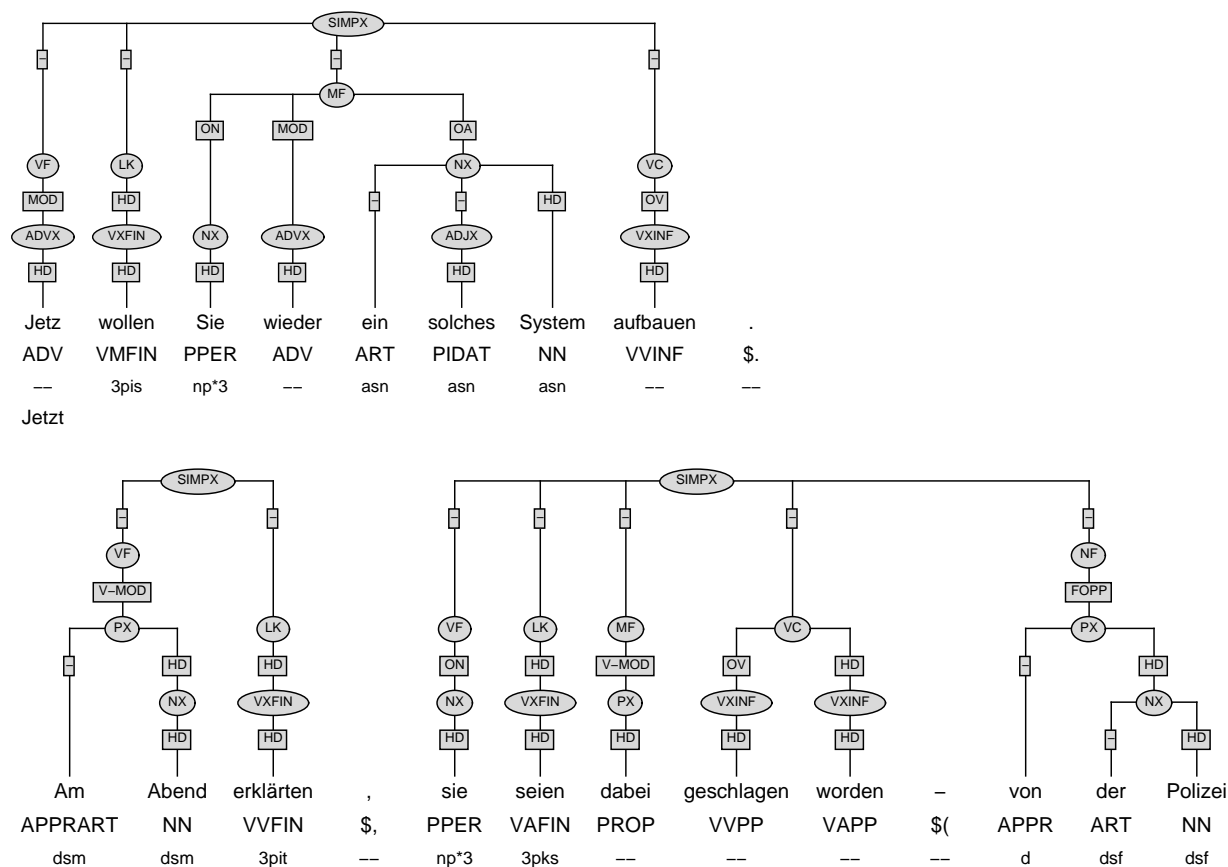
3.4.4 Printing and Spelling Errors

In contrast to spoken language data like in the VERBMOBIL (cf. (Stegmann et al. 2000)) which exhibit fragmentary utterances, false starts, repetitions, interruptions, and hesitation noises as its characteristic properties, data taken from newspaper corpora does not include unintentionally formed syntactic constructions.

Deviations from syntactic wellformedness are either intended by the author or are caused by printing errors. While incorrect writing of words is neglected in the syntactic

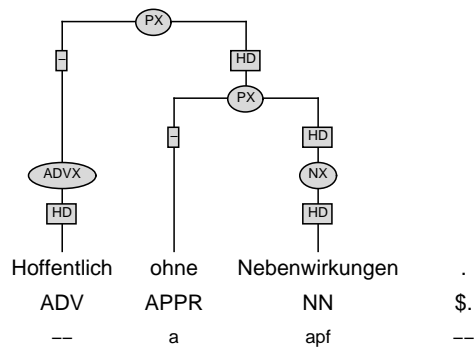
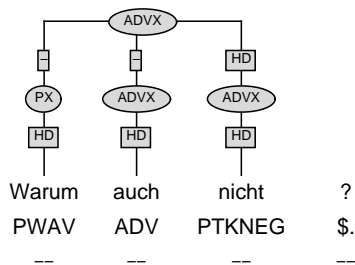
⁵These tree diagrams and all following tree diagrams in this report were generated with the aid of the *Negra Annotate* tool.

analysis (the respective lexical entry is marked with the correct writing of the word in a comment line below), lexical elements which do not belong to the syntactic construction (intentional or unintentional) are structured as much as possible, but are not attached to the surrounding constituents:



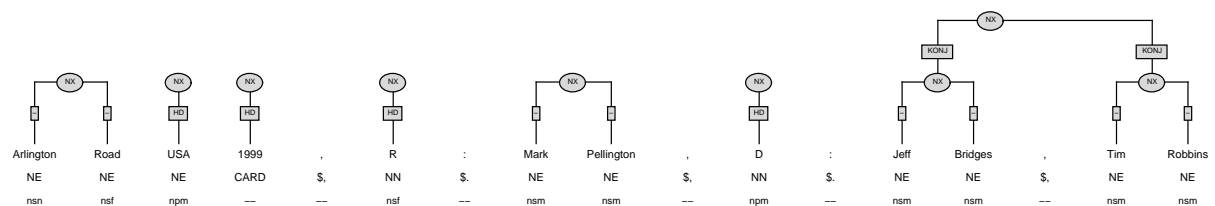
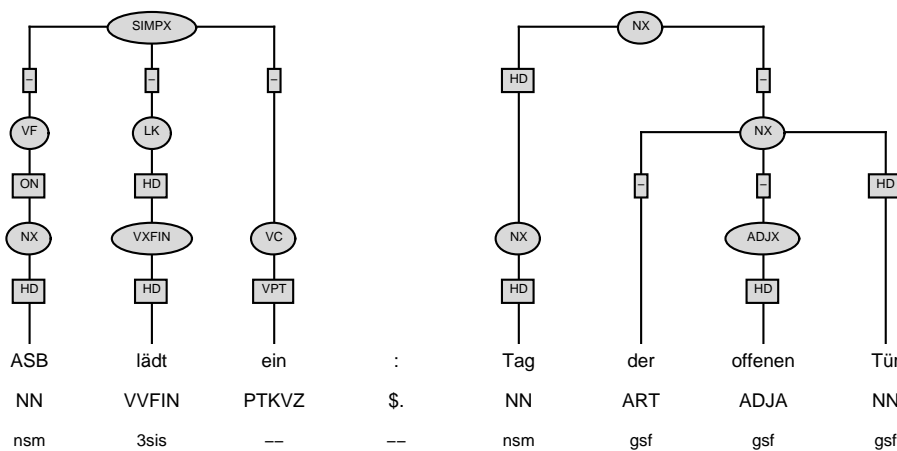
3.4.5 Isolated Phrases

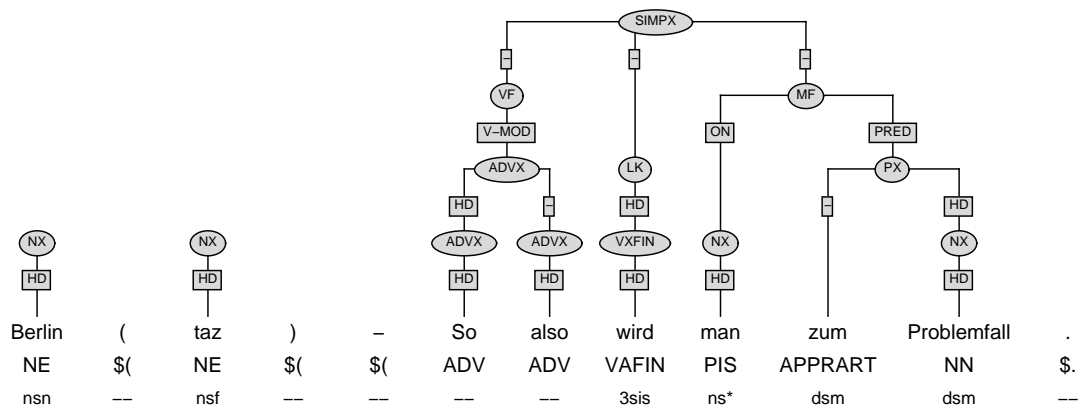
There are textual fragments in newspaper data which cannot be analysed as a SIMPX or as a constituent of a SIMPX because they are lacking a verbal constituent or they are not assigned a specific grammatical function within a well-formed sentence. These fragments are annotated as isolated phrases. The isolated elements are structured as much as possible (mostly up to the level of phrasal categories), but they are not typically connected to surrounding constituents as a whole, so that a conflict with the topological field analysis is avoided. Their root node carries a phrasal category of their lexical head (NX, PX, ADVX, etc.):



In accordance with the *longest match principle*, as many parts of the fragment as possible are projected to the phrase level and are included into a tree structure. It has to be decided which part of the whole construction is the head and which parts depend on this head.

Phrases within a syntactic unit are not attached on a higher level if they do not show dependency relation. This is often the case with syntactic elements which are separated by a colon or a dash (cf. 5.3.2):



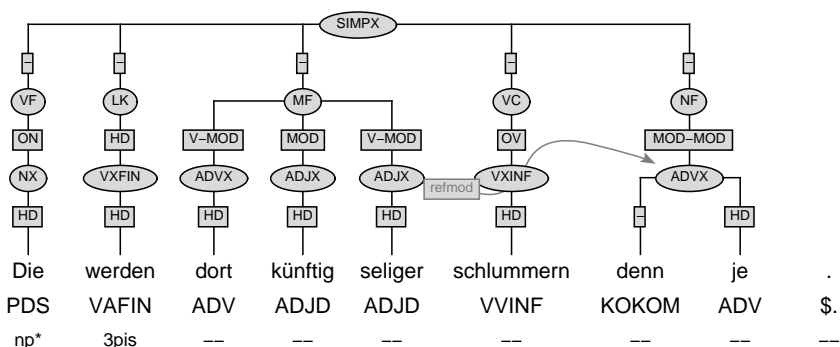


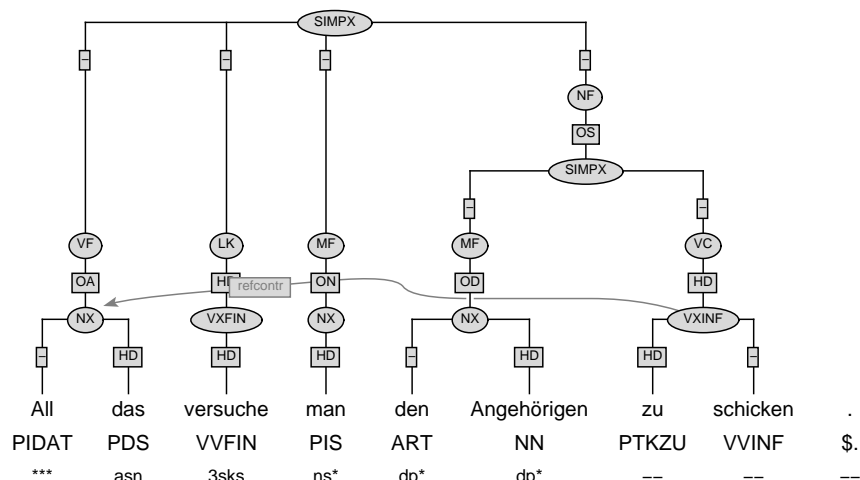
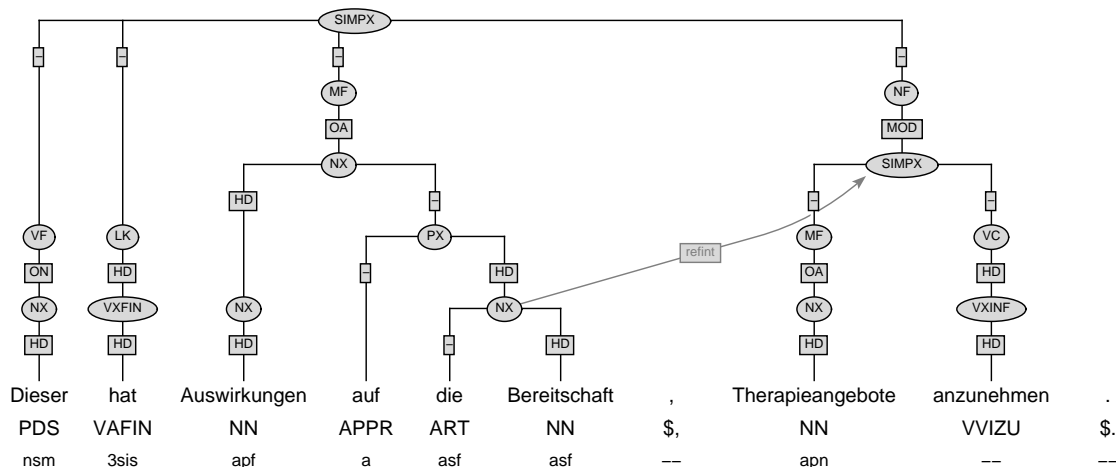
3.4.6 Long-Distance Dependencies

Our annotation scheme facilitates a surface-oriented representation of long-distance dependencies without crossing branches and traces. If a modifying constituent is not adjacent to the modified constituent, their dependency relation, which can even go beyond the border of topological fields, is encoded by special naming conventions for edge labels. We use edge labels such as OA-MOD (referring to OA) or PRED-MOD (referring to PRED) etc. expressing the non-ambiguity of the modifier.

Beyond this, we make use of secondary edge labels for ambiguity resolution. These labels just serve as additional information to the grammatical functions encoded in the edge labels. These secondary edge labels indicate underspecified long distance dependencies in the following cases:

1. If the above mentioned edge labels need further disambiguation, e.g. if there are two OAs or V-MODs below one SIMPX node (refmod).
2. If the dependency relation exists between two nodes of which at least one is phrase internal and therefore carries only head or non-head information (refint).
3. If there is a dependency relation outside of SIMPX in control verb constructions (refcontr).





3.4.7 Empty Categories

In general, an empty category analysis, e.g. for phrases without heads, is being avoided in the TüBa-D/Z treebank.

Empty Edge Labels

Specifiers, prepositions,⁶ complementizers, discourse markers, KOORD and PARORD constituents, conjunctions, and unambiguous modifiers (that are attached to phrases immediately rather than to topological fields) are not labelled with grammatical functions. Furthermore, the edges below the SIMPX node are empty. They are not labelled in order to speed up annotation where the information is unnecessary or self-evident.

Furthermore, empty edge labels are used in elliptical phrases, e.g. noun phrases only consisting of an article and an attributive adjective (cf. 6.6).

⁶In order to facilitate the identification of dependencies between verbs and their nominal complements and adjuncts and in keeping with basic assumptions in Dependency Grammar, the annotated head of a prepositional phrase is the NX (or complement) rather than the preposition itself. Therefore, prepositions carry no edge label.

3.5 Lemma Information

The trees in the TüBa-D/Z are enriched with lemma information for all tokens. Morphology and lemmatization are performed by an automatic pre-tagging, which makes use of the existing syntactic annotation of the treebank. The output of this pre-tagging is manually disambiguated and corrected. For a detailed description of the pre-tagging system see Versley et al. (2010); for an overview of lemmatization problems see Schnorr (1991).

3.5.1 Lemmatization Rules for POS-Tags

In the following Table 3.10, the lemmatization rules applied for open-class words (e.g. nouns, adjectives) and closed-class words (e.g. determiners, pronouns) in the TüBa-D/Z are described with respect to the STTS POS-tag of the token.

Table 3.10: Lemmatization rules for POS-tags

POS-tag	lemmatization rule	examples
ADJA ADJD	<p>base form: mapping to the predicative form exceptions:</p> <p>comparative: mapping to the comparative form for attributive adjective adverbial adjective predicative adjective</p> <p>superlative: stem without ending for attributive adjective</p> <p>deverbal adjective: mapping to the predicative form</p>	<p><i>(der) hohe (Anteil) → hoch</i> <i>(das ist) gut → gut</i></p> <p><i>besondere (Sorgfalt) → besonder</i> <i>andere (Menschen) → ander</i></p> <p><i>bessere (Chancen) → besser</i> <i>(es dauert) länger → länger</i> <i>(es sei) besser → besser</i></p> <p><i>(der) schnellste (Schwimmer) → schnellst</i></p> <p><i>gespannt, zerstritten, brennend</i></p>
ADV	invariant form	<i>schon, bald, doch</i>
APPR APPO APZR	invariant form	<i>in</i> <i>zufolge</i> <i>an</i>
APPRART	reduced to preposition	<i>im → in</i> <i>zur → zu</i>

POS-tag	lemmatization rule	examples
ART	<p>base form: nom/sg definite article (sg/pl): lemmata: <i>der, die, das</i></p> <p>indefinite article (sg): lemmata: <i>ein, eine</i> plural: zero article</p>	<p>masc.: <i>der, des, dem, den, die</i> → <i>der</i> fem.: <i>die, der, den</i> → <i>die</i> neut.: <i>das, des, dem, die, der, den</i> → <i>das</i></p> <p>masc.: <i>ein, eines, einem, einen</i> → <i>ein</i> fem.: <i>eine, einer</i> → <i>eine</i> neut.: <i>ein, eines, einem</i> → <i>ein</i></p>
CARD	invariant form	<i>zwei, 2, 10.000</i>
ITJ	invariant form	<i>hallo, aha, hey</i>
KOUI	invariant form	<i>um</i>
KOUS	invariant form	<i>weil</i>
KON	invariant form	<i>und</i>
KOKOM	invariant form	<i>als, wie</i>
NE	base form: nom/sg	<i>Hans</i> → <i>Hans</i> <i>Bremerhavens</i> → <i>Bremerhaven</i>
NN	<p>base form: nom/sg</p> <p>gender remains unchanged</p> <p>deadjectival nouns: lemmatized to the form of the strong declension of adjectives in German</p> <p>deverbal nouns:</p> <p>plural nouns: base form nom/sg if a singular form exists</p> <p>base form nom/pl if a singular form does not exist (pluralia tantum)</p> <p>homonyms and polysemes keep their base form</p> <p>compounds are not split</p>	<p><i>Schränke</i> → <i>Schrank</i> <i>Ideen</i> → <i>Idee</i> <i>Lehrerin</i> → <i>Lehrerin</i> <i>Kaufmann</i> → <i>Kaufmann</i></p> <p>masc.: <i>(der) Schöne</i> → <i>Schöner</i> fem.: <i>(die) Schöne</i> → <i>Schöne</i> neutr.: <i>(das) Schöne</i> → <i>Schönes</i></p> <p><i>(das) Reisen</i> → <i>Reisen</i></p> <p><i>Daten</i> → <i>Datum</i> <i>Medien</i> → <i>Medium</i></p> <p><i>Leuten</i> → <i>Leute</i></p> <p><i>Schlösser</i> → <i>Schloß</i> <i>Flügeln</i> → <i>Flügel</i></p> <p><i>EU-Kommissar</i> → <i>EU-Kommissar</i> <i>Senioren-Bahncard</i> → <i>Senioren-Bahncard</i></p>

POS-tag	lemmatization rule	examples
PDS PDAT	base form: nom/sg one lemma each for masc., fem., neut.	masc.: <i>dieser/dieses/diesem/diesen</i> → <i>dieser</i> fem.: <i>jene/jener</i> → <i>jene</i> neut.: <i>das/dem/den</i> → <i>das</i>
PIS PIAT PIDAT	base form: nom/sg one lemma each for masc., fem., neut. or one general lemma	masc.: <i>keiner/keinen/keinem</i> → <i>keiner</i> fem.: <i>letztere/letzteren</i> → <i>letztere</i> neutr.: <i>jedes/jedem</i> → <i>jedes</i> <i>beiden</i> → <i>beide</i> <i>allen</i> → <i>alle</i> <i>man</i> → <i>man</i>
PPER	base form: nom/sg polite form	<i>ich/meiner/mir/mich</i> → <i>ich</i> <i>du/deiner/dir/dich</i> → <i>du</i> <i>er/seiner/ihm/ihn</i> → <i>er</i> <i>sie/ihrer/ihr/sie</i> → <i>sie</i> <i>es/seiner/ihm/es</i> → <i>es</i> <i>wir/uns/uns/uns</i> → <i>wir</i> <i>ihr/euer/euch/euch</i> → <i>ihr</i> <i>sie/ihrer/ihnen/sie</i> → <i>sie</i> <i>Sie/Ihrer/Ihnen/Sie</i> → <i>Sie</i>
PPOSS	base form: nom/sg lemma according to the gender of the possession	masc.: <i>meiner/meiner/meinem/meinen</i> → <i>meiner</i> fem.: <i>meine/meiner/meiner/meine</i> → <i>meine</i> neutr.: <i>mein(e)s/meine(e)s/meinem</i> <i>/mein(e)s</i> → <i>mein(e)s</i>
PPOSAT	base form: nom/sg lemma according to the gender of the possession	masc.: <i>mein/meines/meinem/meinen</i> → <i>mein</i> fem.: <i>meine/meiner/meiner/meine</i> → <i>meine</i> neutr.: <i>mein/meines/meinem/mein</i> → <i>mein</i>
PRELS	base form: nom/sg	<i>der, dessen, den, dem</i> → <i>der</i> <i>die, derer, der, die</i> → <i>die</i> <i>das, dessen, dem, das</i> → <i>das</i>
PRELAT	base form: nom/sg	masc./neut.: <i>dessen</i> → <i>dessen</i> fem.: <i>deren</i> → <i>deren</i>
PRF	reflexive pronouns are labeled as #refl	<i>mir/mich</i> → #refl <i>dir/dich</i> → #refl <i>sich</i> → #refl <i>uns</i> → #refl <i>euch</i> → #refl

POS-tag	lemmatization rule	examples
PWS PWAT	base form: nom/sg invariant form	<i>wer, wessen, wem, wen</i> → <i>wer</i> masc.: <i>welcher, welchen, welchem</i> <i>welchen</i> → <i>welcher</i> fem.: <i>welche, welcher, welcher, welche</i> → <i>welche</i> neut.: <i>welches, welchen, welchem,</i> <i>welches</i> → <i>welches</i> <i>was</i>
PWAV	invariant form	<i>wo, wie, warum, womit, worauf</i>
PROP	invariant form	<i>damit, davor, seitdem, stattdessen</i>
PTKA PTKANT PTKNEG PTKZU	invariant form	<i>am</i> <i>ja</i> <i>nicht</i> <i>zu</i>
PTKVZ	no lemma for verb particles, the lemma of the verb is represented as particle#verb (see Table 3.11)	<i>ein</i> → - - <i>warf</i> → <i>ein#werfen</i> (<i>er warf etwas ein</i>)
TRUNC	lemma is the complete word suffixed with %n, %v, %a, %c, %p for the respective part of speech	<i>In- und Ausland</i> → <i>Inland%n und Ausland</i> <i>hin- und herzieht</i> → <i>hinziehen%v und herziehen</i>
VVFIN VVIMP VVINP VVIZU VVPP VAFIN VAIMP VAINP VAPP VMFIN VMINP VMPP	base form: infinitive See Table 3.11 for auxiliary and passive use (%aux, %passiv).	<i>ging</i> → <i>gehen</i> <i>sprich</i> → <i>sprechen</i> <i>zahlen</i> → <i>zahlen</i> <i>aufzufallen</i> → <i>auf#fallen</i> <i>getroffen</i> → <i>treffen</i> <i>ist</i> → <i>sein%aux, ist</i> → <i>sein</i> <i>seid</i> → <i>sein%aux, seid</i> → <i>sein</i> <i>haben</i> → <i>haben%aux, haben</i> → <i>haben</i> <i>gewesen</i> → <i>sein%aux, gewesen</i> → <i>sein</i> <i>will</i> → <i>wollen%aux, will</i> → <i>wollen</i> <i>möge</i> → <i>mögen%aux, möge</i> → <i>mögen</i> <i>gekonnt</i> → <i>können</i>
FM	foreign language material is invariant	<i>ad hoc, goes, areas</i>
XY	non-words are invariant, lemmata in lower-case letters	<i>18a</i> → <i>18a</i> <i>H2O</i> → <i>h2o</i>
\$(invariant form	, . ? ... (

3.5.2 Lemmatization Rules for Specific Linguistic Phenomena

The following Table 3.11 describes the lemmatization rules applied for specific linguistic phenomena in the TüBa-D/Z.

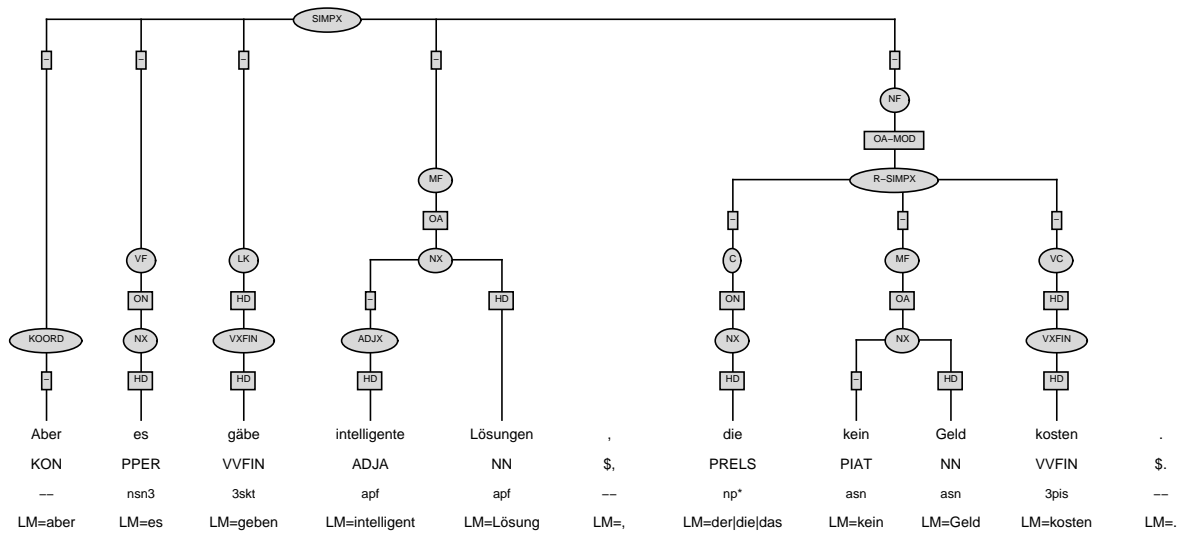
Table 3.11: Lemmatization rules for specific linguistic phenomena

phenomenon	lemmatization rule	examples
abbreviation or acronym	abbreviations and acronyms are invariant	<i>z. B., usw., Dr.</i> <i>TSV, FDP</i>
spelling errors	mapping to the correct spelling of the lemma	<i>wolte</i> → <i>wollte</i> <i>Durchamtmen</i> → <i>Durchatmen</i>
multiword term	one lemma for each multiword token	<i>New York</i> → <i>New York</i> <i>Orang Utan</i> → <i>Orang Utan</i>
dialect	the base form of dialect words is the respective standard German word with an underscore appended	<i>es jütt</i> → <i>es geben_</i> <i>snakt</i> → <i>sprechen_</i> <i>Dag</i> → <i>Tag_</i>
contraction of words	mapping to a complex lemma with an underscore between the base forms of the contraction parts exception: APPRART reduced to the preposition	<i>Glaubense</i> → <i>glauben_Sie_</i> <i>isser</i> → <i>sein_er</i> <i>zur</i> → <i>zu</i>
non-standard use of lower-case and upper-case letters	mapping to the correct writing of the lemma based on German orthography polite form with upper-case letters	<i>seele</i> → <i>Seele</i> <i>KOMMENTAR</i> → <i>Kommentar</i> <i>Sie</i> → <i>Sie</i>
spelling variations	are annotated as distinct lemmata	<i>fantastische</i> → <i>fantastisch</i> <i>phantastische</i> → <i>phantastisch</i>

phenomenon	lemmatization rule	examples
ambiguous plural forms	for plurals unmarked for gender, all possible lemmata are listed separated by a diacritic ‚ ‘, e.g. lemmata of deadjectival plural nouns or plural pronouns with underspecified gender	<i>Jugendliche</i> → <i>Jugendlicher Jugendliche Jugendliches</i> <i>die</i> (PDS np*) → <i>der die das</i> <i>denen</i> (PRELS dp*) → <i>der die das</i>
auxiliaries: <i>sein, haben, werden</i> modal verbs: <i>müssen, sollen, können, wollen, dürfen, mögen</i>	the lemma is suffixed with the tag %aux if used as auxiliary	<i>ist</i> → <i>sein%aux</i> <i>darf</i> → <i>dürfen%aux</i>
auxiliaries and modal verbs used as main verbs	base form: infinitive without %aux suffix	<i>ist</i> → <i>sein</i> (... <i>es ist höchste Zeit</i> ...) <i>kann</i> → <i>können</i> (... <i>wer kann das überhaupt noch</i> ...)
passive werden	the lemma is suffixed with the tag %passiv	<i>wird (geehrt)</i> → <i>werden%passiv</i>
verbs with a separable prefix	the verb lemma is denoted as prefix#verb , whether the prefix is separated or not (See Table 3.10 for verb particles (PTKVZ))	<i>stellen ... ein</i> → <i>ein#stellen</i> <i>eingestellt</i> → <i>ein#stellen</i>

The following tree diagram illustrates the TüBa-D/Z lemma annotation below the morphological feature combinations marked as "LM=lemma" for each token of the sentence:

Aber es gäbe intelligente Lösungen, die kein Geld kosten.



pronouns, nouns, determiners (base form nom/sg):
 LM=es, LM=Lösung, LM=kein, LM=Geld, LM=der|die|das

verbs (base form infinitive):
 LM=geben, LM=kosten

adjective (base form predicate):
 LM=intelligent

conjunction, punctuation marks (invariant):
 LM=aber, LM=, LM=.

Chapter 4

The Annotation of the Internal Structure of Phrases

4.1 Premodification and Postmodification in Phrases

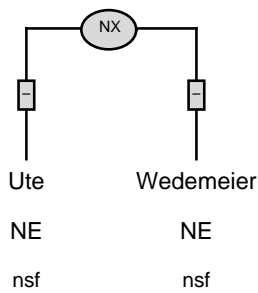
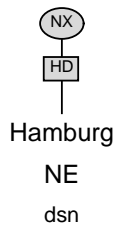
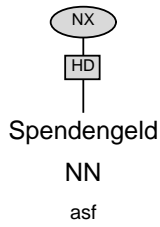
The annotation of phrases is also carried out following the *flat clustering principle* in order to keep the number of hierarchy levels in a syntactic structure as small as possible. As will be shown in the following sections, phrases may include adjectival or nominal premodifiers and/or postmodifiers of any syntactic category. Both kinds of modifiers are in principle projected to their phrase levels. Since the modification scope of premodifiers is unambiguous, they are directly attached to the head of the phrase which they modify. By contrast, postmodifiers are always attached on a higher level to preserve ambiguity. This decision, referred to in 3.3 as the *high attachment principle*, was made to avoid the problematic distinction whether a postmodifier is a free adjunct or a complement of the modified phrase. The attachment strategy for premodifiers and postmodifiers is applied for all categories of phrases.

4.2 Noun Phrases

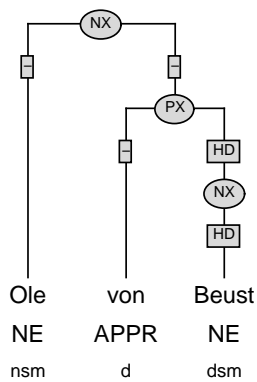
A simple noun phrase (NX) consists of a head noun (noun, proper noun, or a pronoun), (optionally) a determiner and (optionally) an adjectival or a nominal premodifier of any complexity preceding the head noun. A complex noun phrase is a simple noun phrase with a postmodifier of any syntactic category and complexity.

4.2.1 Noun Phrases without Modifiers

Simple noun phrases without modifiers are single nouns, proper nouns, pronouns or proper nouns consisting of more than one NE. All of them are directly projected to their phrase level. While single nouns, proper nouns and pronouns carry the edge label HD, the NE-tagged tokens of a complex proper noun are attached on the same level without head information:

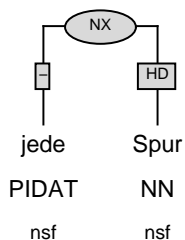
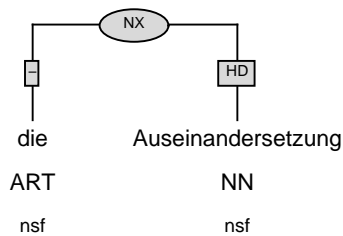


If proper nouns include other parts of speech than NEs, these parts are tagged according to their distribution. Therefore, proper nouns with a preposition include a prepositional phrase.

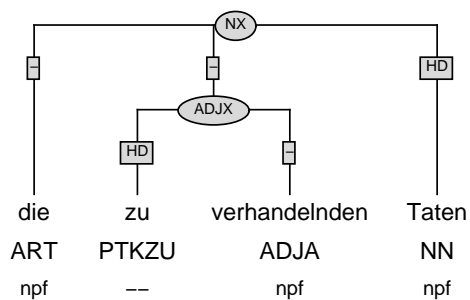
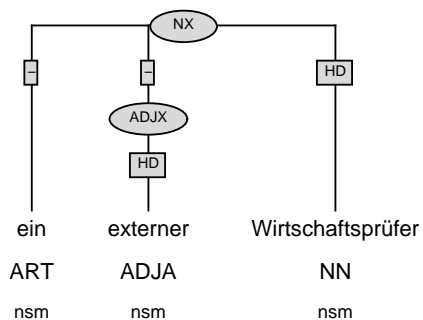


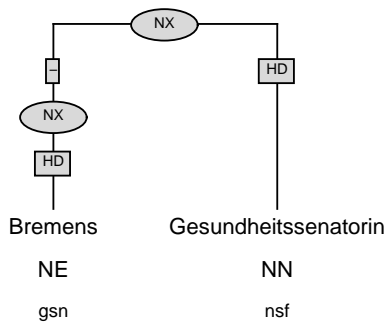
4.2.2 Prenominal Modification

In a simple noun phrase, both the determiner and the head noun are directly attached on the same level to NX so that the label of the head noun carries the edge label HD and the edge label of the determiner is empty.

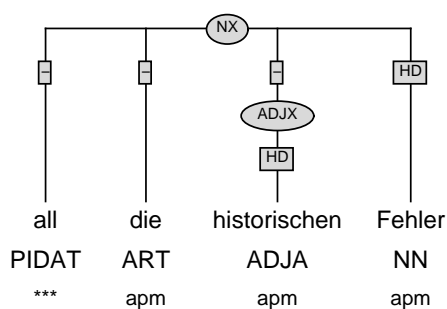


Since prenominal modifiers are directly attached to the head noun on the same level, their edge labels are empty (whereas the edge labels of modifiers that are attached to topological fields are non-empty (cf. 8.4)). Prenominal modifiers are either attributive adjectives or preceding genitive phrases:

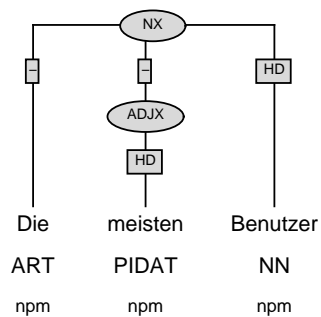


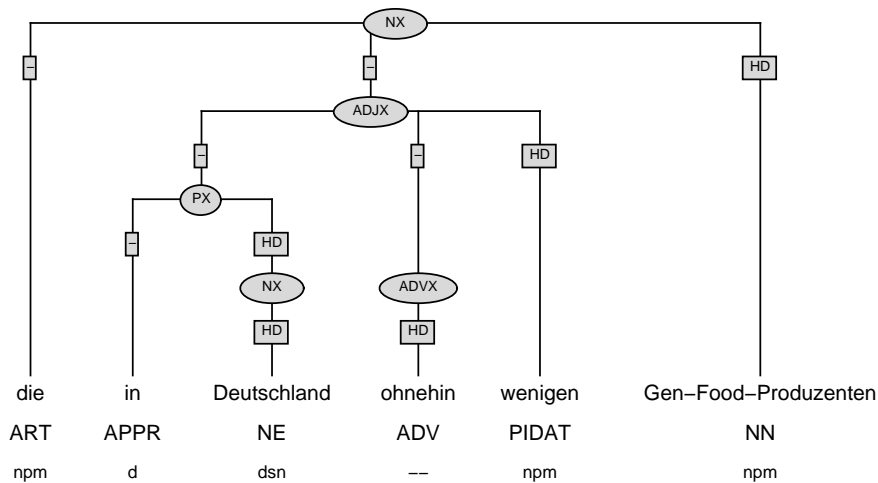


If there is a PIDAT preceding the article it is directly attached to the noun phrase.

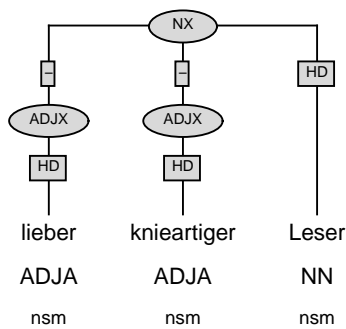


If a PIDAT is following the article in adjective position it is projected to its phrase level (ADJX) with possible premodifiers and then directly attached like an attributive adjective to the noun phrase.

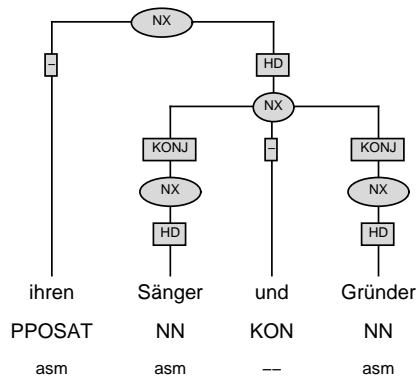




If there is more than one prenominal modifier, the one on the left hand side of the noun is modifying the following noun, the one on the left hand side of the modifier is modifying both, the modifier and the noun, and so on. All of these modifiers are attached to the head noun on the same level which yields a rather flat noun phrase structure. This strategy is justified by the fact that these modifiers have a scope of modification beyond the adjectival phrase, e.g. as in coordinated noun phrases like *insgesamt 12.000 Studienplätze und 15.000 Lehrstellen*, the adverb *insgesamt* modifies *12.000 Studienplätze* as well as *15.000 Lehrstellen*.

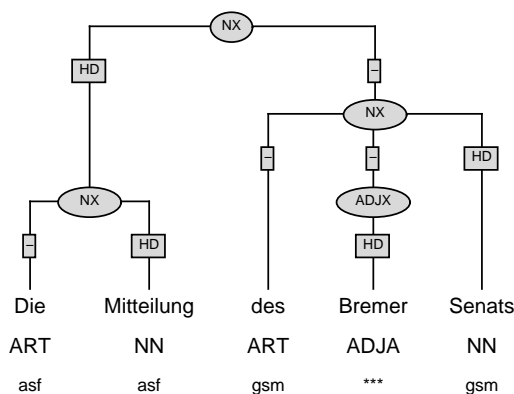
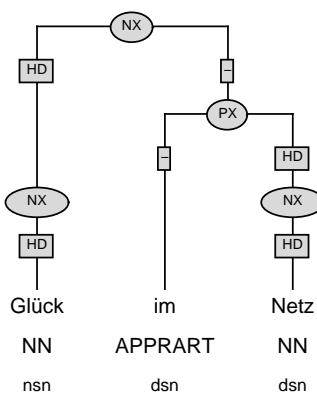


In case of complex head nouns, e.g. complex (proper) nouns consisting of two nominal parts or coordinated head nouns (cf. 6.5.5), first the complex noun respectively the coordination (cf. 6.5) is annotated with its own internal dependency structure. Afterwards, the determiner and possible premodifying adjectival phrases are attached on a higher level.



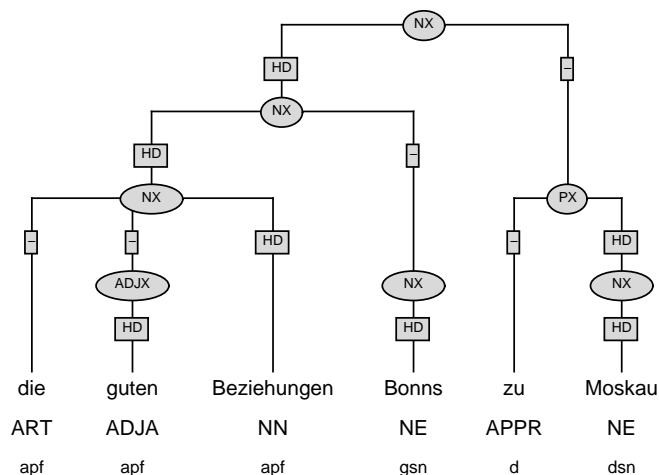
4.2.3 Postnominal Modification

Whereas prenominal modifiers are always directly attached to the head noun on the same level, postnominal modifiers are attached to the head noun on a higher level. Postnominal modifiers are also always first projected to the phrase level before they are attached to the head noun on a higher level. Phrase internal postmodifiers can be of any phrasal category. The following tree structures show a prepositional phrase (PX) and a genitive phrase (NX) as postmodifiers. See section 6.4, page 109 for the analysis of relative clauses.

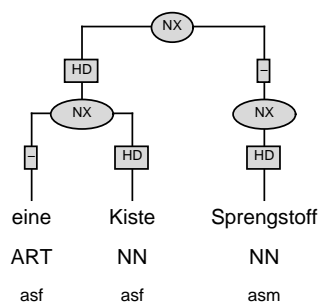


In case a noun has more than one postmodifier, these modifiers usually show a hierarchical structure, for example, the first modifier modifies the head noun, the second

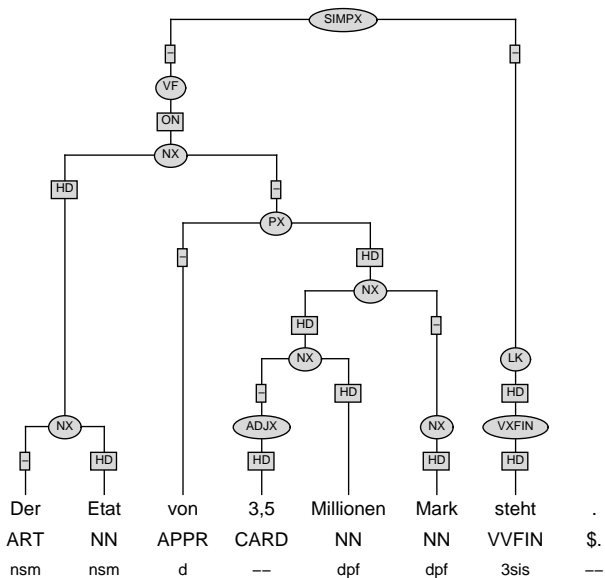
modifier modifies the complete preceding noun phrase structure, and so on.

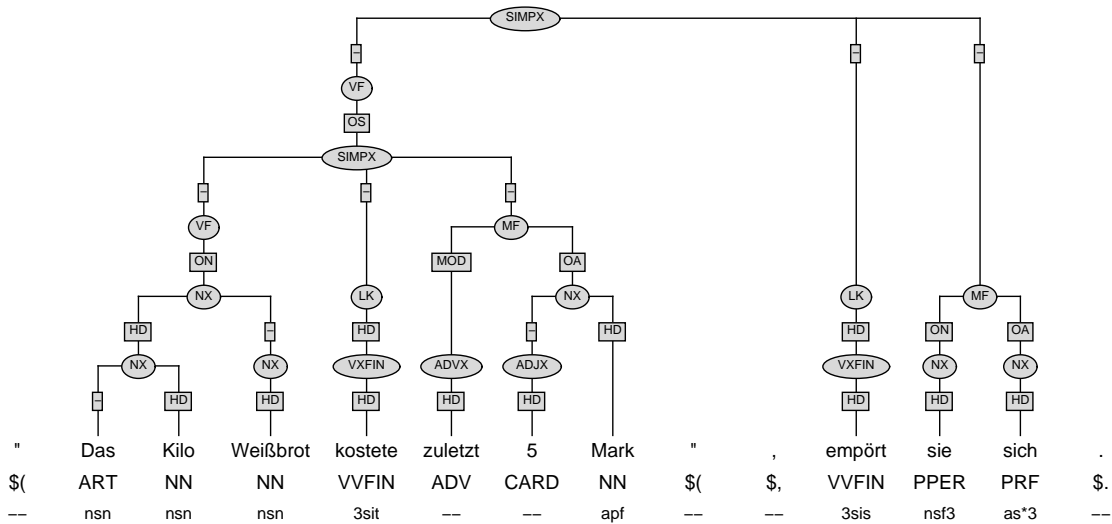


Attributes of degree and quantity nouns are also defined as postnominal modifiers:

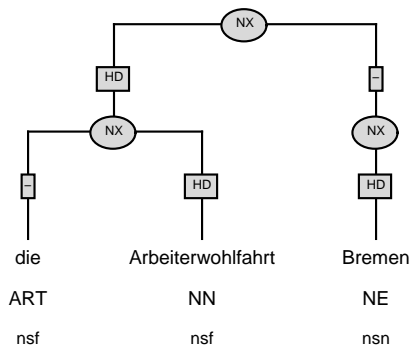
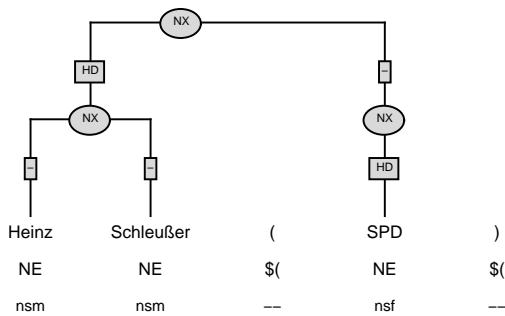


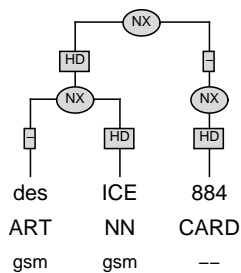
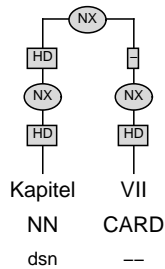
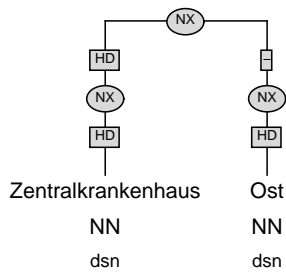
Cardinal numbers either appear as quantity nouns or premodifying adjectival attributes, e.g. the cardinal number *1,000,000* can also be expressed by the quantity noun *eine Million*. Therefore, we have to distinguish the following two ways of annotation:





For nominal postmodifiers apart from genitive phrases the same attachment rule is applied. This kind of postmodifiers which may also appear in brackets, e.g. *Heinz Schleußer (SPD)*, is semantically closely related to the preceding head noun phrase. *die Arbeiterwohlfahrt Bremen*, for instance, means *die Arbeiterwohlfahrt* which is located in *Bremen*, but does not mean *die Arbeiterwohlfahrt* which is called *Bremen*. Hence, these postmodifiers have to be distinguished from appositions (cf. 4.2.4) and complex named entities (cf. 4.2.6).



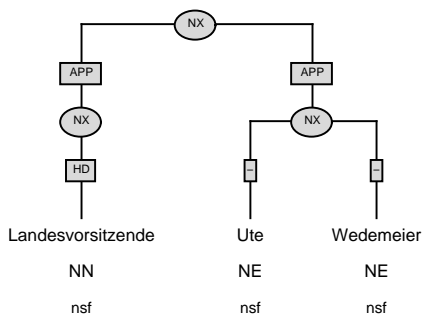
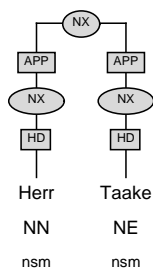
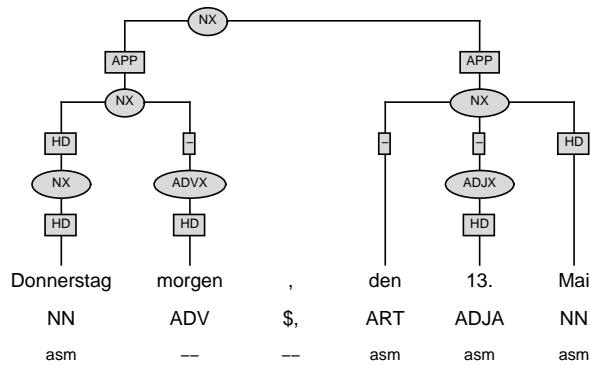


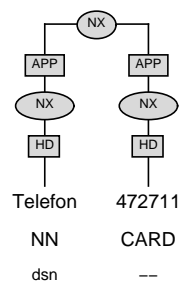
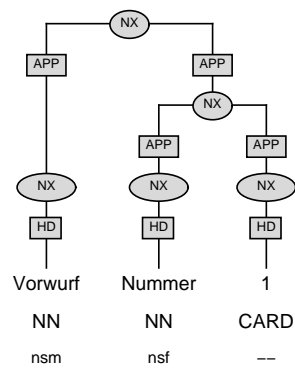
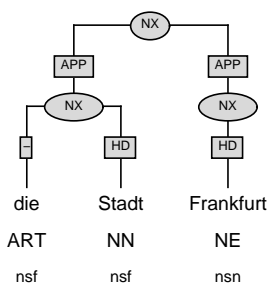
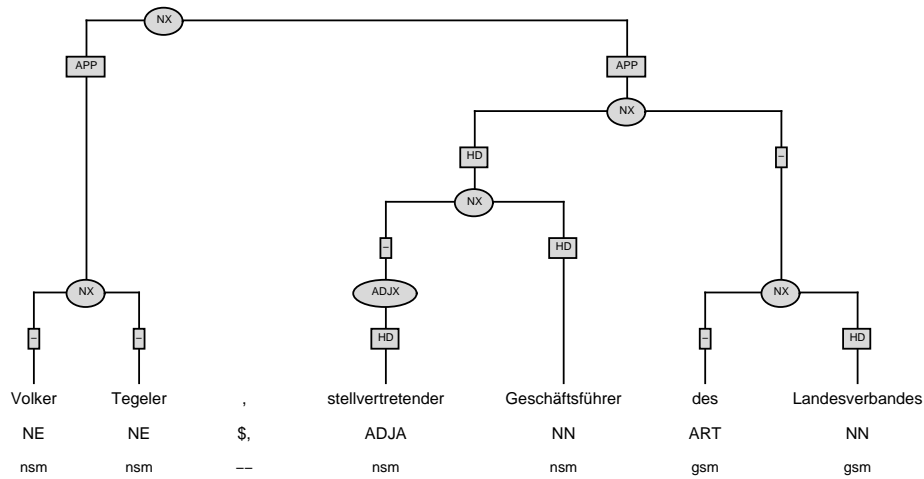
4.2.4 Appositional Constructions

An apposition is a specific kind of attribute to a noun, which normally agrees in case with this noun and does not change its overall meaning. There is no consensus among grammarians of what is exactly meant by the notion *apposition* (cf. (Eisenberg 1999 2001)). Eisenberg (1999 2001), for instance, claims that, e.g. *Ute Wedemeier die Landesvorsitzende* and *die Landesvorsitzende Ute Wedemeier* are both appositions but it is not clear which part is the apposition and which part is the head noun. The Duden Grammar (1995) distinguishes between loosely constructed appositions (*lockere Apposition*) (e.g. *Ute Wedemeier, die Landesvorsitzende*), which follow the head noun separated by a comma, and tightly constructed appositions (*enge Apposition*) (e.g. *(die) Landesvorsitzende Ute Wedemeier*), which precede the head noun (cf. (Drosdowski 1995)). According to Helbig/Buscha (1998) there is case agreement between loosely constructed appositions and head nouns which are separated by a punctuation mark. By contrast, Engel (1996) thinks that only loosely constructed appositions can be regarded as appositions. He treats tightly constructed appositions as *nomen varians* or *nomen invariants*.

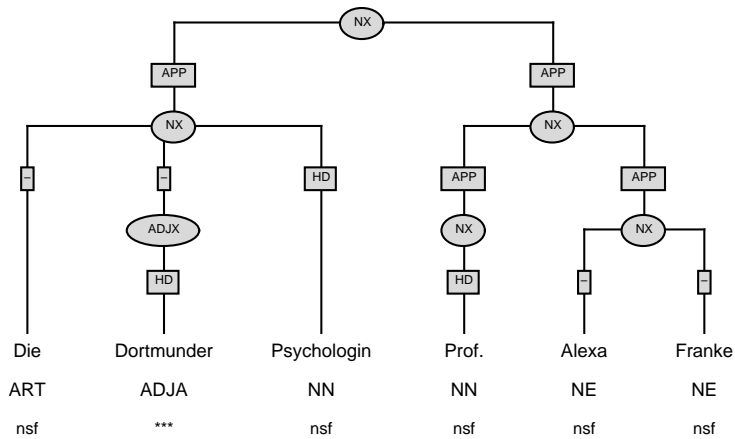
Because of these different definitions of the notion of apposition, we do not decide on what is the head noun and what is the apposition. We assume referential identity between the two parts. Loosely constructed appositions as well as tightly constructed appositions are treated as appositional constructions, i.e., the head noun and its apposition

form a complex structure which does not give any information about head assignment. Therefore, both parts are first projected to their phrase level and then coordinated on a higher level, each of them labelled as apposition (APP), i.e. as a part of an appositional structure. What is important is the referential identity in meaning. Thus, *Nummer 1* is an appositional construction, whereas *Seite 1* is a noun phrase with the postmodifier *1*. Forms of address for persons and titles, e.g. *Herr*, *Frau*, *Doktor (Dr.)*, *Professor (Prof.)*, *Bundeskanzler*, are also treated as appositional constructions. Here are some examples:

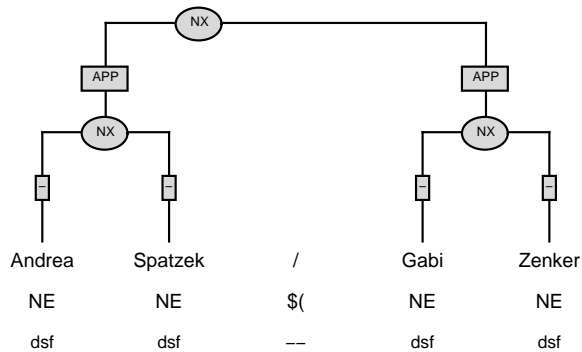




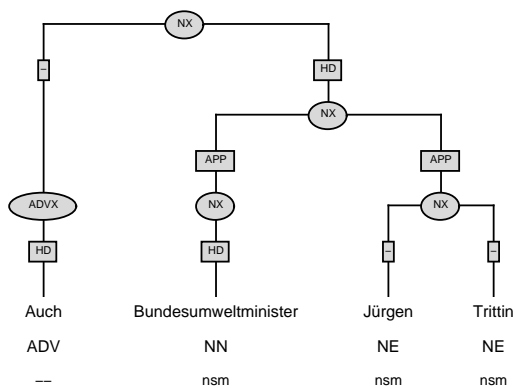
In case of a form of address combined with one or more titles preceding a name, we annotate an embedded appositional construction:



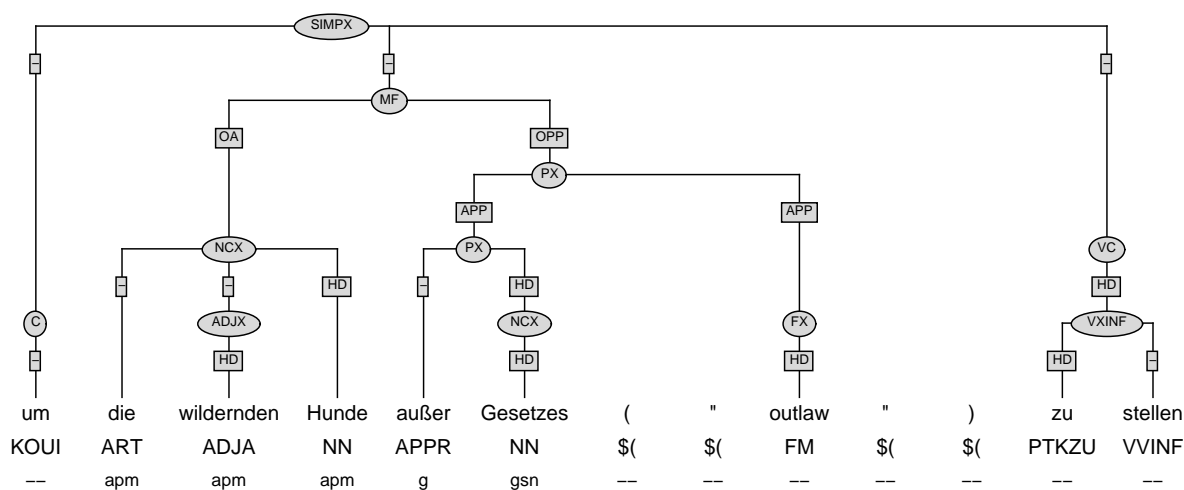
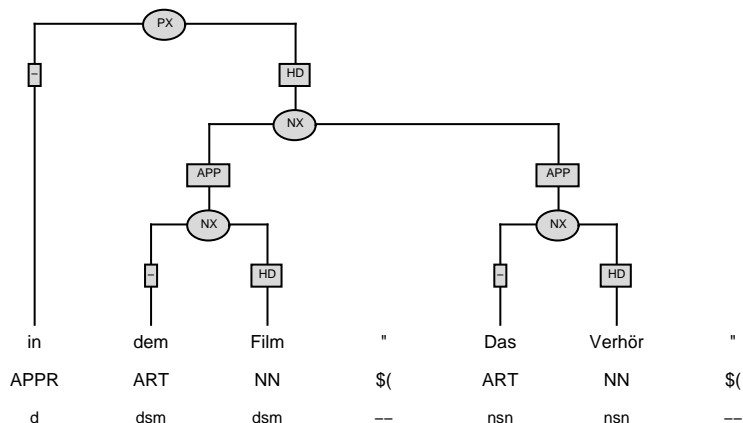
The same way, we treat proper nouns which are identical to a preceding proper noun, for example, an actor's name and role:



Premodification of the whole appositional construction is attached to an additional NX level.

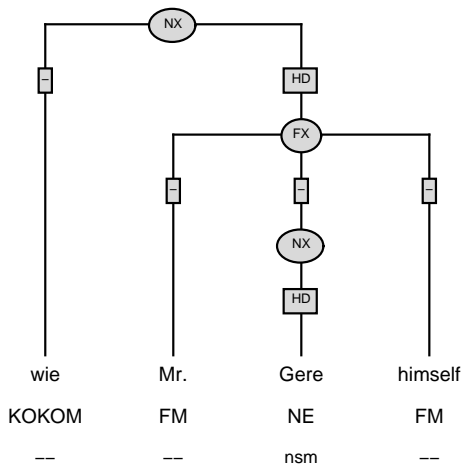
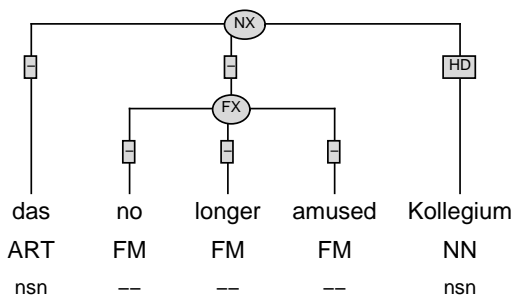
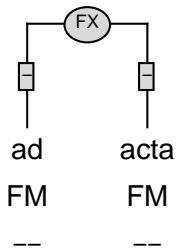
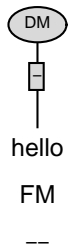


There are some examples in which the appositional construction does not agree in case. These are postnominal titles of books, movies, etc. and translations interspersed in the sentence. In the latter type, we extend the appositional construction also to non-nominal phrases.

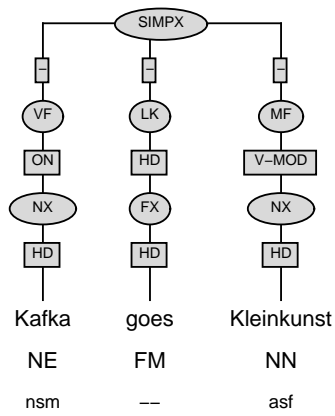


4.2.5 Foreign Language Material

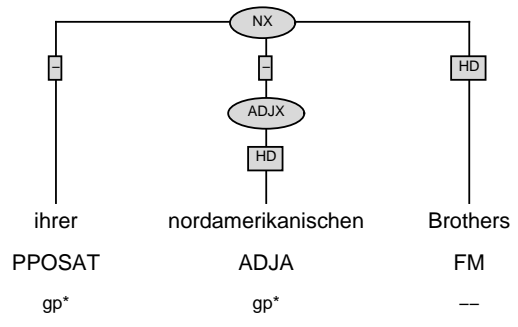
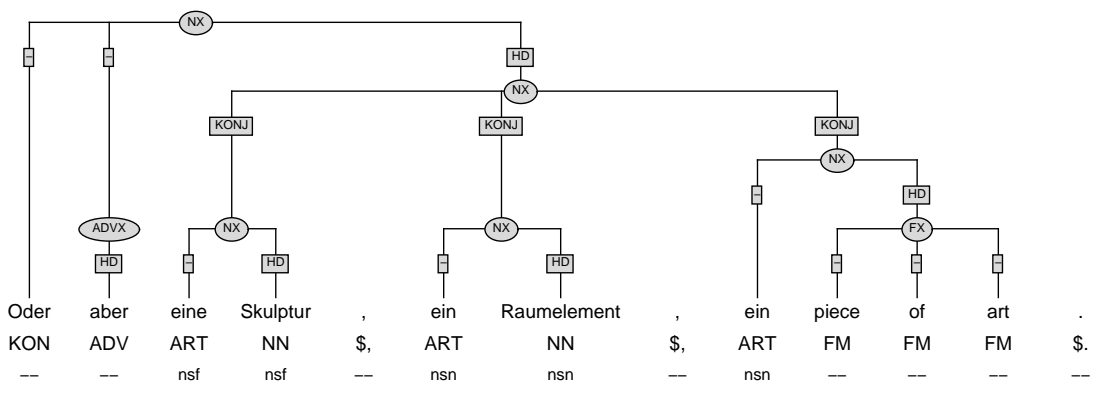
Words or parts of a text written in a foreign language except foreign language proper nouns are tagged as foreign language material (FM), e.g. *hello* (FM), *no* (FM) *longer* (FM) *amused* (FM). All parts of foreign language proper nouns are tagged as NE (e.g. *Mary* (NE) , *New* (NE) *York* (NE), *University* (NE) *of* (NE) *Illinois* (NE)). Single foreign words are projected to a syntactic level assigned the node label FX, which is an universal label for any syntactic category (phrasal and sentential) in the respective foreign language. More complex parts of a text tagged as FM are attached on the same level without any internal syntactic structure and head assignment. Their mother node is also assigned the label FM, e.g. *no longer amused*. For foreign language constructions containing a proper noun, the annotation strategy is the following: in a first step, all NEs are projected to the phrase level (NX), in a second step, these phrase node labels together with all FMs are projected to the next higher level with the node label FX. Again, there is no head assignment directly below the FX node, e.g. *Mister Gere himself*.



Often, foreign language material is a part of a German syntactic construction and plays the role of a grammatical function. Therefore, the FX node is attached as a constituent to the tree structure. If it is directly attached to a field or a sentence bracket, the edge label above the FX node denotes its grammatical function within the clause, e.g. *Kafka goes Kleinkunst* (head of the clause).



If a FX or a single FM is head of a phrase which can be identified as a German phrase, e.g. by an article and/or an adjective (noun phrase), it is projected to the specific phrasal category, e.g. NX instead of FX in constructions like *ein piece of art* or *ihrer nordamerikanischen Brothers*.



If FX is modified by a postmodifier the mother node of the complex phrase is also FX, which again may be preceded by another phrase, e.g. *Unter der Überschrift 'user als loser'*.

Table 4.1 lists the commonly occurring semantic subclasses for named entities in the TüBa-D/Z with examples:

Table 4.1: Semantic Classes and Subclasses for Named Entities

Semantic Classes	Common Semantic Subclasses	Examples
PER	persons surnames names of animals (personified)	<i>Hans Winkler</i> <i>(Familie) Feuerstein</i> <i>(Schweinchen) Babe</i>
ORG	organisations companies institutes museums newspapers, journals clubs theaters, cinemas universities TV and radio stations restaurants, hotels forces fashion labels sporting events bands	<i>Nato, EU</i> <i>Microsoft, Bertelsmann,</i> <i>Institut für chinesische Medizin</i> <i>Pergamonmuseum</i> <i>Süddeutsche Zeitung, Der Spiegel</i> <i>VfB Stuttgart</i> <i>Metropol-Theater, CinemaxX</i> <i>Freie Universität</i> <i>Arte, Radio Bremen</i> <i>Sassella, Adlon</i> <i>Blauhelme</i> <i>Chanel</i> <i>Olympische Spiele, Wimbledon</i> <i>Beatles, Die Fantastischen Vier</i>
LOC	districts sights, churches planets geographical areas streets, places mountains, lakes continents	<i>Schöneberg</i> <i>Brandenburger Tor, Johanniskirche</i> <i>Mars</i> <i>Königsheide</i> <i>Sögestraße, Alexanderplatz</i> <i>Alpen, Viktoriasee</i> <i>Europa, Asien</i>
GPE	countries, states (incl. historical) cities (incl. historical)	<i>Frankreich, Hessen, Assyrien</i> <i>Berlin, Babylon</i>
OTH	operating systems titles of books, movies, etc. mottoes, slogans wars	<i>DOS</i> <i>Faust, Schlaflos in Seattle</i> <i>Zwischen Himmel und Erde</i> <i>Zweiter Weltkrieg</i>

In order to annotate the semantic classes, **syntactic-semantic node labels** of the pattern *syntactic category = semantic class* are defined for the mother node of named entities (see Table 3.9). The syntactic-semantic nodes indicate that the structure below represents a (complex) named entity of a certain syntactic category belonging to one of the five semantic classes (cf. 3.4.2).

The former node label 'EN-ADD' and the secondary edge label 'EN' are deleted.

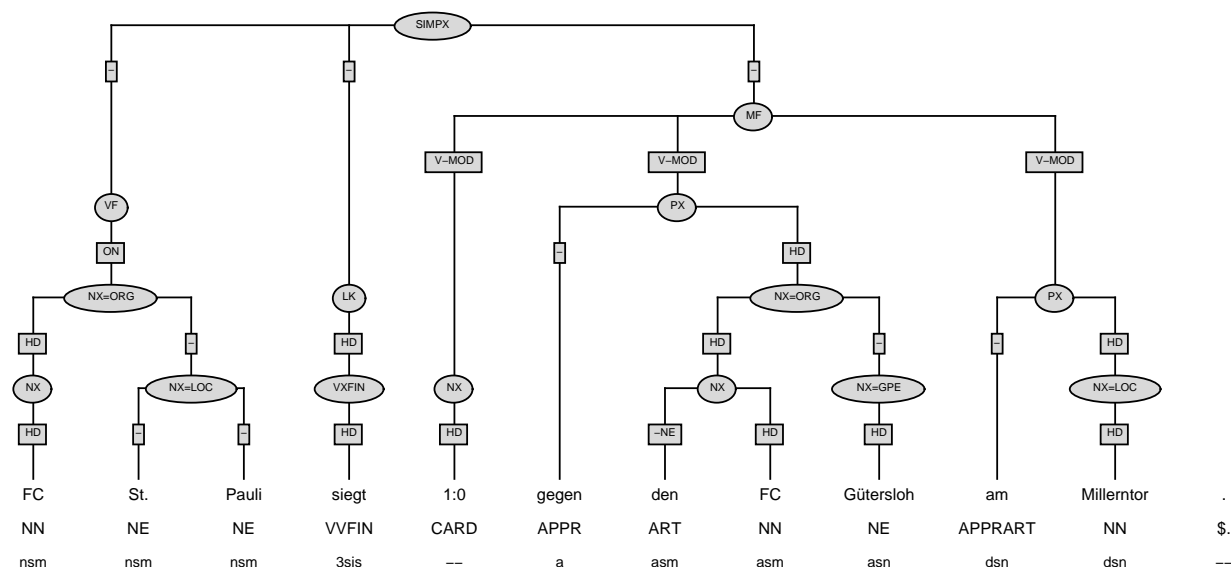
Our annotation strategy for named entities is shown in the following tree examples:

Named entities may consist of one or more lexical elements tagged as NE. In case of a single NE, this NE is projected to its phrase level, carrying the respective syntactic-semantic node label (*Gütersloh* (NX=GPE)). Named entities consisting of two or more NEs are attached on the same level. None of them carries a head label in order to indicate that there is no obvious dependency relation between them (*St. Pauli* (NX=LOC)).

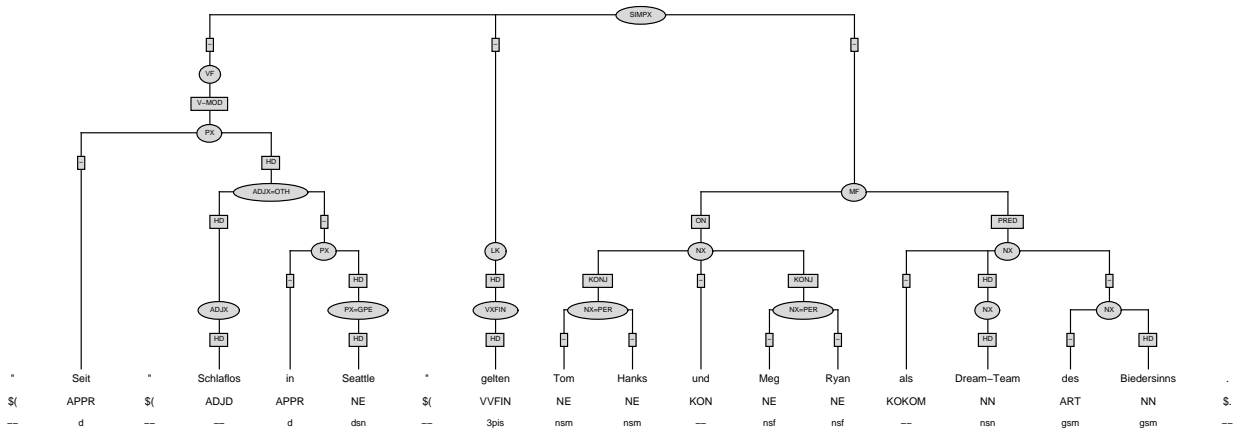
There can occur postmodifiers within a named entity which are a named entity themselves like *St. Pauli* (NX=LOC) in *[[FC (NX)] [St. Pauli (NX=LOC)]]* (NX=ORG).

Parts which do not belong to a named entity are marked with the edge label ‘-NE’ as in *[[den (-NE) FC (NX)] [Güthersloh (NX=GPE)]]* (NX=ORG).

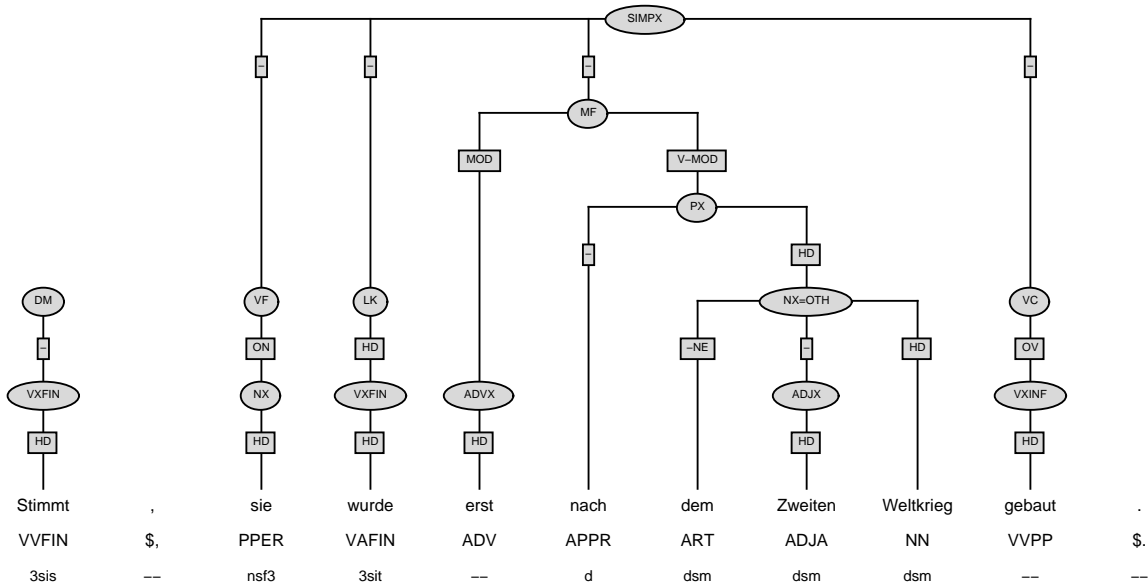
Named entities which are not tagged as NE, e.g. *Millerntor* (NX=LOC) are also assigned a syntactic-semantic node label.



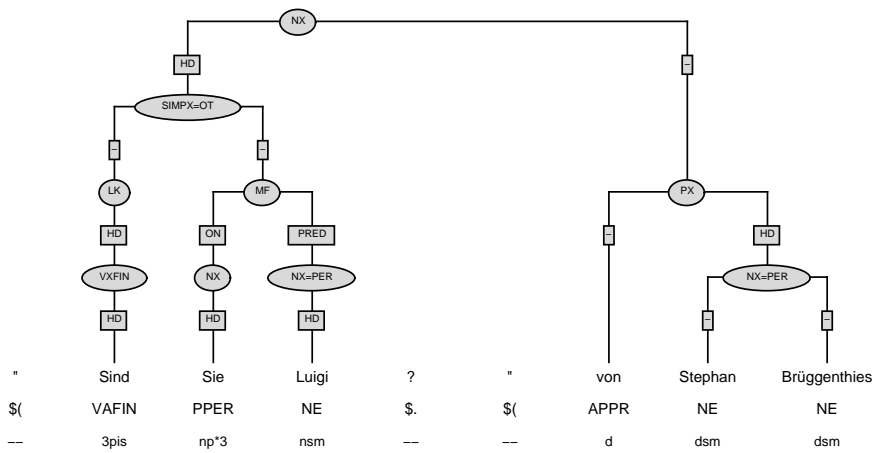
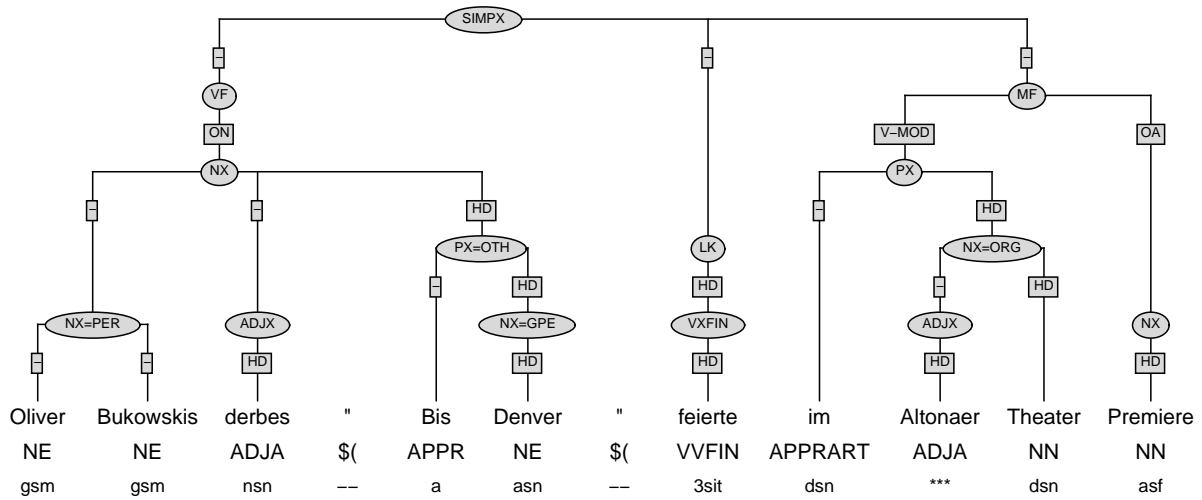
As mentioned above, all elements of German named entities which consist of a complex syntactic structure, e.g. a phrase or a sentence, are always tagged according to their distribution and annotated with their internal syntactic structure as noun phrases, prepositional phrases, adjectival phrases, clauses, etc., e.g. the movie title (*Schlaflos* (ADJD) *in* (APPR) *Seattle* (NE)) (ADJX=OTH) in the following tree example. If two named entity nodes are coordinated like *Tom Hanks* (NX=PER) and *Meg Ryan* (NX=PER) their mother node is NX which represents the nominal status of the named entity.



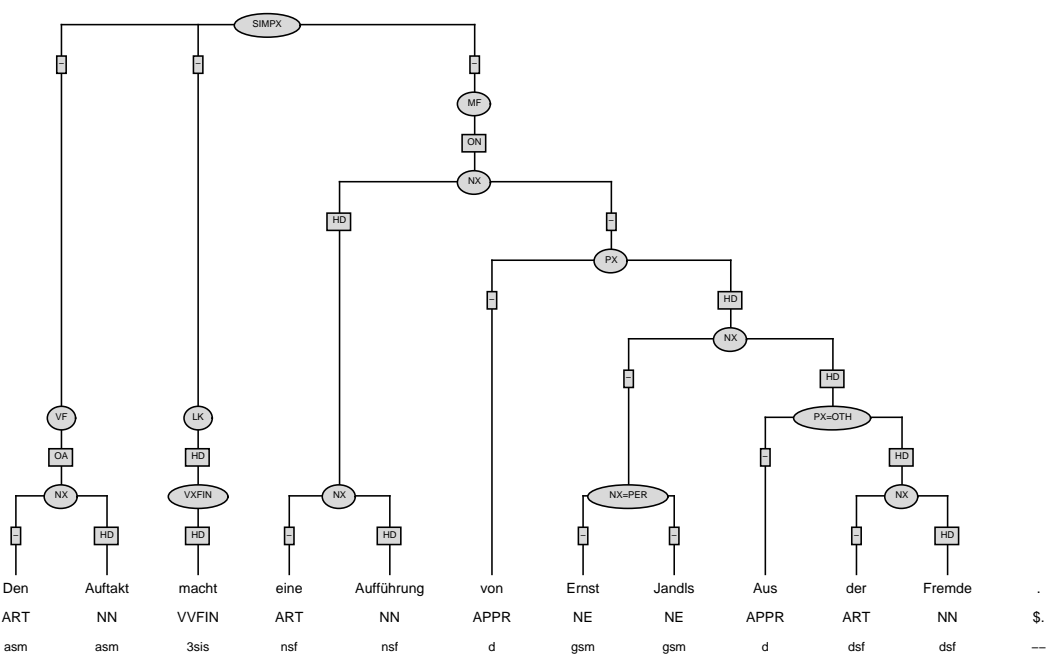
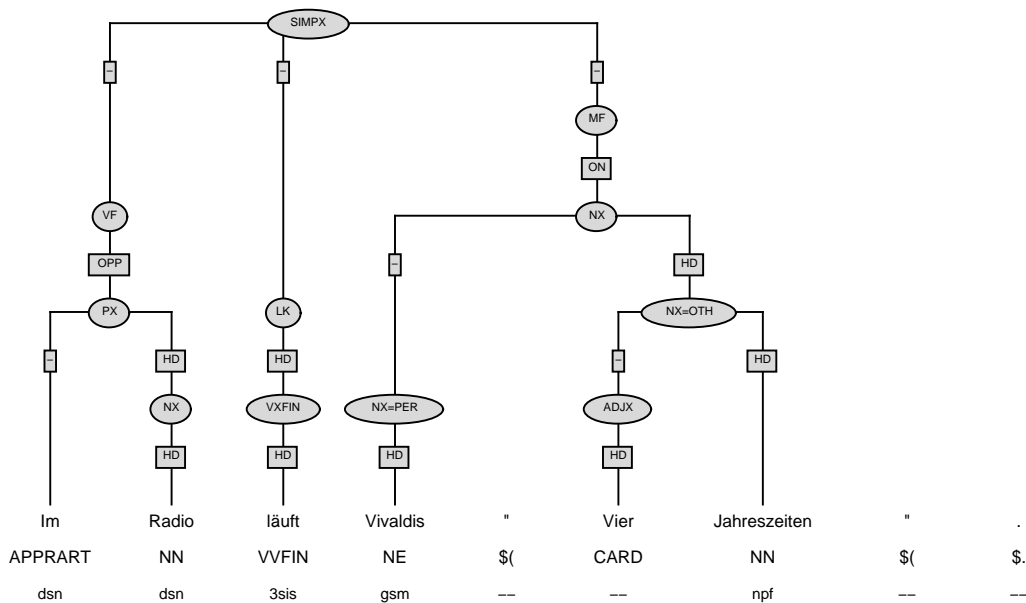
If the original form of a named entity (e.g. *Zweiter Weltkrieg*) is inflected and/or premodified by an article and/or attributive adjective like in the following example tree (*dem Zweiten Weltkrieg* (NX=OTH)), the mother node of the named entity carries the semantic class information and all parts that do not belong to the named entity are assigned the edge label ‘-NE’.



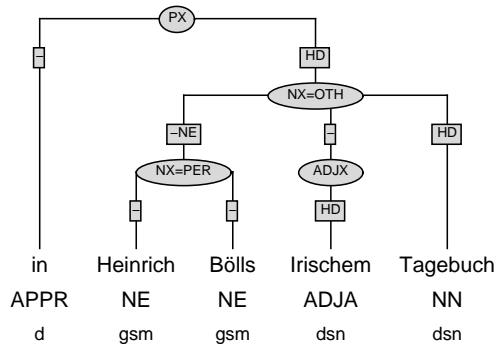
If a named entity of a syntactic category other than NX has a premodifier or a post-modifier, (both can be a named entity itself) the mother node of the whole constituent is always NX which represents the nominal status of the named entity:



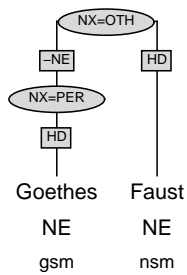
If a named entity is a syntagma with its own internal syntactic structure, i.e. it does not agree with the inflection of another constituent of the sentence (e.g. uninflected titles of books, movies, etc.) all premodifiers are attached on a higher level:



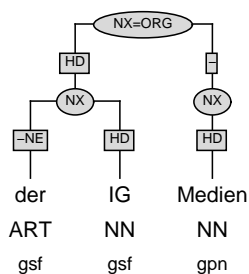
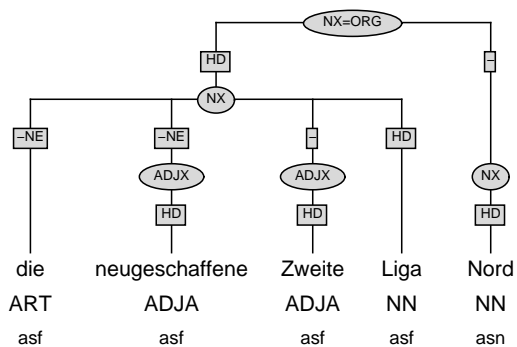
If the named entity is inflected, all premodifiers are directly attached to the head noun, labeled as '-NE':



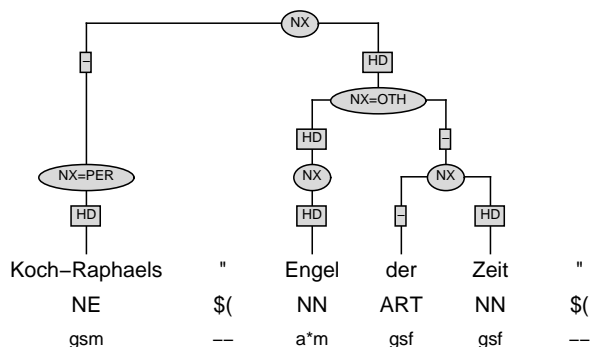
Premodifiers of single word named entities are also directly attached to the head noun:



If a postmodifier is part of a named entity which is not a title, all premodifiers are directly attached to the NX of the head noun. If the premodifiers are not part of the named entity itself, they are assigned the edge label '-NE'.

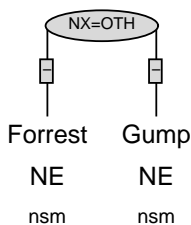


In the case of titles including a postmodifier, all premodifiers of the named entity are attached on a higher level:

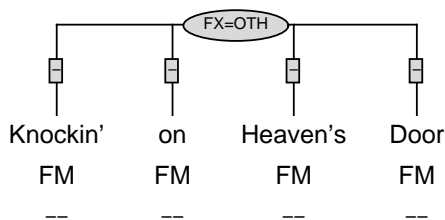


Foreign Language Named Entities

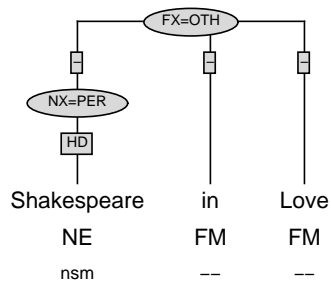
The syntactic annotation of foreign language named entities differs from the annotation of German named entities in the following aspects. According to the STTS guidelines, foreign language proper nouns are tagged as NE, while all other lexical elements of a foreign language are tagged as foreign language material (FM). A foreign language named entity which consists of a proper noun, e.g. the title of a movie is assigned a syntactic-semantic node label of the category NX (*Forrest Gump* (NX=OTH)).



If a foreign language named entity consists of only FM tagged tokens, these tokens are directly attached on the same level without internal syntactic structure. The mother node of the phrase is marked as a syntactic-semantic node label of the category FX, e.g. *Knockin' on Heaven's Door* (FX=OTH).

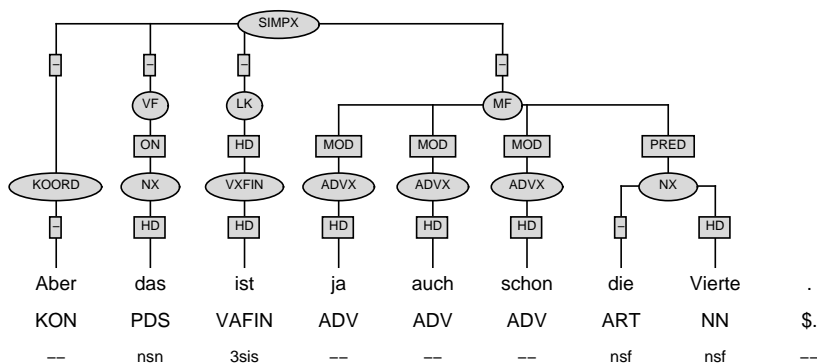
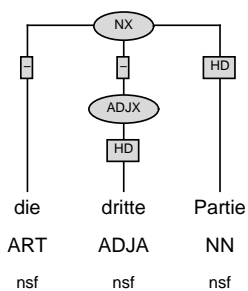


If a foreign language named entity consists of NE as well as FM tagged tokens, e.g. *Shakespeare* (NE) *in* (FM) *Love* (FM), NE is projected to NX=PER. The NX=PER node and all FM tagged tokens are attached directly on the same level.



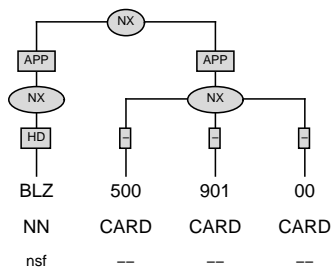
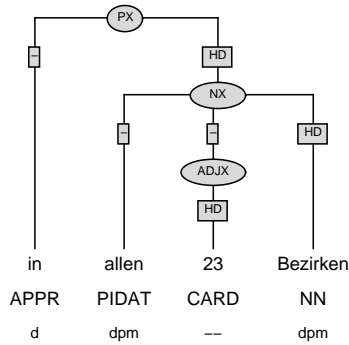
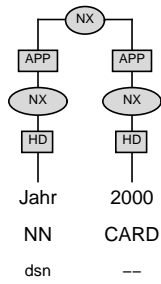
4.2.7 Ordinal Numbers

According to their distribution, ordinal numbers occur either as a premodifying attributive adjective (e.g. *die dritte (ADJA) Partie*) or as a head noun (e.g. *er ist der sechste (NN)*). In the first case, the premodifier is projected to an adjectival phrase, in the latter case it is projected to a noun phrase.

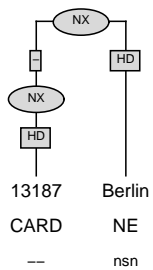


4.2.8 Cardinal Numbers

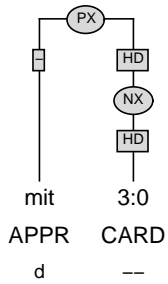
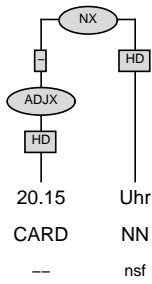
According to their syntactic function (nominal or adjectival), cardinal numbers (CARD), are either projected to NX or ADJX. If their numerals are written separately or in groups, e.g. numbers of bank accounts, they are attached on the same level like proper names without internal head assignment.



A premodifying cardinal number is nominal if it does not express a quantity like in the example above, but a characteristic of the following noun, e.g. the number of a zip code:

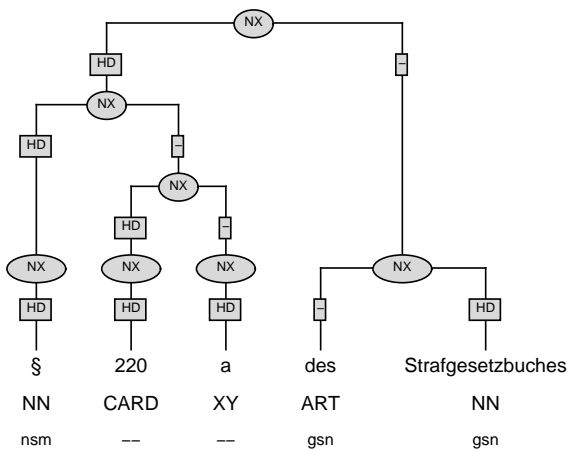
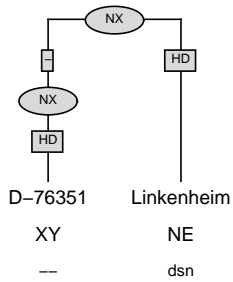


Complex time expressions or results of competitions are also treated as cardinal numbers:



4.2.9 Letters and Non-Words

Letters and non-words are tagged as XY. They are projected to their phrase level and assigned the syntactic category to which they belong in the construction. Signs which represent a lexical element, e.g. the sign for paragraph, are tagged with the respective part-of-speech tag:



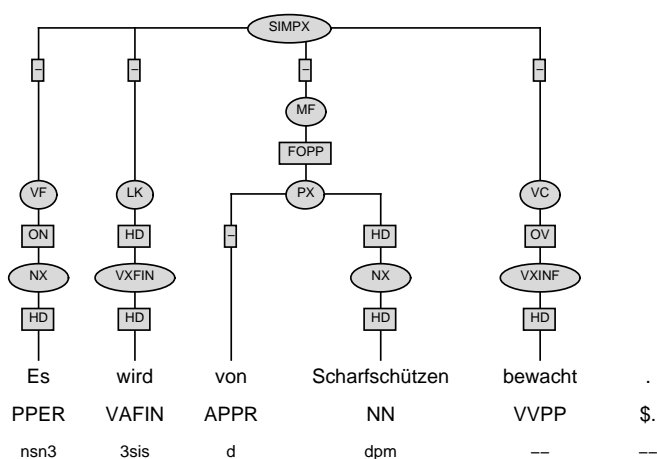
4.2.10 Expletive and Other Uses of *es*

The pronominal form *es* functions as expletive element in German. Three different expletive usages are traditionally distinguished: formal subject or object, correlate of an extraposed clausal argument, and *Vorfeld-es* (cf. (Eisenberg 1999 2001), (Pütz 1986)). For sake of completeness, the following list begins with an example of *es* as a referential personal pronoun.

Personal Pronoun

The pronoun functions as an argument of the verb and refers to some person, object, or event that is salient in the context. It can be tested, whether *es* is used as a pronoun by replacing it by another noun or pronoun (such as *das* or *er/ihn*).

In the example tree *es* refers to the neuter noun *Gästehaus* in the preceding sentence: *Die italienische Regierung hat die Familie im staatlichen Gästehaus Casino dell'Algardì untergebracht.*

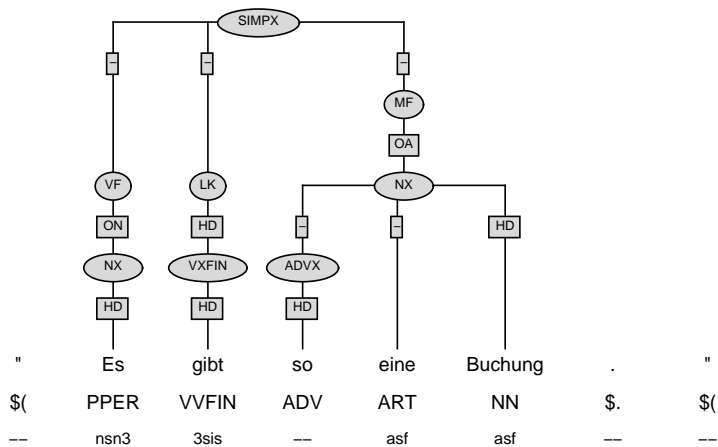


Formal Subject or Object

The formal subject obligatorily occurs with *weather verbs*, e.g. *Es regnet* and impersonal or agentless constructions such as *Es gibt so eine Buchung* or *Es geht um populäre Unterhaltung*. Some verbs optionally permit an expletive subject but also occur with referential subjects such as *Max/Es klopft an der Tür*. A formal object is found in constructions like *jmd. legt es an auf etw.* or *jmd. verdirbt es mit jmdm.* In all examples mentioned, *es* functions as a grammatical argument without semantic contribution, i.e. it does not refer to a person, object, or event.

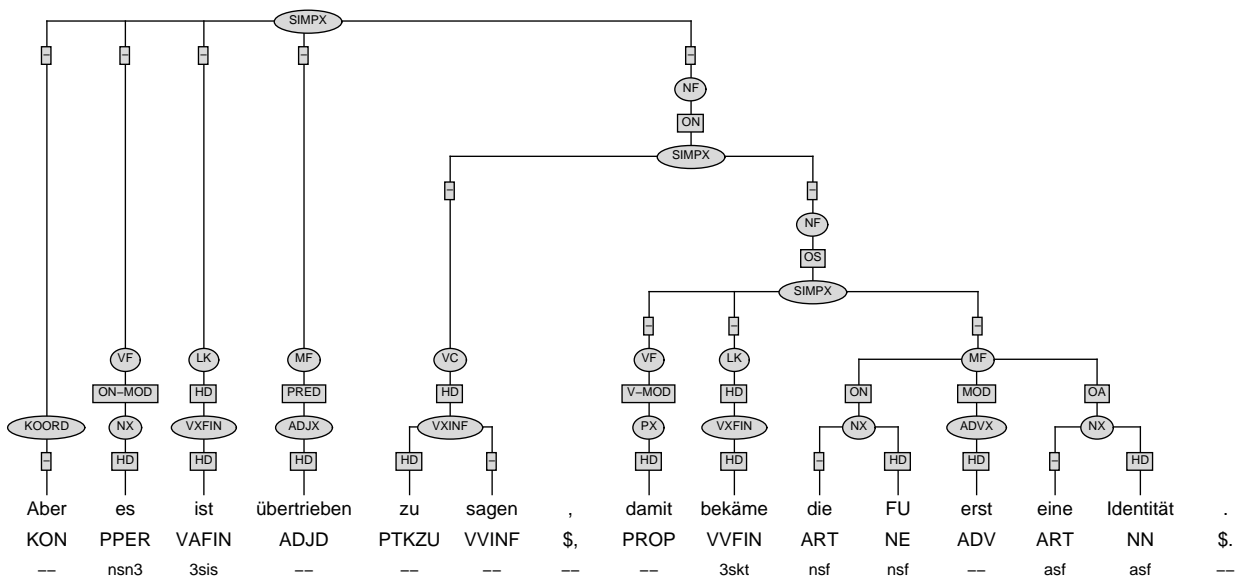
In TüBa-D/Z formal subjects and objects are treated like referential pronouns and are labelled alike, e.g. with edge labels ON or OA.

Formal arguments are obligatory and may occur in the Mittelfeld. In case of doubt, it is a good test to paraphrase the sentence such that another element occupies the Vorfeld, e.g. *Natürlich gibt es so eine Buchung* versus **Natürlich gibt so eine Buchung*.



Correlate of an Extraposed Clausal Argument

If a clausal argument is extraposed in the *Nachfeld*, it is optionally doubled by an expletive *es* in the *Vorfeld* or *Mittelfeld*. The expletive is labelled ON-MOD or OS-MOD depending on the function of the clausal argument.



Vorfeld-es

The last type is a purely structural dummy element. It occurs in *Vorfeld* position only and is not correlated with any argument of the clause. It does not agree with the verb which becomes evident if there is a plural subject in the *Mittelfeld*, which is illustrated in the example tree below. It is ungrammatical in the *Mittelfeld*, e.g. **... dass es ihn die Völker zahlen*. *Vorfeld-es* is labelled ES to indicate its purely structural function. In the first release of TüBa-D/Z, 12/2003, *Vorfeld-es* was integrated by means of ON-MOD.

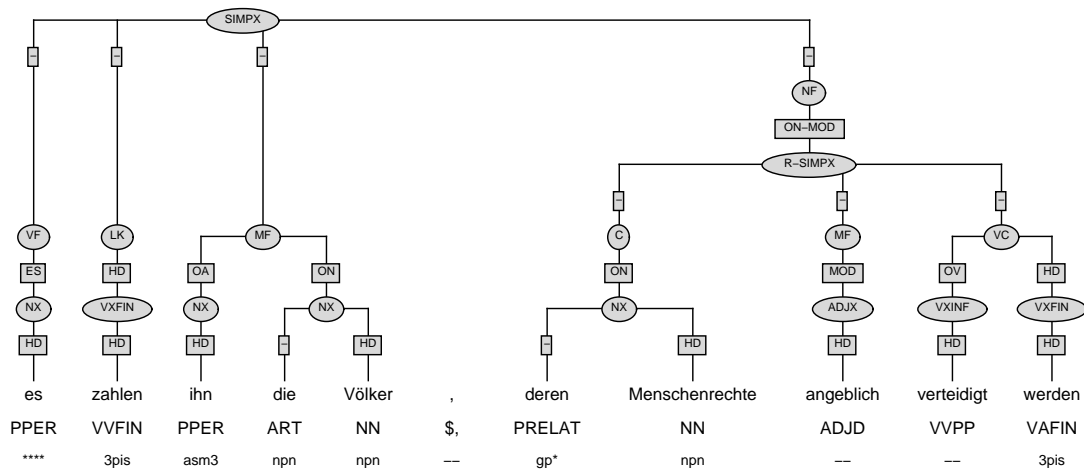


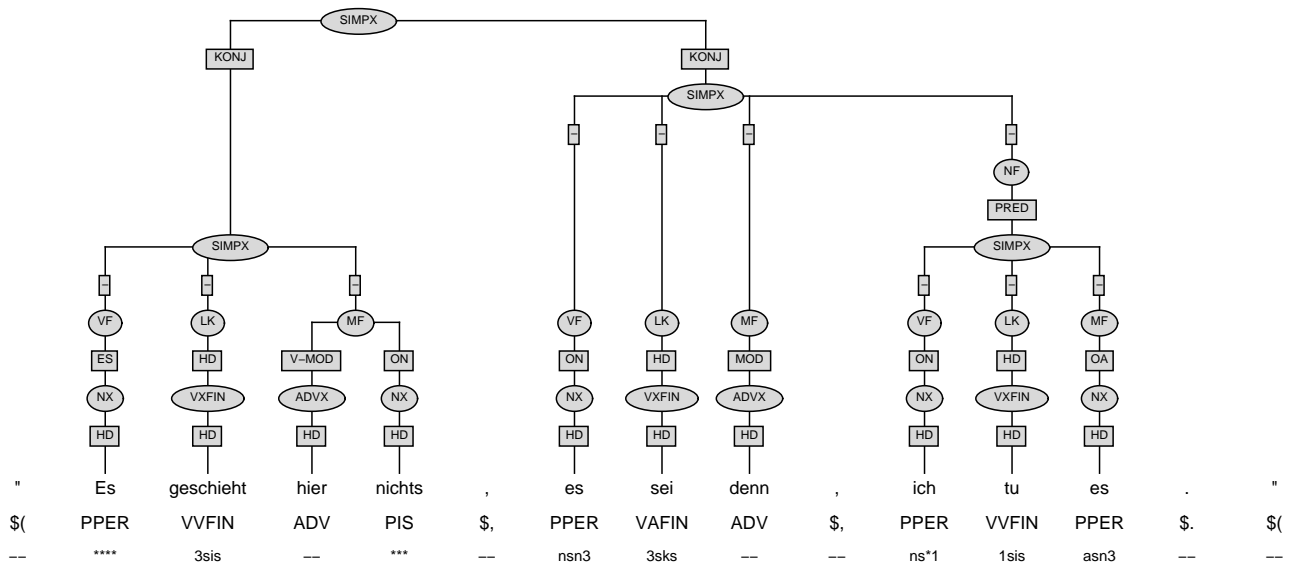
Table 4.2 summarizes tests and labels for the different uses of *es*.

Table 4.2: Types of *es*

test	type	referential pronoun	formal argument	correlate	Vorfeld- <i>es</i>
substitutable by other pronouns		yes	no	no	no
optional		no	no	yes	no
correlates with clausal argument		no	no	yes	no
ungrammatical in Mittelfeld		no	no	no	yes
edge label		ON, OA, etc.	ON, OA	ON-MOD, OS-MOD	ES

Es sei denn

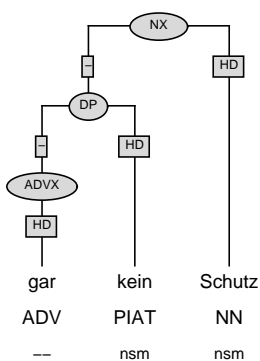
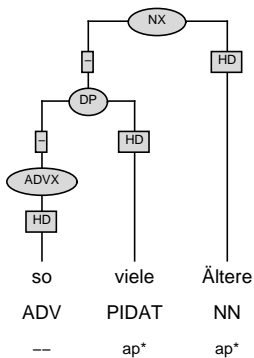
The lexicalized phrase *es sei denn*, meaning *außer*, is analysed as a copula construction.



4.3 Determiner Phrases

Certain pronouns serving as determiners in noun phrases may be premodified, for instance, by degree adverbs such as in *so viele Ältere, gar kein Schutz*, etc.

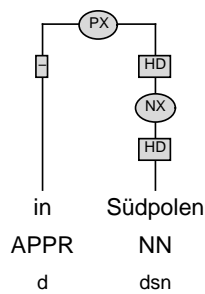
In the case of *so viele Ältere*, the premodifying adverb *so* is attached to the indefinite pronoun *viele*. Together, they form a determiner phrase (DP), which is attached to the head noun *Ältere* on the same level:



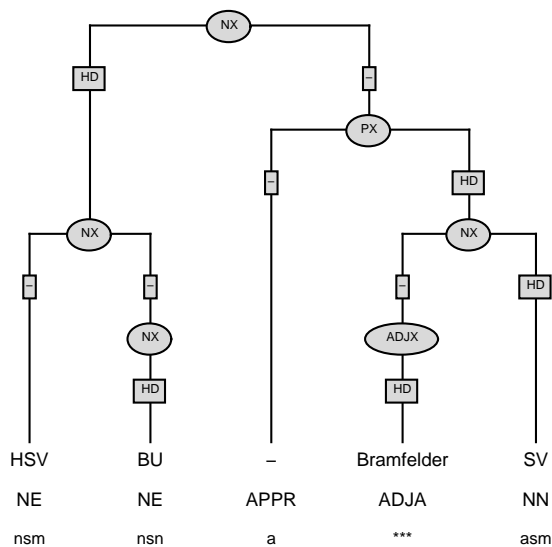
4.4 Prepositional Phrases

4.4.1 Prepositions

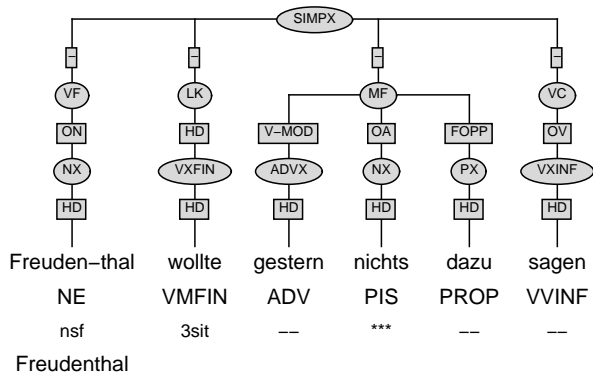
Considering prepositional phrases, it turns out to be appropriate not to annotate the preposition as the head of the phrase. It is rather reasonable to annotate the complement within the prepositional phrase as the head. This decision facilitates the identification of dependencies between verbs and their nominal complements and adjuncts. Moreover, it is in accordance with basic assumptions in Dependency Grammar .



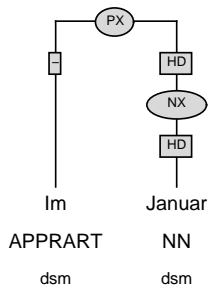
If the preposition is realized as a non-alphabetic sign, e.g. - (*bis*, *gegen*), this sign is tagged as APPR and annotated like a preposition:



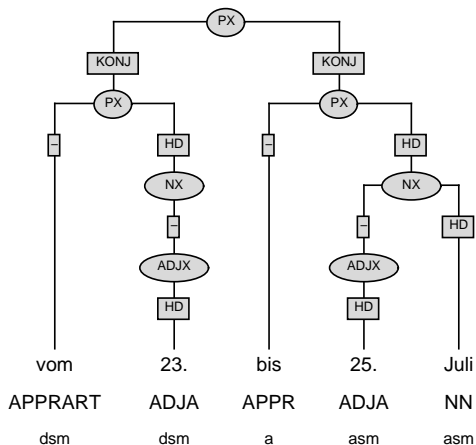
Since pronominal adverbs (PROP) are pronominal forms of a prepositional phrase, they are directly projected to PX:

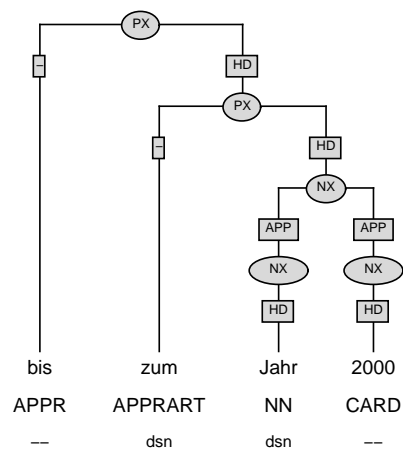
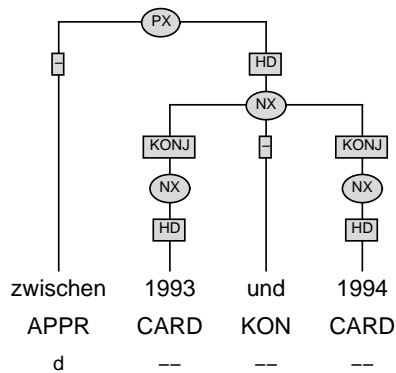


In German, there are so-called *Verschmelzungsformen*, i.e. merged forms of a preposition and a determiner, e.g. *in dem Januar* amalgamates to *im Januar*. The merged form is assigned the part-of-speech tag APPRART (including richer morphological annotation). In terms of syntax, it is annotated like a preposition:



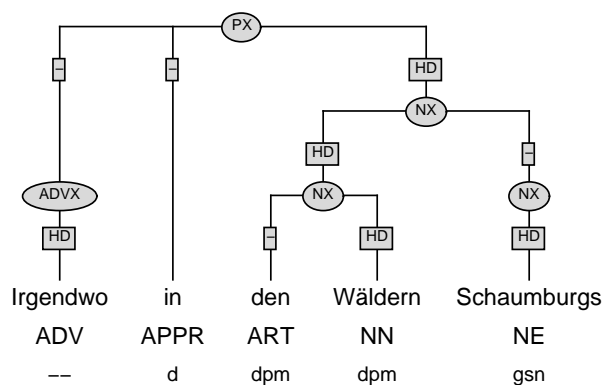
Prepositional phrases expressing intervals, e.g. with *von/bis*, *von/bis zu* or *zwischen*, are annotated in the same way as coordinate structures (cf. 6.5.1), i.e. without head assignment on the level of coordination, since the two phrases are assumed to be conjuncts. If two prepositions follow each other (e.g. *bis zum*), the result is an embedded structure of a prepositional phrase taking another preposition. The first preposition does hereby not receive a morphological case feature.



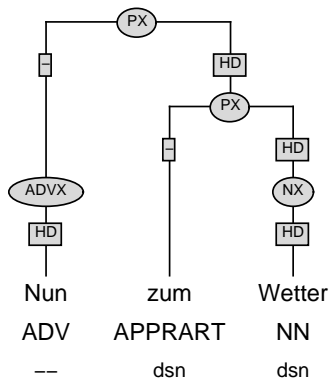


As opposed to the case with two prepositions, intervals like *dritter bis fünfter November* are annotated as a coordinate attributive adjective phrase within a simple noun phrase (cf. 6.5.1).

Premodification of non-isolated prepositional phrases follows the general principle of low attachment.

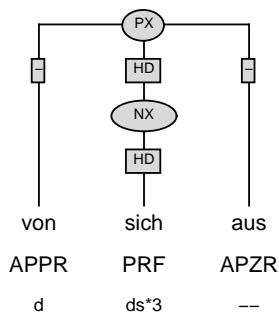


There is one exception to the low attachment principle: isolated phrases in which a preceding adverb does not semantically modify the prepositional phrase. In this case the adverbial phrase is high attached to an additional level of PX.

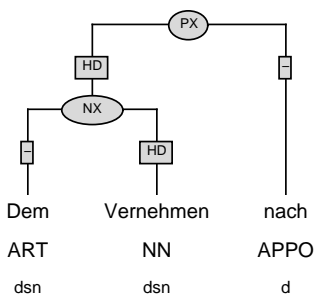


4.4.2 Circumpositions and Postpositions

Circumpositions are treated as ternary branching prepositional phrases. The circumposition on the left hand side is tagged as APPR and the circumposition on the right hand side as APZR:



Postpositions are tagged as APPO. The complement of the postposition occurs on the left side and constitutes the head of the prepositional phrase:



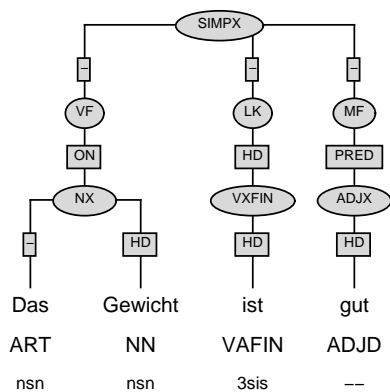
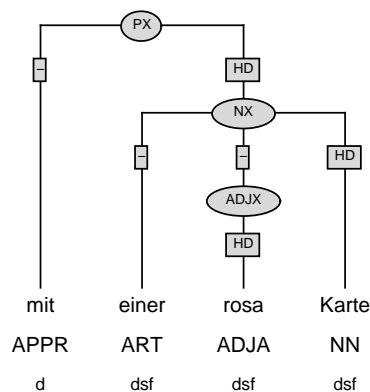
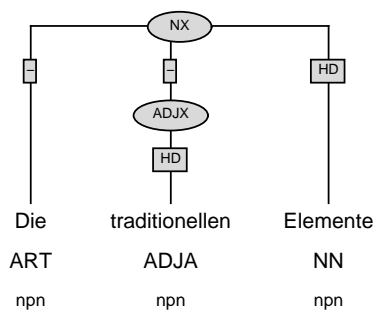
4.5 Adjectival Phrases

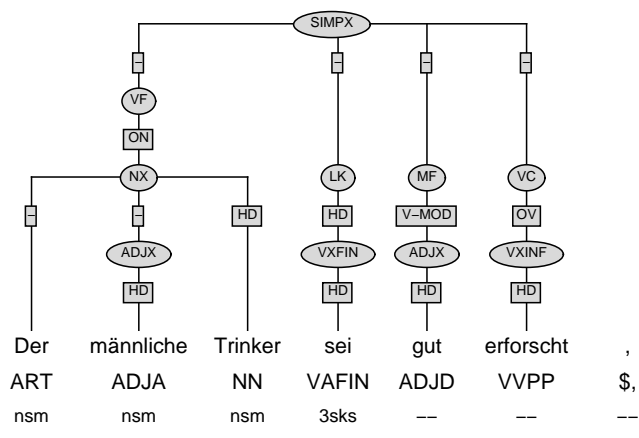
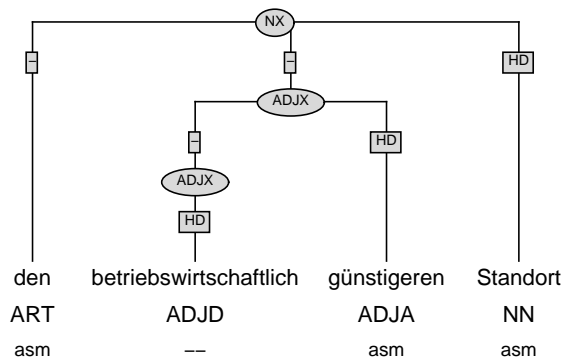
We distinguish between attributive adjectives on the one hand and adverbial or predicative adjectives respectively on the other hand. Attributive adjectives are tagged as

ADJA (*die traditionellen Elemente*) or CARD (*20.15 Uhr*), whereas adverbial or predicative adjectives are tagged as ADJD (*das Gewicht ist gut*; *den betriebswirtschaftlich günstigeren Standort*) or PWAV (*wie wirke ich*).

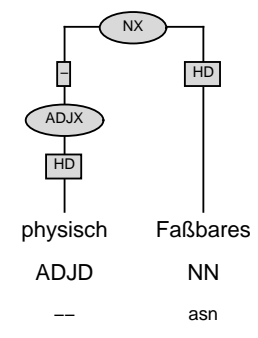
The annotation of superlative and comparative forms is explained in section 7.1 on page 130.

In general, German adjectives are inflected when they are an attribute of a noun. They are not inflected either when they function as a predicative adjective or a premodifier of an adjective or an adverb or when they belong to a small class of noninflected adjectives, e.g. some ancient form such as *gut Wetter* or *lieb Mütterlein* or some adjectives denoting a colour (*mit einer rosa Karte*). All adjectives have to be projected to their phrase level before they are attached to another phrase or to a field.

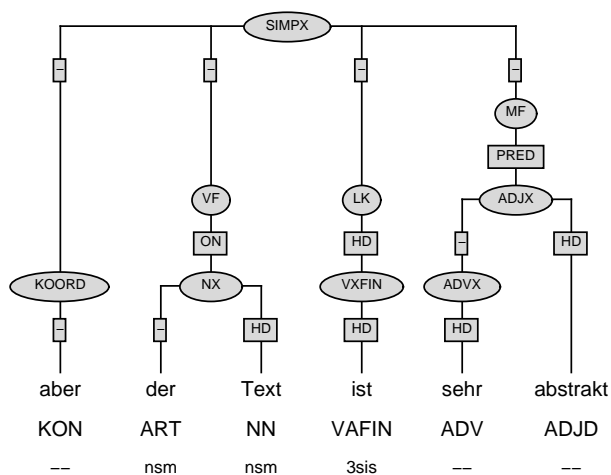
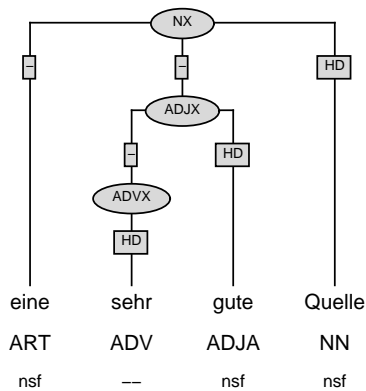




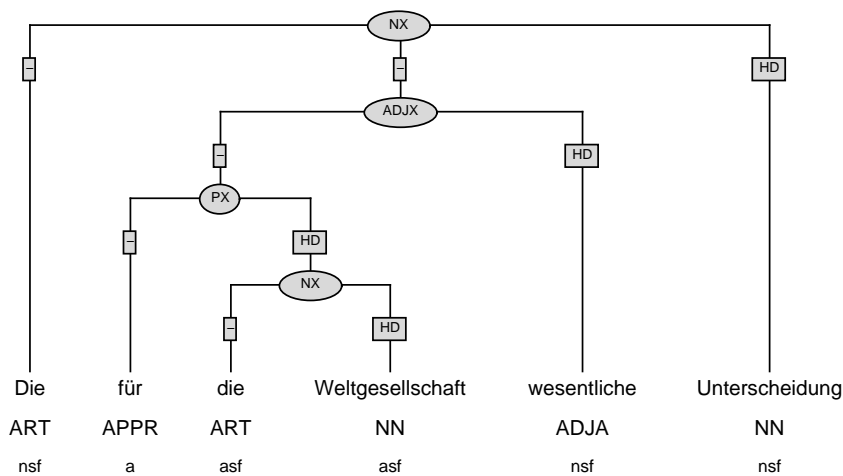
A nominalized adjective like *Fassbares* might be premodified by an adverbial adjective (ADJD) instead of an attributive adjective (ADJA). The former ones do never inflect.



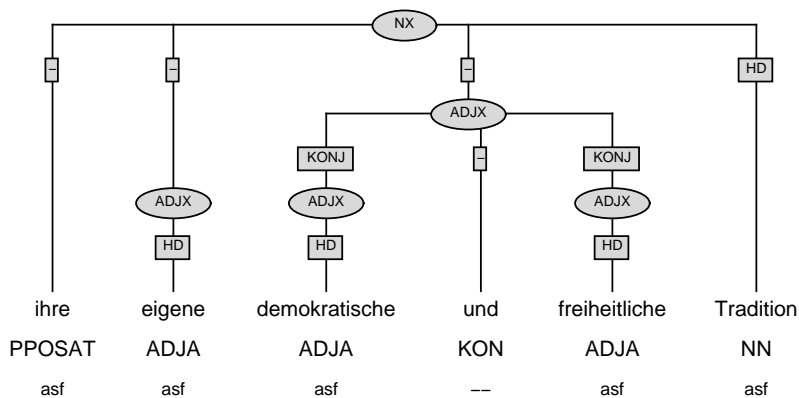
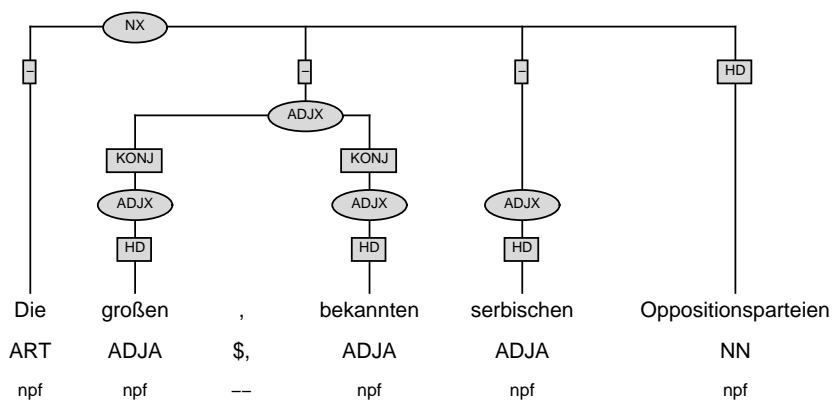
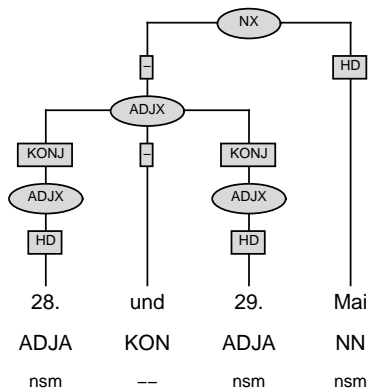
Whenever an adjective is modified by another modifier, the same annotation strategy as for noun phrases is applied, i.e., the modifier is directly attached to the adjectival phrase. The adjectival phrase as a whole is the premodifier of the noun phrase. For instance:



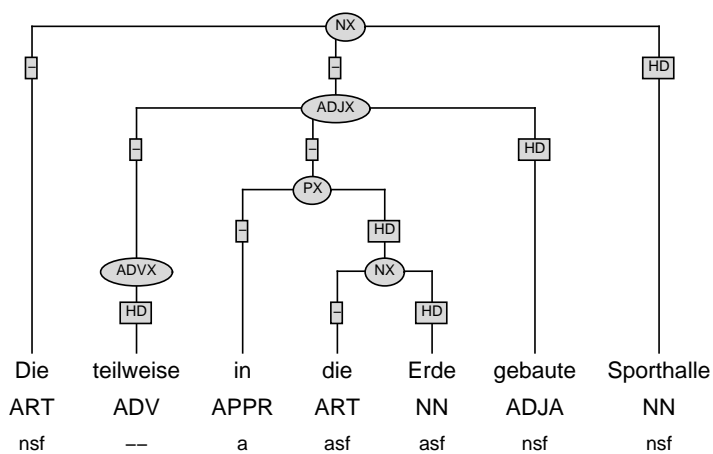
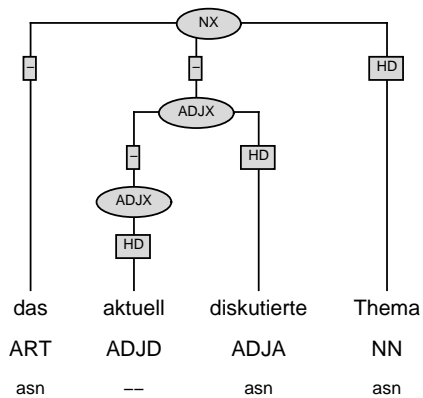
The same holds if an adjective selects an argument. *Für die Weltgesellschaft* is the facultative argument of *wesentlich*. It is directly attached to the adjectival phrase.



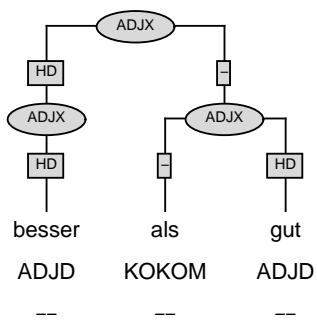
Premodifying adjectives may occur in a linear order and/or as a coordination (cf. 6.5.1) of attributive adjectives:



If the premodifying adjective is deverbal, the adjectival phrase can be of any complexity. In this case, the adjectival phrase has its own internal dependency structure. All elements which depend on the adjective are annotated as its premodifiers. Deverbal adjectives are either attributive or adverbial and predicative respectively, and occur as the present participle or past participle form of a verb.



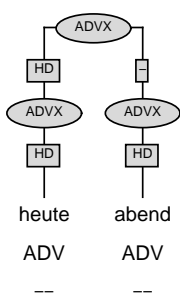
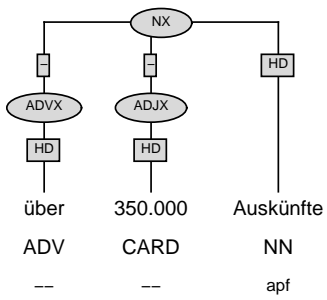
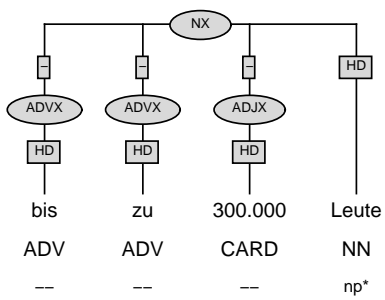
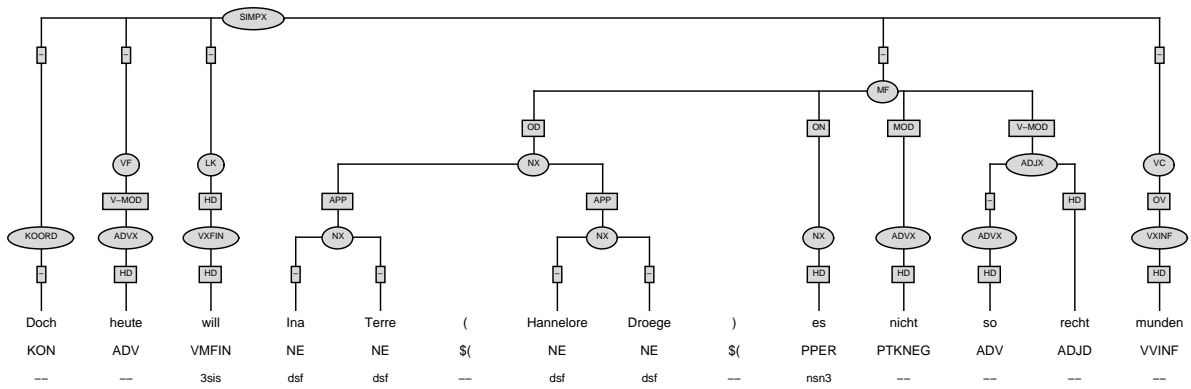
In the following example, postmodification of an adjectival phrase is shown:

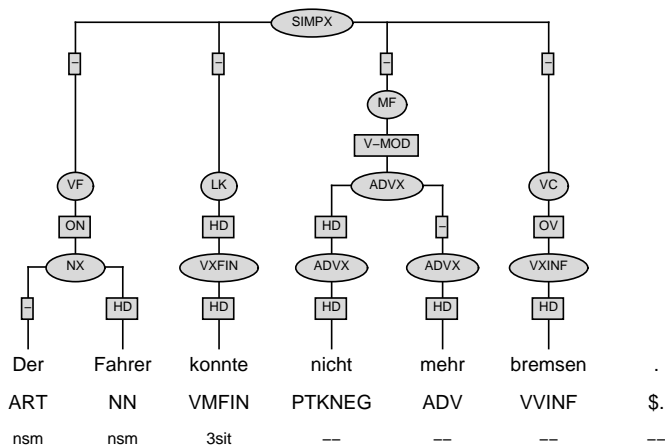


4.6 Adverbial Phrases

Besides adverbials also negation particles (PTKNEG) project to an adverbial phrase. They either occur as premodifiers¹ or postmodifiers or they are directly attached to a field.

¹ *bis zu*, *über* are considered to be ADV rather than APPR because of their semantic meaning.





4.7 Verb Phrases

Whereas finite verb phrases are labelled VXFIN, non-finite verb phrases are labelled VXINF.

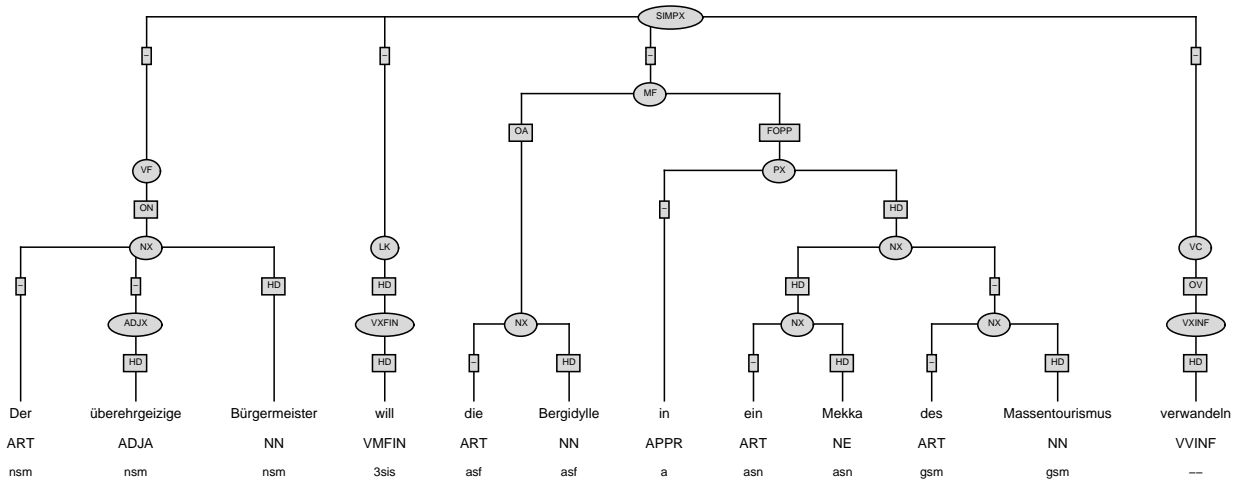
Since infinitives and past participles share certain properties (e.g. exchangeability in *Man hat nur noch das eigene Herz schlagen hören/gehört.*), they are assumed to carry the same phrase label (VXINF). The finite verb in LK as well as the non-finite verbs in VC are always projected to their phrase level. All verb phrases of the verb complex are attached on the same level to form the verb complex. In order to follow the *flat clustering principle*, no internal hierarchy of the verb complex is annotated.

4.7.1 Head of a Sentence and Verb Complex

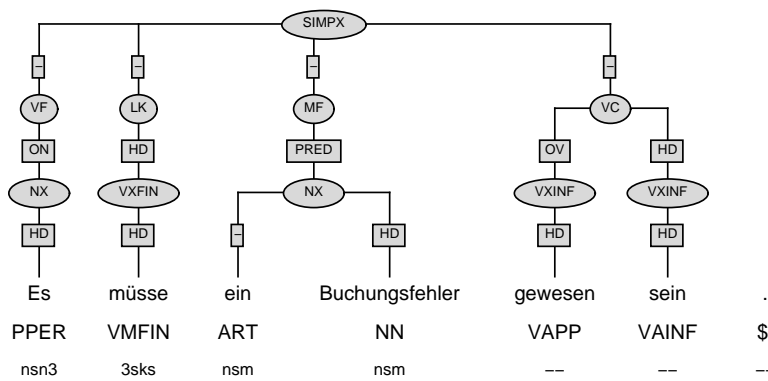
The finite verb which can either appear in LK (verb-first clauses and verb-second clauses) or in VC (verb-final clauses), is always the head of the entire sentence. Non-finite verbal elements belong to VC. If the finite verb is located in LK and if there is more than one non-finite element in VC, the non-finite element which is selected by the finite verb is denoted as the head of VC. All other elements of VC are verbal objects. The head of VC selects the verbal object OV. This verbal object may select another verbal object OV, and so on. In order to denote the dependency relations between verbal objects within the verb complex, we attach a secondary edge label *refvc* between their phrase nodes.

4.7.2 Verb Complexes in Verb-second and Verb-final Clauses

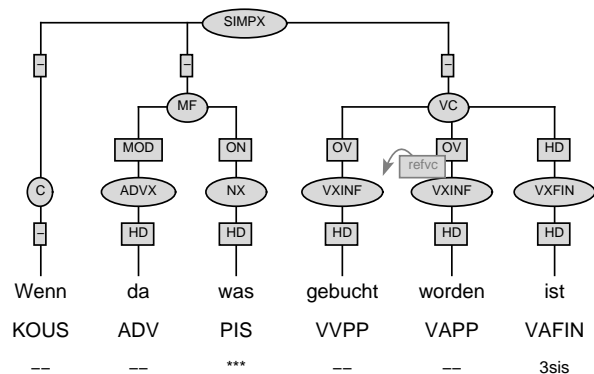
The following example shows a verb-second clause with the head of the sentence in LK and a verb complex consisting of a single non-finite element.

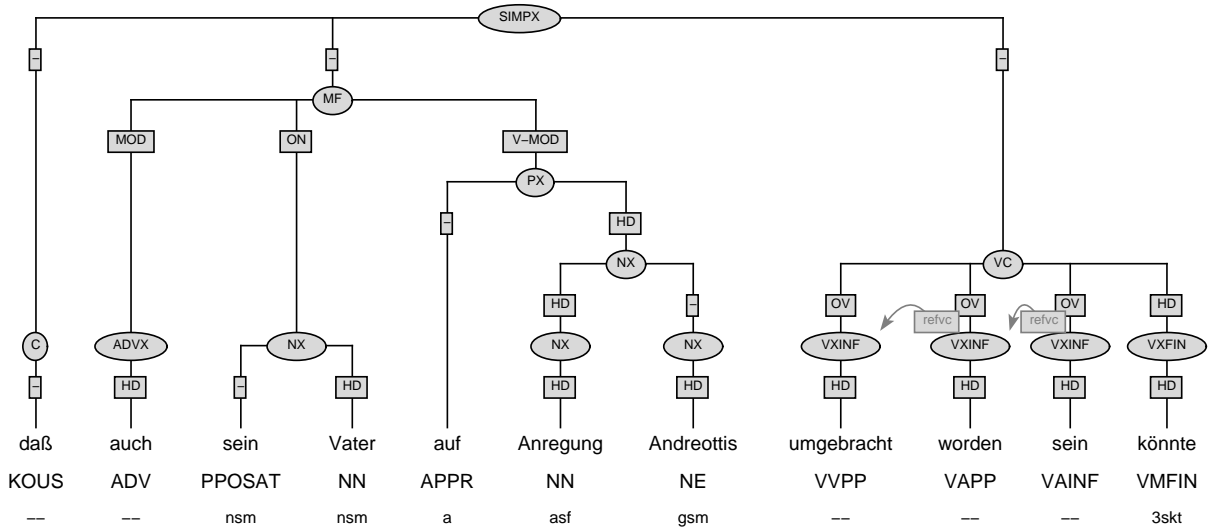


If the verb complex comprises more than one immediate daughter, the one that is selected by the finite verb is the head of VC.

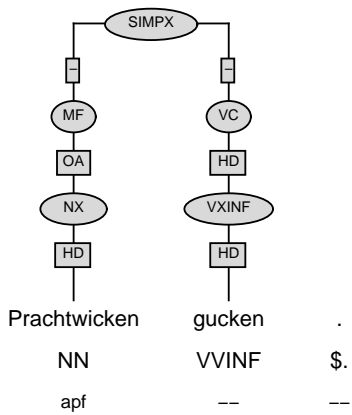


The following trees demonstrate verb complexes with two or more verbal objects. The secondary edge label *refvc* is pointing from the selecting OV to the depending OV.



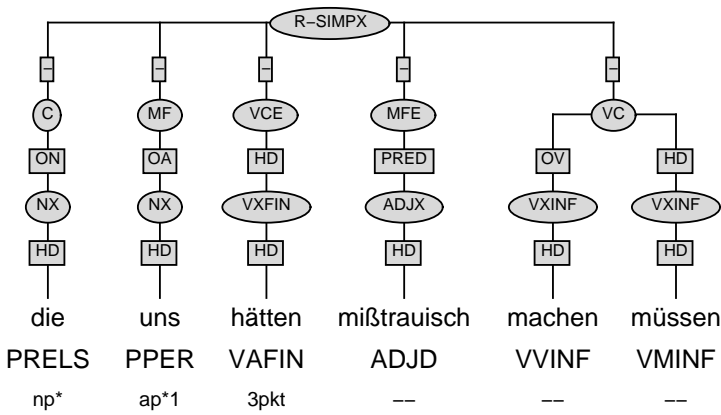
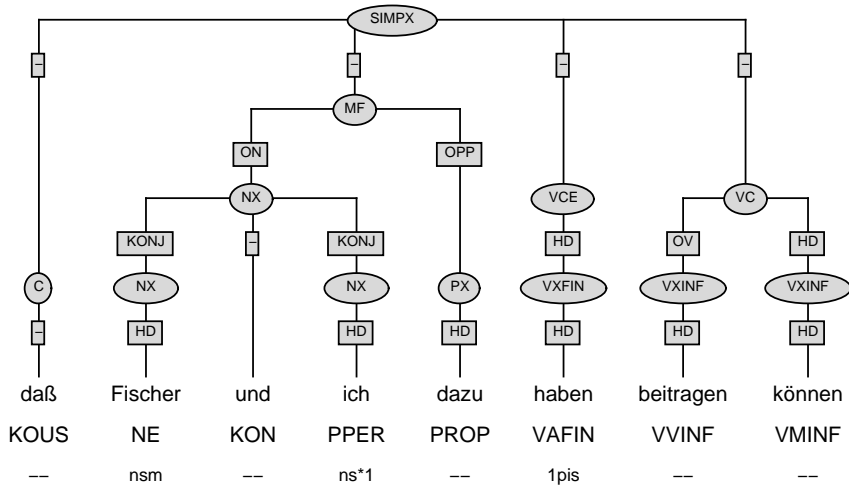


If there is no finite verb at all, the rightmost element of the verb complex (if there is more than one element) is annotated as the head of the sentence. This often occurs in headlines (cf. 5.2 and 7.4).

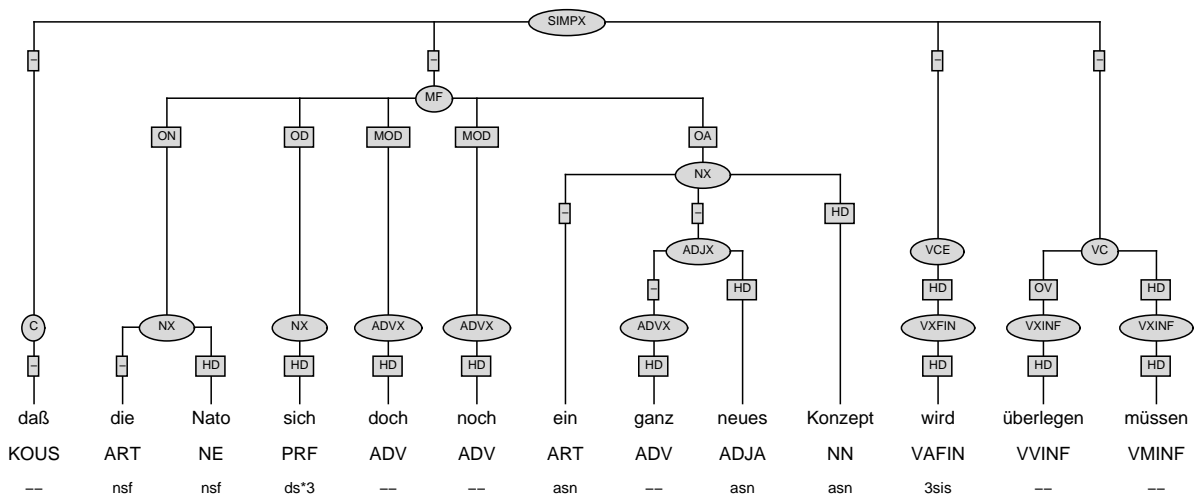


4.7.3 Ersatzinfinitiv Constructions

In order to indicate *Ersatzinfinitiv* constructions, two specific field node labels are introduced. VCE is the node label for the part of the verb complex consisting of the finite verb which subcategorizes for the *Ersatzinfinitiv*. MFE is the node label for the second part of MF between VCE and the second part of the verb complex VC (e.g. [C *die*] [MF *uns*] [VCE *hätten*] [MFE *mißtrauisch*] [VC *machen müssen*]).

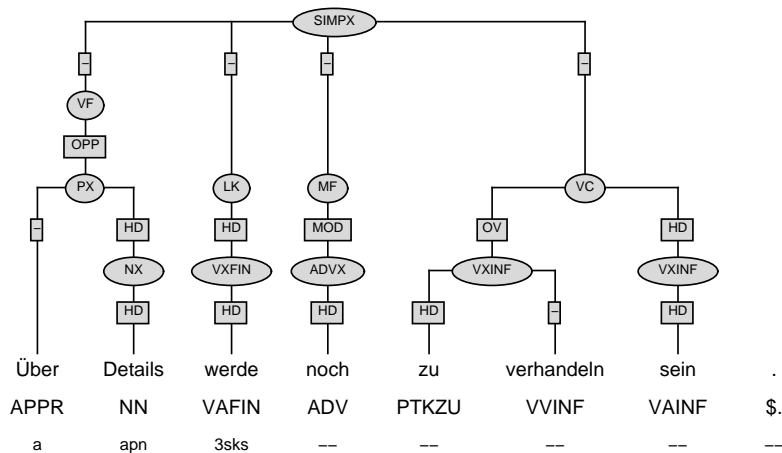
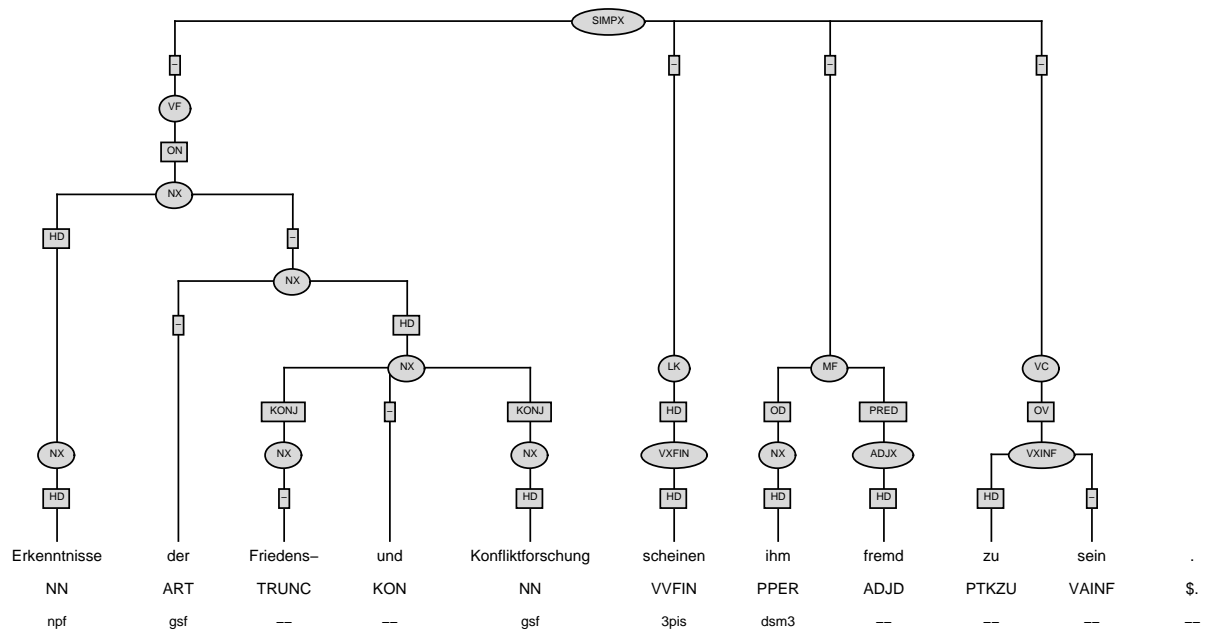


In the example below, the finite verb precedes the non-finite verbs although *müssen* is no *Ersatzinfinitiv*. Since its position corresponds to the position of the finite verb in real *Ersatzinfinitiv* constructions and here also a second middle field is possible, we follow the same annotation strategy.

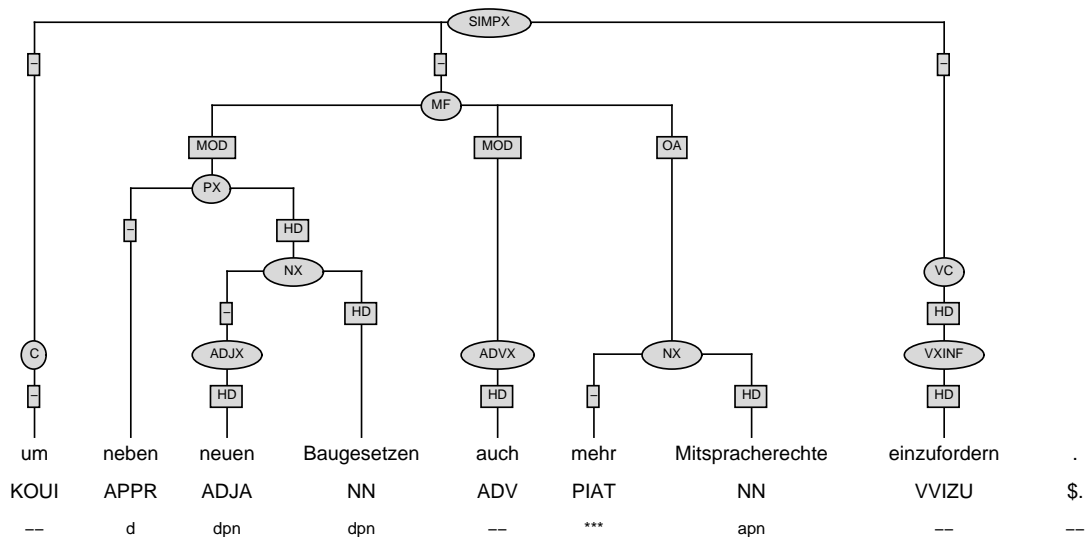


4.7.4 Infinitives with *zu*

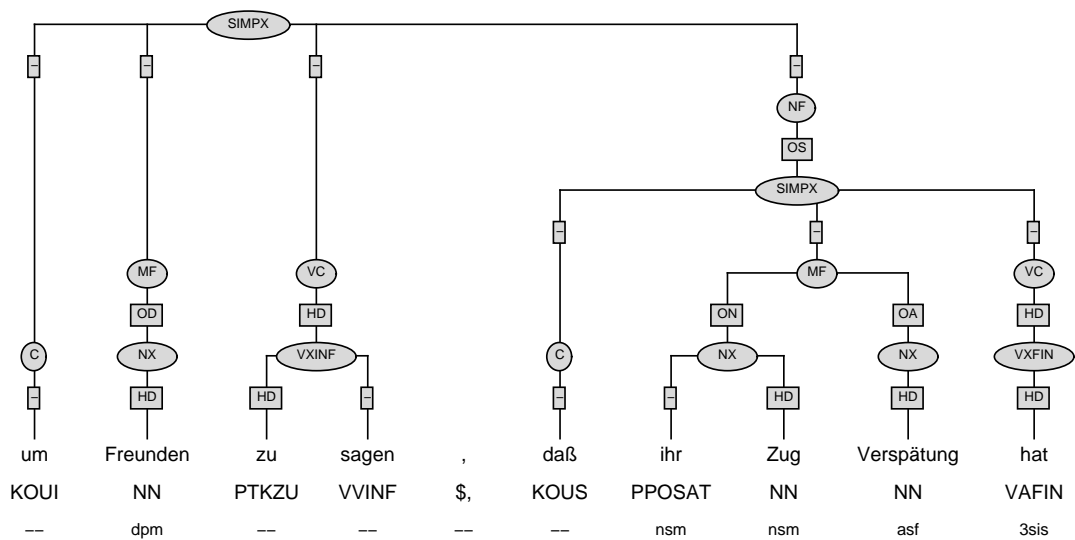
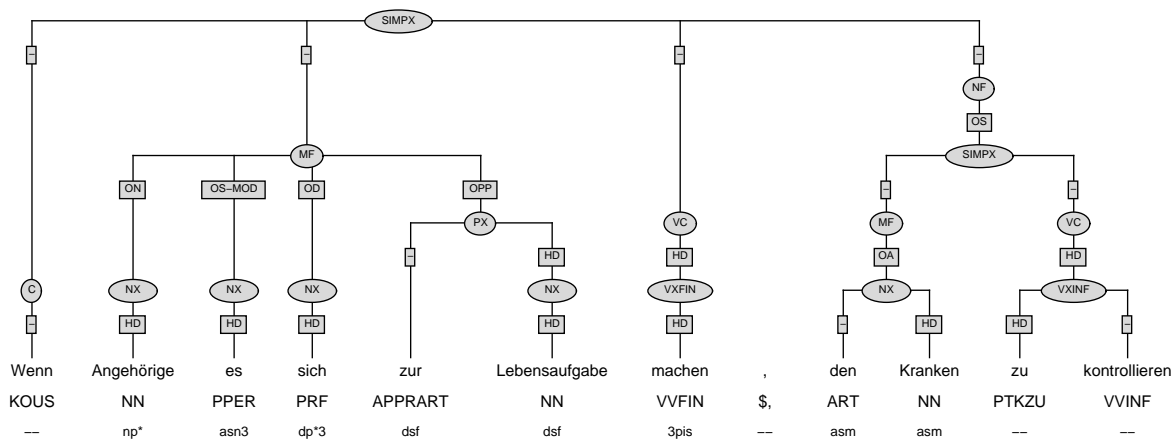
Regarding infinitives with *zu*, *zu* determines the non-finiteness of the verb on its right hand side. This is the reason why *zu* is considered the head of the VXINF whereas the infinitive is assumed to be the complement. Like other infinitives, they occur in the verb complex:



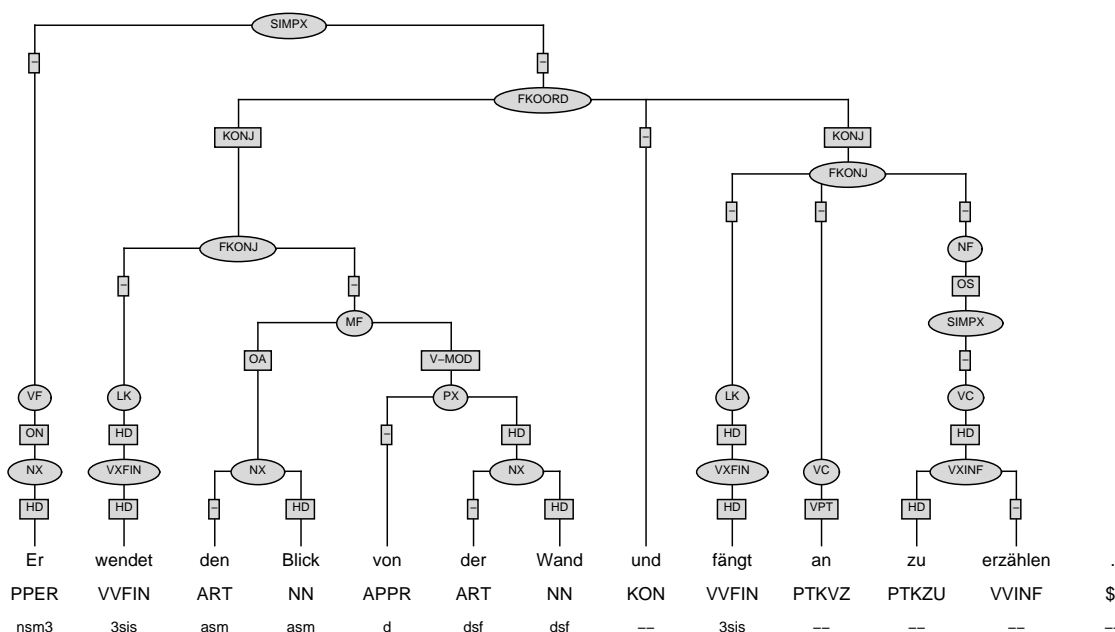
The infinitive with *zu* can also be realized as an infix of the verb. In this case, the verb is tagged as VVIZU. Moreover, it is projected to VXINF with the grammatical function HD:



Beside the examples above, the infinitive with *zu* occurs in optional (in most cases with *um zu*) and obligatory infinitive clauses.



Infinitive clauses can consist of only one verb complex:



4.7.5 Coherency and Incoherency of Verbal Constructions

The notion of coherency attributed to Bech (1955: 57) covers the relation of dependency between adjacent verbal elements, i.e. the relation of subcategorization between a verb and a non-finite verbal complement. Kiss (1995) calls this relation *infinitive Komplementation* (non-finite complementation). Bech (1955: 57) distinguishes between three different modi of obligatory and optional coherency:

1. verbs constructing coherently and incoherently, e.g. *versprechen*, *versuchen*

coherent, extraposition possible:

- a. [*wie er mit kritischen politischen Gegenpositionen umzugehen versteht*]

incoherent, extraposition:

- b. [*wie er versteht,*][*mit kritischen politischen Gegenpositionen umzugehen*]

2. verbs constructing only coherently, e.g. *wollen*, *möchten*

coherent, no extraposition possible:

- a. [*wie er mit kritischen politischen Gegenpositionen umgehen will*]

- b.*[*wie er will mit kritischen politischen Gegenpositionen umgehen*]

3. verbs constructing only incoherently, e.g. *überreden*, *überzeugen*

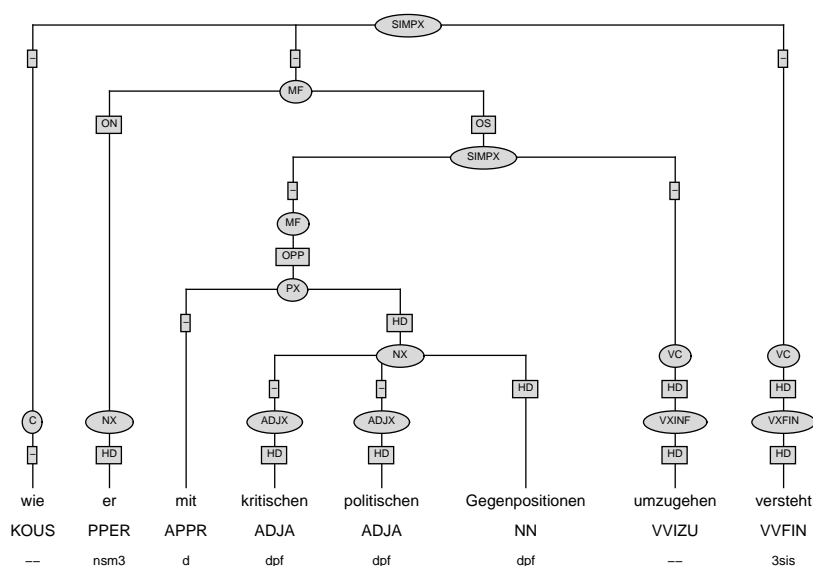
incoherent, extraposition obligatory:

- a. [*wie er sie überredet,*][*mit kritischen politischen Gegenpositionen umzugehen*]

- b.*[*wie er sie [mit kritischen politischen Gegenpositionen umzugehen] überredet*]

Coherent and incoherent constructions of verbs are annotated differently. In case of coherency, the verbal complement is part of the verb complex. In the clause *wie er mit kritischen politischen Gegenpositionen umzugehen versteht*, for instance, the infinitive with *zu* is the verbal object of the finite verb. While in case of incoherency, the verbal complement is annotated as a sentential complement, i.e., *mit kritischen politischen Gegenpositionen umzugehen* in the clause *wie er sie überredet, mit kritischen politischen Gegenpositionen umzugehen* is a sentential object in NF.

We define that a construction is incoherent, if extraposition in NF is possible. That is, whenever it is possible to shift the infinitival complement together with a constituent of MF, which it subcategorizes for, into NF, these elements are annotated as sentential objects. Therefore, the coherent example above (*wie er mit kritischen politischen Gegenpositionen umzugehen versteht*) is annotated with a sentential object in MF since extraposition is possible (cf. the incoherent example 1.b.).



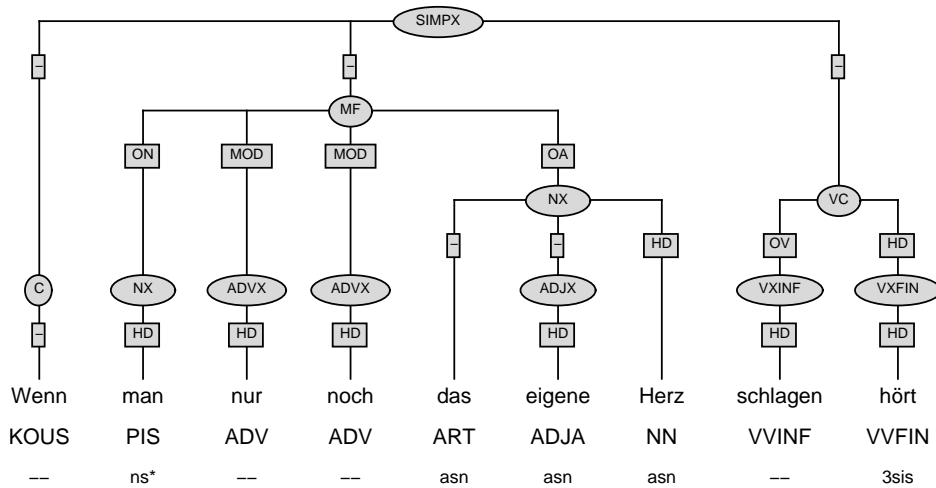
If a complement of the verb within the sentential object is located out of the sentence boundaries, e.g. in the C-field, the secondary edge label *refcontr* gives additional information about the dependency relation (cf. 3.4.6).

4.7.6 AcI Constructions

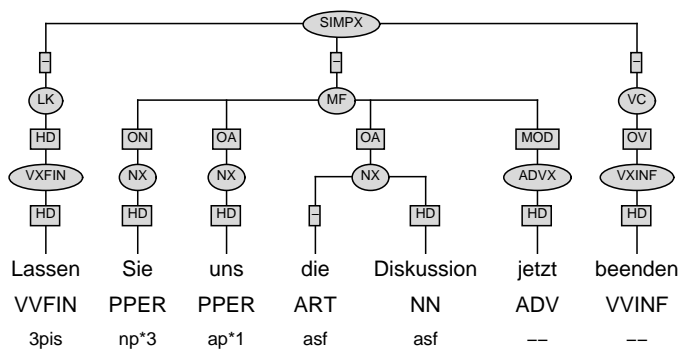
AcI (accusativus cum infinitivo) verbs are a small group of *verba sentiendi* (e.g. *sehen, hören, fühlen, spüren*) which subcategorize for an accusative and an infinitive. The verbs *lassen, machen, heißen* have a modal verb like reading in which they also select an accusative and an infinitive.

The infinitive itself subcategorizes for complements with respect to its valency but its subject is realized by an accusative which is the direct object of the AcI verb.

Since AcI constructions are coherent infinitive constructions in which extraposition is not possible (cf. (Eisenberg 1999 2001), p.355), the AcI is not annotated as a sentential object (* *wenn man nur noch hört das eigene Herz schlagen*). The infinitive as the verbal object of the AcI verb is located in the verb complex and the accusative is realized as OA in MF.



As a consequence of this analysis, we annotate two accusative objects (OA) if the AcI construction comprises a transitive infinitive verb such as *beenden* in the following example. *Uns* functions as its subject and *die Diskussion* as its direct object. Both are in accusative case and both are labelled OA.



4.7.7 Imperatives

Imperative verbs have only one singular and one plural form and are not inflected concerning the grammatical category person. Their form corresponds to second person singular and plural verbs which are tagged as VVIMP or VAIMP.

Warte mal! (warte/VVIMP:s)

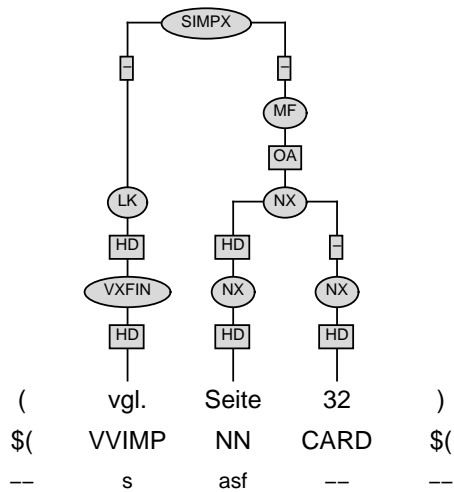
instead of

Wartest du mal? (wartest/VVFIN:2sis)

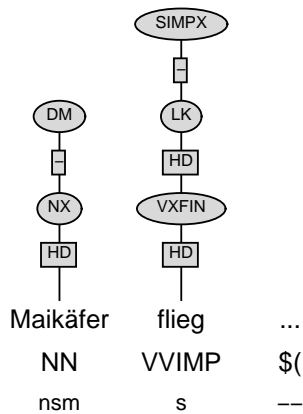
It is important to keep apart imperative sentences from imperative verbs. An imperative sentence does not need to comprise an imperative verb form as is shown in the following examples

Warten Sie mal bitte! (warten/VVFIN:3pis)

Bitte warten! (warten/VVINF:-)

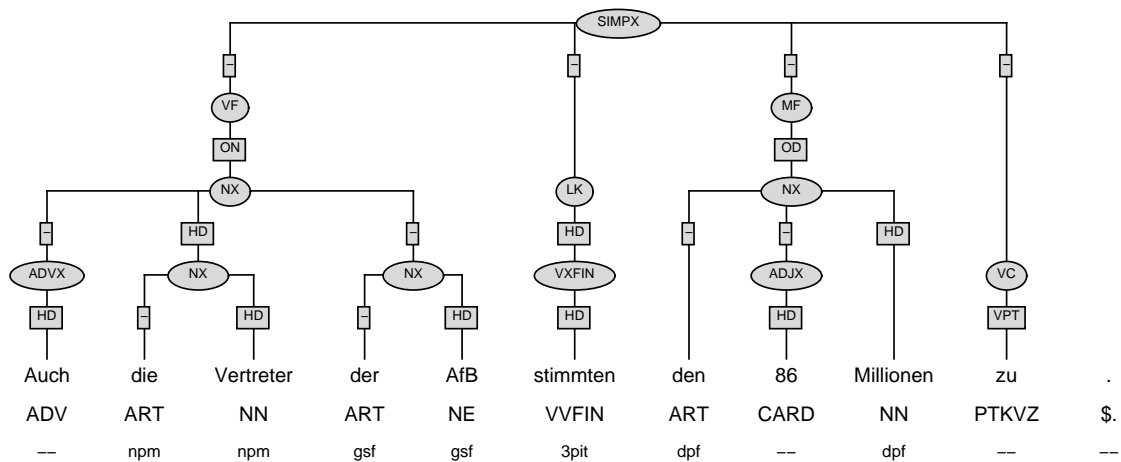


Normally imperative verbs are lacking the subject, but the addressed person can also be mentioned to stress the utterance:

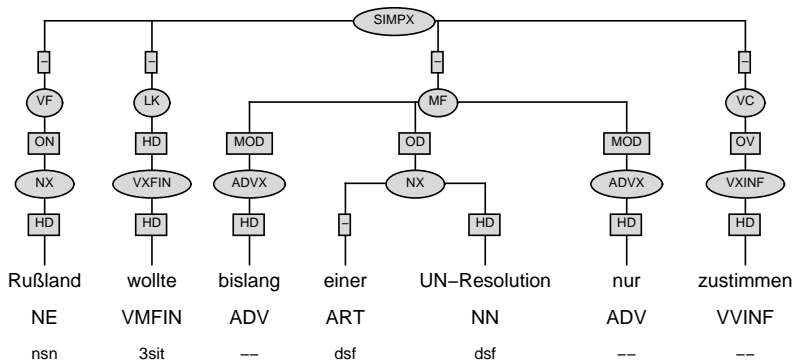


4.7.8 Particle Verbs

Separable verb particles are tagged as PTKVZ and annotated with the edge label VPT:



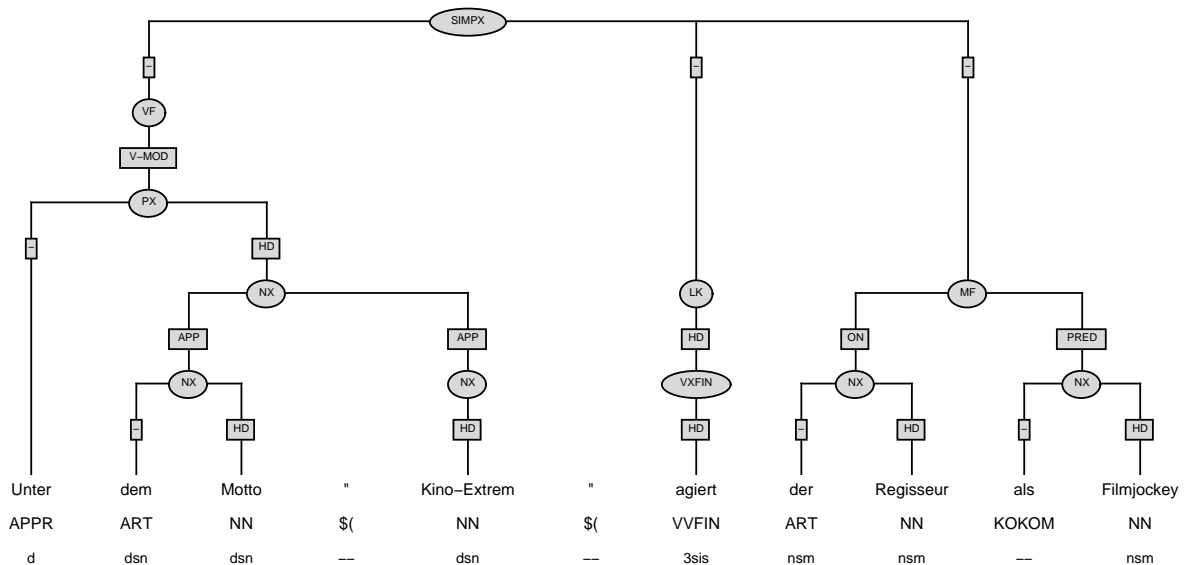
In verb-final clauses, the particle verb occurs unseparated within the verb complex:



4.7.9 Verbs with Predicate

Typically, the complement type PRED (predicate) occurs with verbs like *sein*, *haben*, *scheinen*, *aussehen*, *sich anhören*, *klingen*, etc. PRED is annotated, if the following conditions apply:

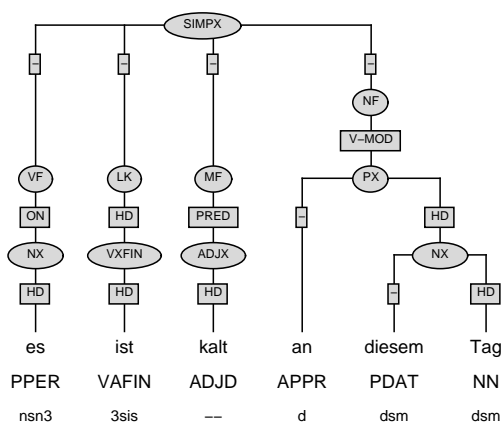
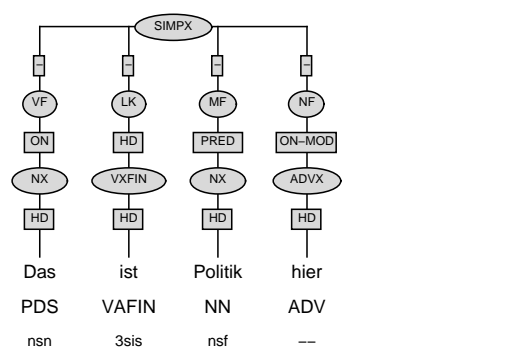
- if it is not possible to determine the case of the constituent in question properly (e.g. *gut* in *Das ist gut.*)
- if the constituent in question actually predicates the subject, i.e. the subject is characterized as having the property expressed by PRED (e.g. in *Die Ursache war unklar. Die Ursache* is characterized by the property of being unclear)
- many PRED verbs are raising-verbs (subject without theta-role)
- if *als*-phrases are selected by the verb they are labelled as PRED (e.g. *Unter dem Motto Kino-Extrem agiert der Regisseur als Filmjockey.*)



Some examples for verbs that take predicates: *recht sein*, *recht haben*, *leid tun*, *frei sein*, *fertig sein*, *sich gut/schlecht treffen*, *gut/schlecht finden*, etc.

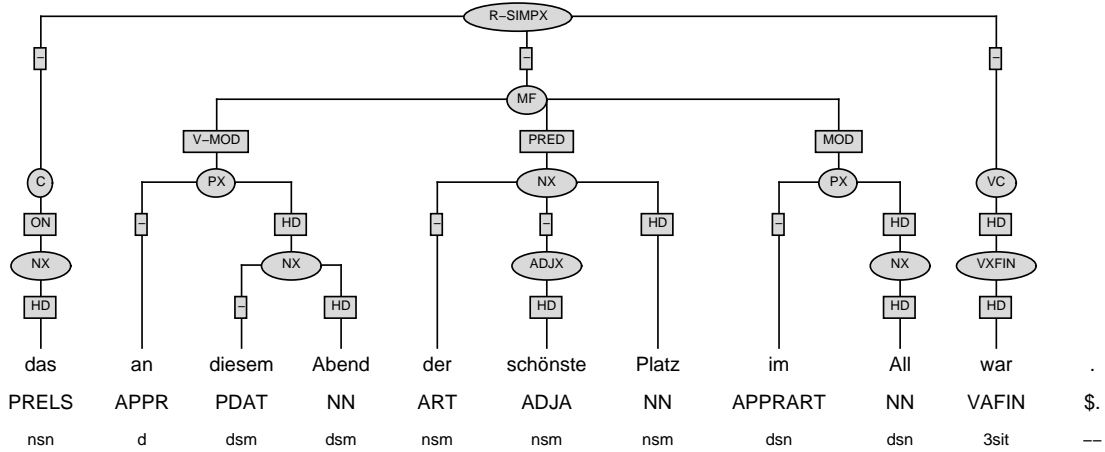
PRED verbs have to be distinguished carefully from verbs occurring with ordinary modifiers (e.g. ON-MOD, V-MOD) such as *gut passen*.

With respect to topological fields, note that PRED usually marks the border between MF and NF, i.e., whatever constituent occurs on the right hand side of PRED belongs to NF. In general, this constituent is an adjunct which PRED does not subcategorize for:

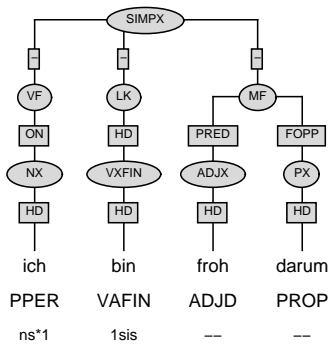
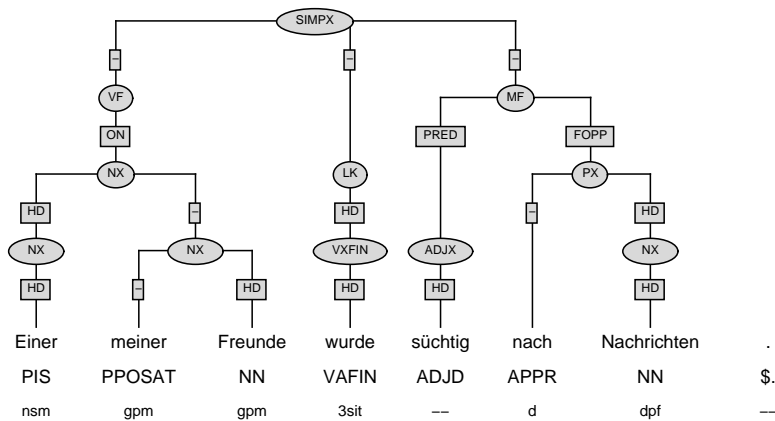


But there are exceptions in which PRED does not necessarily constitute the border between MF and NF:

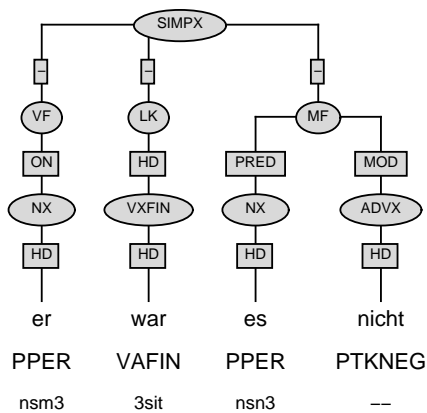
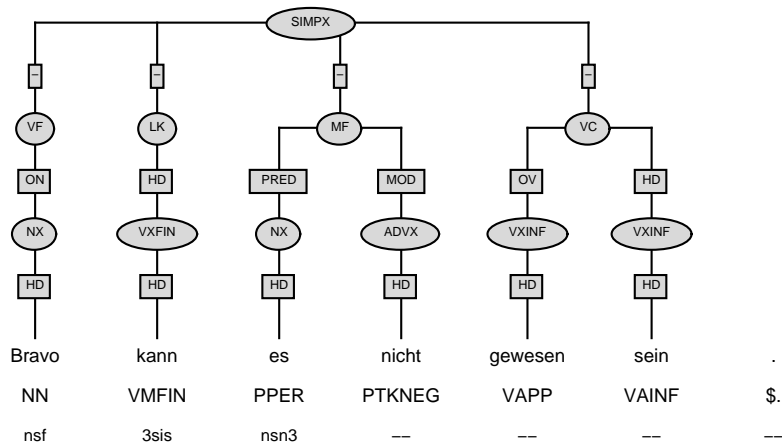
- Another constituent may occur between PRED and VC, for instance, if an ambiguous modifier follows PRED.



- PRED subcategorizes for the constituent that follows it. Complements of PREDs are always attached to a field since they are assigned a grammatical function within the sentence structure (cf. 8.1):

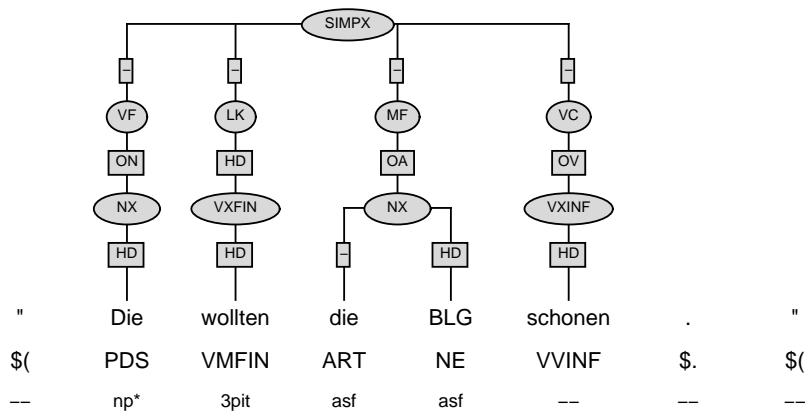


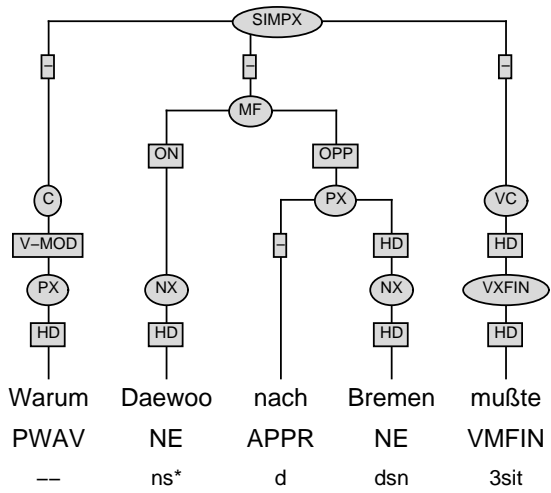
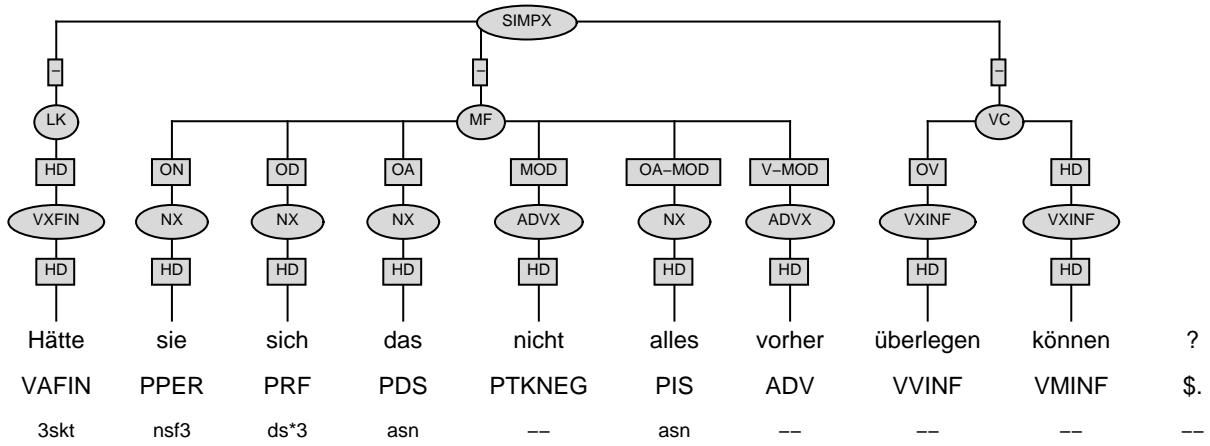
- Because of the word order rule that pronouns in MF have to precede other constituents, PRED might not be the last element in MF if it is a pronoun:



4.7.10 Modal Verbs

Modal verbs are always tagged as VMFIN or VMINF regardless of their use as an auxiliary or a main verb. If a modal verb functions as an auxiliary verb, it is projected like any other auxiliary verb. If a modal verb is the main verb of a sentences, verbal modifiers refer to the modal verb in the same way as they refer to other main verbs:





Chapter 5

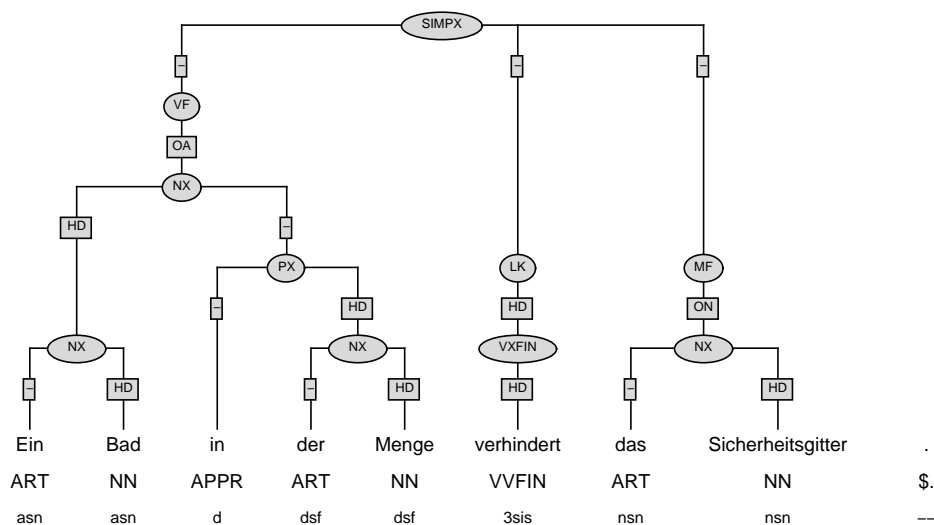
Attachment Principles for Phrases

5.1 Attachment to Fields

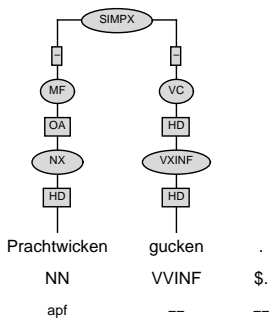
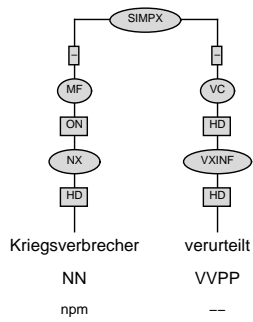
Phrases are attached to the topological field in which they occur. Their edge labels denote their grammatical function within the sentence structure. In LK and VC there can only occur verb forms, separable verbal prefixes, or infinitive particles. LK and VC mark the beginning and the end of MF (cf. 3.2).

5.2 Attachment of Ambiguous Complements

The partially free word order and the morphological properties of German can cause ambiguity concerning the grammatical function of a constituent. In the following example, the syntactic structure does not give any information about case assignment. Both noun phrases can be identified as ON or OA:



Headlines like the following are lacking the finite verb. Therefore, in the first example it cannot be decided if it is an active or a passive construction, i.e., if the noun phrase is ON or OA. The second example is an active construction, but again the noun phrase can be both, ON or OA:



Since we do not assign specific edge labels for ambiguous complements, we formulate the following *preference principle for case assignment*:

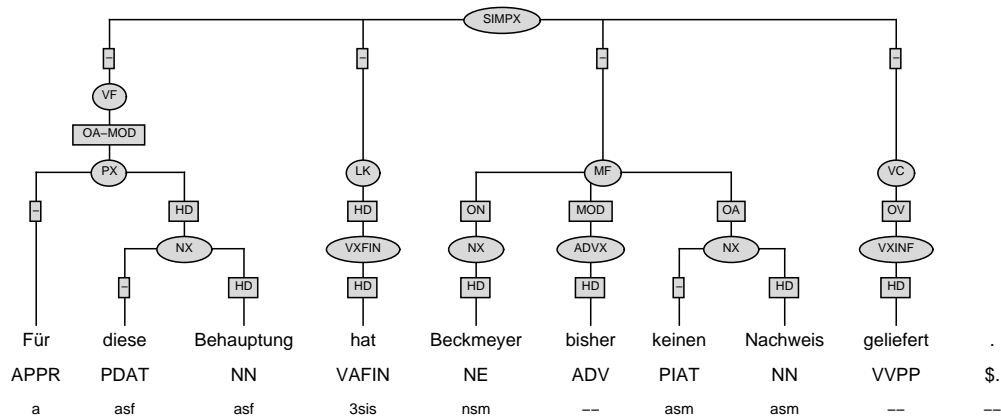
Preference principle for case assignment:

If case assignment is ambiguous, we decide on the more plausible grammatical function and on the more plausible sequence of grammatical functions respectively. The main criteria for the decision are the unmarked word order and the semantic content.

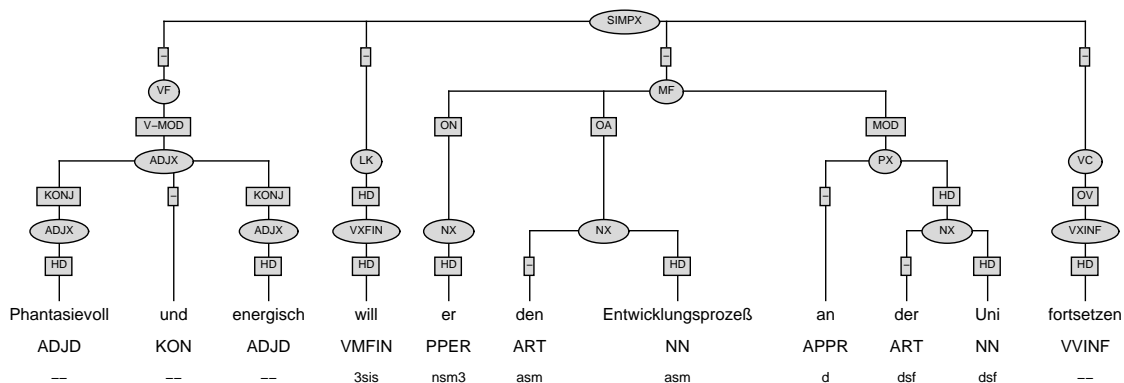
Therefore, in the first example above, OA appears in VF whereas ON has its position in MF. For elliptical headlines, we assume a passive construction if the verb in VC is a past participle and an active construction if the verb in VC is an infinitive (cf. 4.7.2 and 7.4).

5.3 Modifier Attachment

Modifiers either modify one specific constituent or more than one constituent. The scope of modification can even range over the whole sentence structure. Therefore, they are either unambiguous or ambiguous. An unambiguous constituent that modifies just one other constituent within a tree structure is either adjacent or discontinuous. In the first case, it is immediately attached to the constituent which it modifies, concerning the attachment rules for phrases. In the second case, the dependency, which can even go beyond the border of topological fields, is indicated by X-MOD edge labels, which express the non-ambiguity of the modifier (e.g. OA-MOD is the modifier of OA). Thus, edge labels like OA-MOD, V-MOD, OPP-MOD, MOD-MOD, etc. express that the respective constituent modifies only one other constituent in the sentence (OA, V, OPP, a modifier, etc.) which is not adjacent:



If a modifying constituent is ambiguous (i.e. it modifies more than one constituent, the entire sentence, or a constituent that occurred in previous sentences), it is attached to its topological field and given the ambiguous edge label MOD to preserve ambiguity. In the following example *an der Uni* either modifies the accusative object *den Entwicklungsprozeß* or the verb *fortsetzen*:



We formulate the following definitions for MOD and X-MOD:

Definition of MOD:

A constituent is called MOD, if it cannot be assigned a more specific label, either because it is ambiguous or because there is no more specific label (e.g. for sentence modifiers or for constituents that refer to some sentence external expression). Sometimes it is difficult to determine whether a modifier is definite or not. In cases of doubt, modifiers are marked as ambiguous (MOD) rather than as definite modifiers.

Definition of X-MOD:

X is a variable that can be replaced by labels for syntactic categories like OA, OPP, MOD, V. X-MOD marks long-distance modification which is unambiguous, e.g. relative clauses (*Aber es gäbe (intelligente Lösungen OA)*, (*die kein Geld kosten OA-MOD*)).

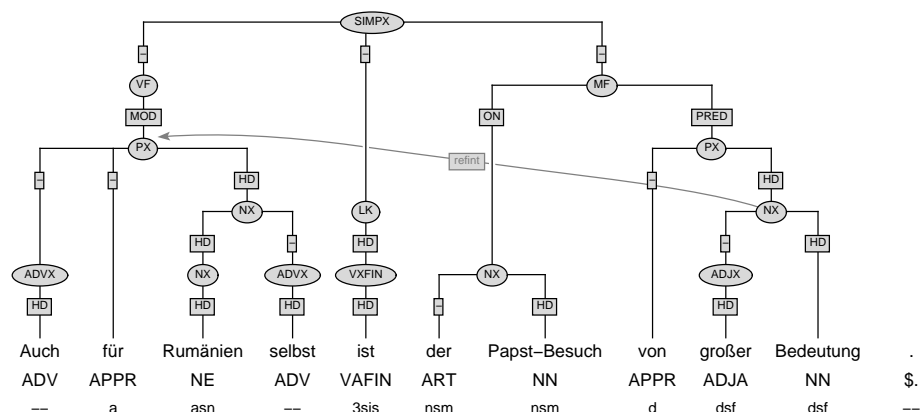
Typical MODs and V-MODs:

Generally, modifying subclauses (e.g. *Katastrophenstimmung herrscht erst, [wenn nichts mehr zu verheimlichen ist] (MOD)*.) are MOD because they modify the complete main

clause. Modifying particles and adverbs like *da*, *dann*, *auch*, *eigentlich*, *ja*, *vielleicht*, *auch*, *natürlich* usually show attachment ambiguity and therefore are annotated as MOD. Only if they unambiguously express the modification of the verb (e.g. *Das Buch liegt da*. or *Er geht auch*.) they carry the edge label V-MOD. Pronominal adverbs (PROP) like *dabei*, *dafür*, *trotzdem*, *deswegen*, *hierauf*, etc. are either ambiguous (e.g. *Dabei (MOD) erscheinen Sie in anderen Verlagen*.) or unambiguous [e.g. *Er achtet dabei (V-MOD) auf alles*.) Non-pronominal adverbs such as *vorher*, *später*, etc. in most cases give temporal or local information. Thus, they are rather V-MOD than MOD.

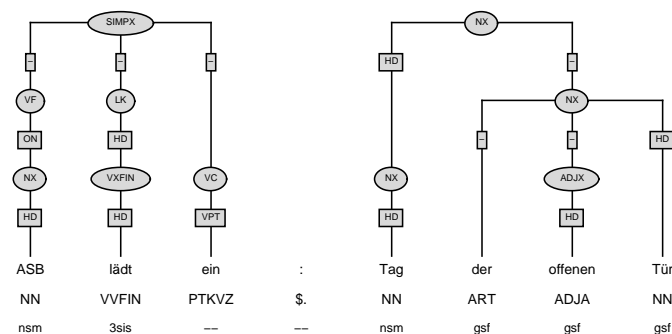
5.3.1 Modifier Attachment in the Initial Field

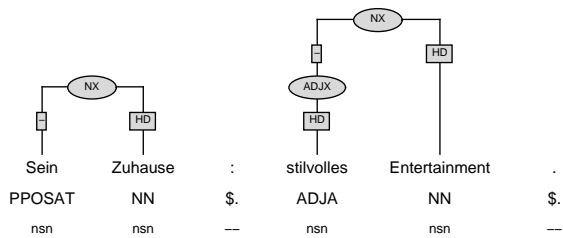
Since only one constituent is allowed in the initial field, all elements preceding and following the head are attached as premodifiers (low attachment) or postmodifiers (high attachment) according to the attachment rules explained in 4.1.



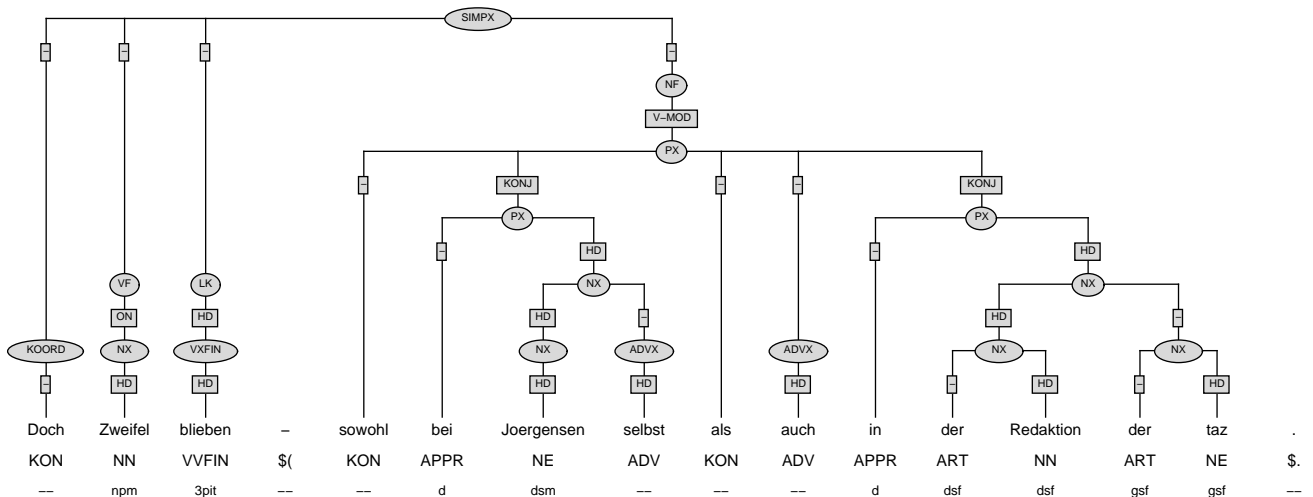
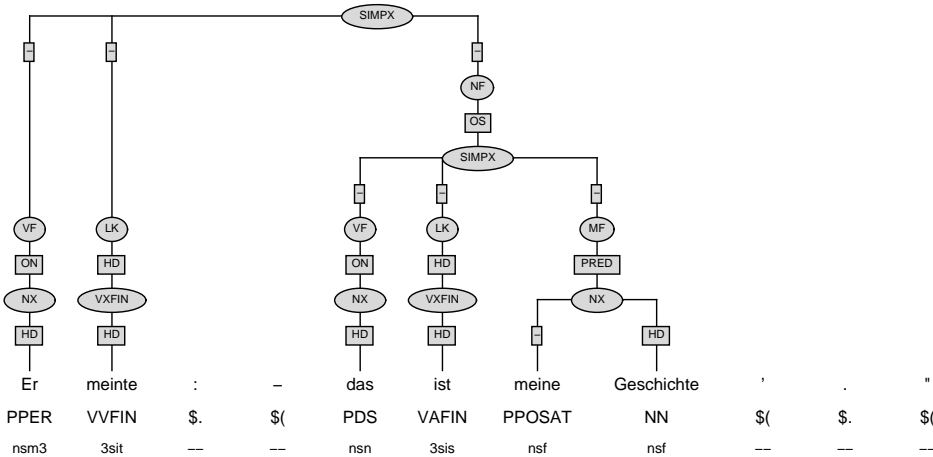
5.3.2 Attachment across Punctuation Marks

The punctuation marks : and - and ... separate a syntactic construction within a unit unless there is no syntactic dependency relation between the two parts (cf. 3.4.5) like in the following:





Attachment is necessary if the part following the punctuation mark has a grammatical function within the sentence structure:



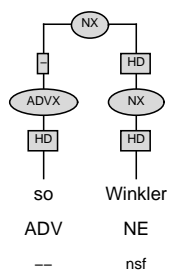
5.3.3 Ambiguous Modifiers in Isolated Phrases

Since isolated phrases (cf. 3.4.5) do not consist of topological fields, ambiguous modifiers (MOD) have to be attached to the phrase itself. The isolated phrase is projected one level higher and the modifier is attached on this higher level. Thus, the information about ambiguity can be preserved even without topological fields or explicit MOD labelling, just by the existence of yet another projection level of the phrase.

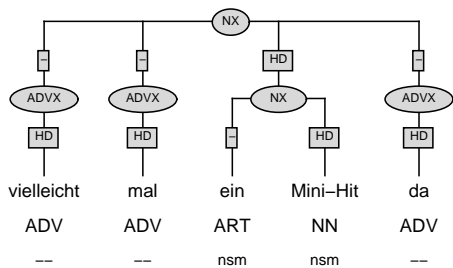
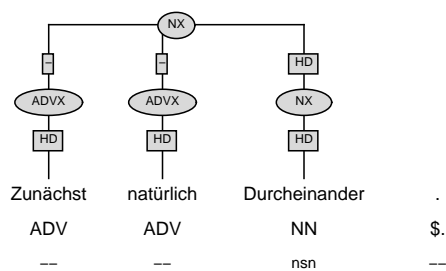
The overall attachment strategy has been chosen in order to keep syntactic structure

flat and to be able to preserve attachment ambiguity where necessary.

In the following examples, *so* may refer to something that is implicit or has been mentioned before:



If there is more than one ambiguous modifier in an isolated phrase, all of them are attached on the next higher level. The mother node of this isolated phrase is marked with the node label of the modified phrase.



Chapter 6

The Annotation of Sentences

The approach of topological fields supports the *flat clustering principle* inasmuch MF and NF allow for more than one constituent being attached to the same field node. The field nodes form a level of annotation between the phrase level and the sentence level. The last step to complete a sentence structure is to attach the field nodes to the highest annotation level of the whole structure: the root node.

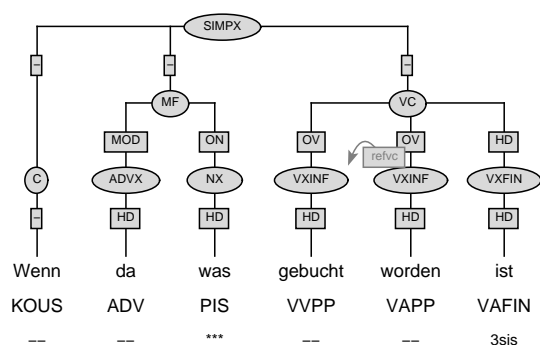
In the following sections, the annotation of sentence structures will be demonstrated.

6.1 Sentence Initial Fields

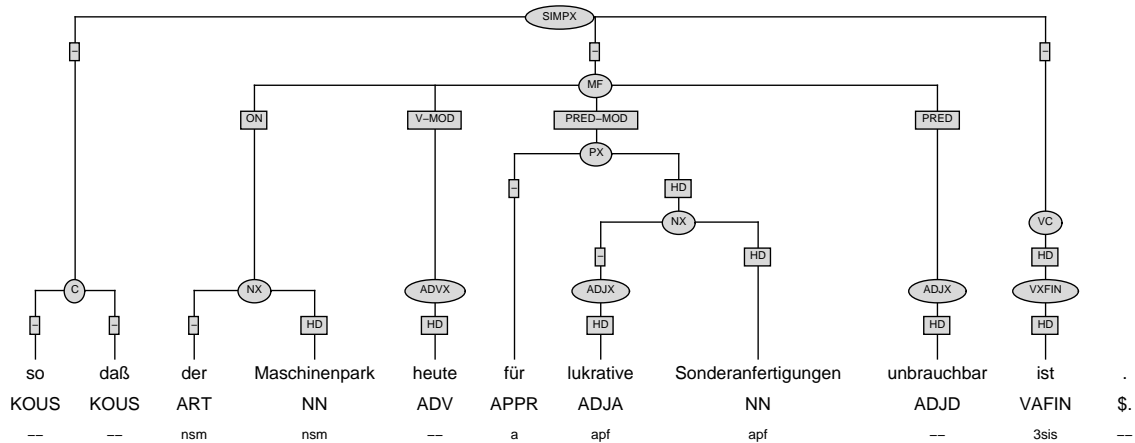
6.1.1 The C-Field in Verb-Final Clauses

The C-field (complementizer field) is the field for subordinating conjunctions KOUS (e.g. *daß, wenn, da, weil, ob*), KOU1 (e.g. *um (+zu)*), relative pronouns (PRELS), interrogative (PWAV) pronouns and (complex) interrogative or relative phrases. Thus, it only occurs in verb-final clauses.

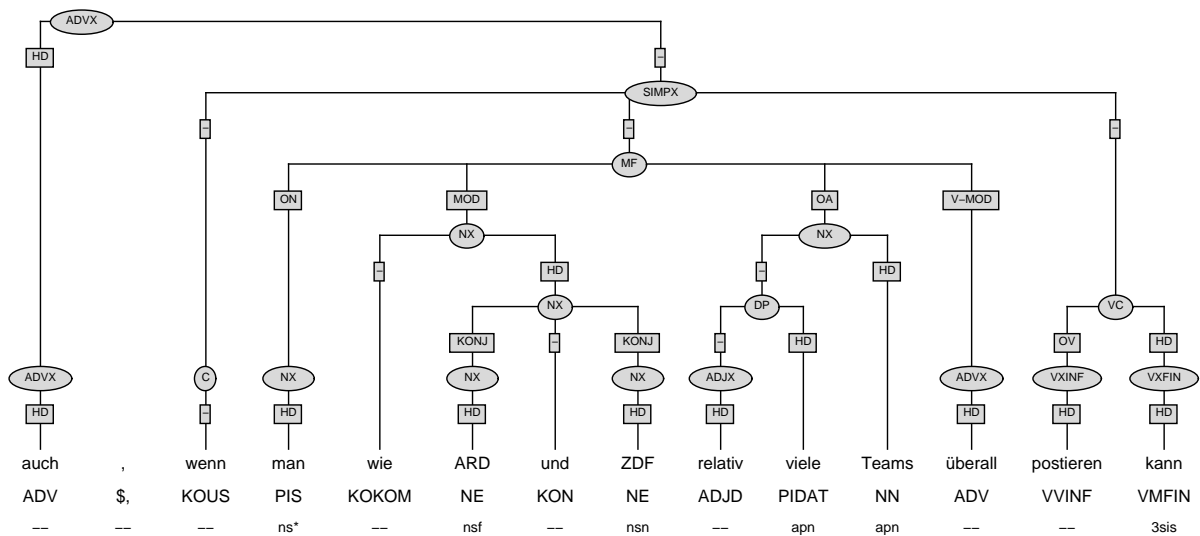
In case of a conjunction, we directly project to the C-field:



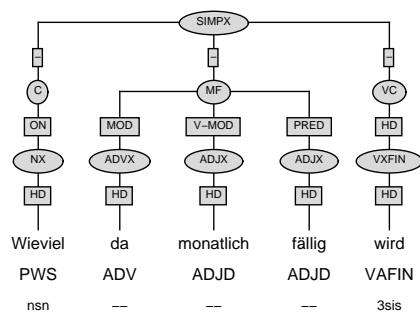
There are conjunctions in German which consist of two elements (e.g. *so daß* and *als ob*). Both of them are also directly attached to the C-field, while none of them carries a head label.

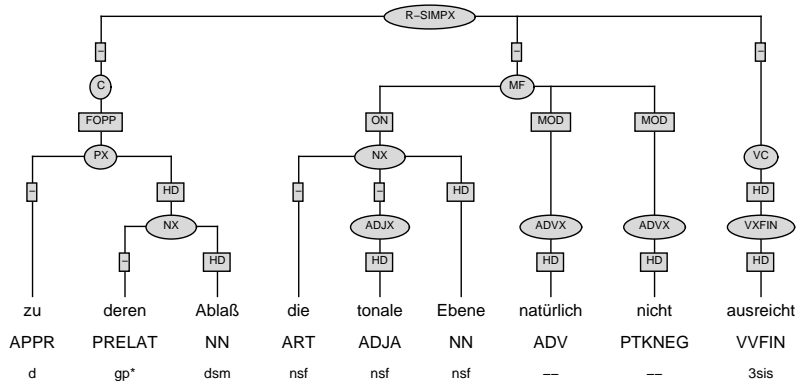


Since C generally does not contain more than one constituent, the adverb *auch* in the following example is not supposed to occur in the C-field together with the conjunction *wenn*. The *wenn*-clause is annotated as the modifier of the adverbial phrase *auch*, i.e., the adverbial phrase subcategorizes for the verb-final clause.



If the constituent in the C-field is a pronoun or a complex phrase, it is first projected to the phrase level and then projected to the C-field. The edge label below the C-Field denotes the grammatical function of this constituent.

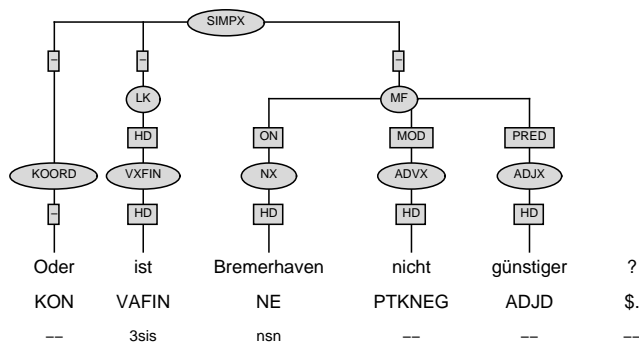
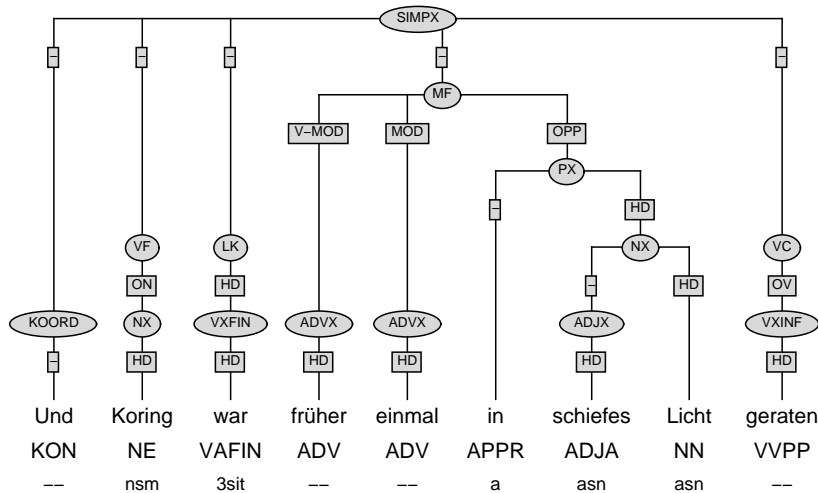




6.1.2 The KOORD-Field in all Clause Types

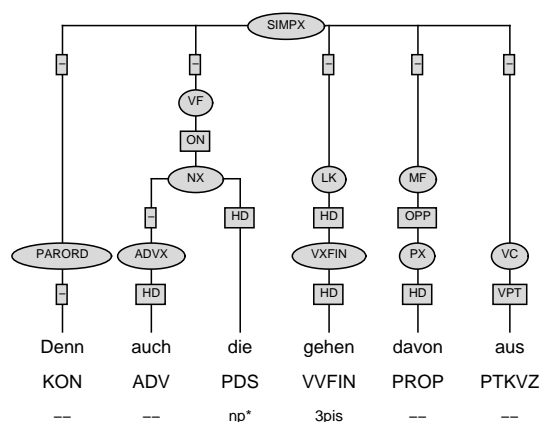
The KOORD-field is optionally the left-most field of all clause types (V-1, V-2, V-end). Therefore, it can only occur at the beginning of a syntactic unit (cf. 3.4.3).

For verb-second clauses, it can be regarded as an alternative field to the PARORD-field. The KOORD-field contains coordinative particles like *und*, *oder*, *aber*, etc. (cf. Höhle (1986)). Here are two examples of different clause types:



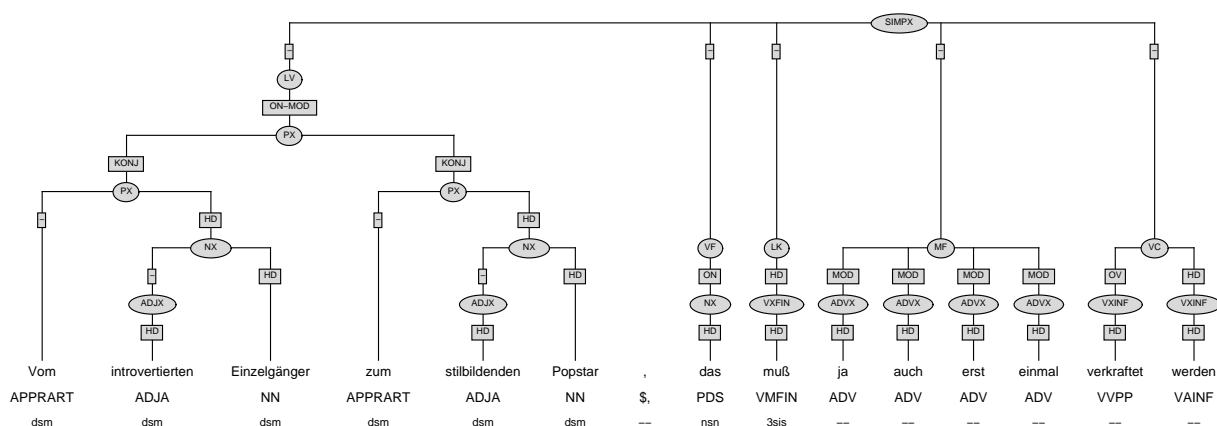
6.1.3 The PARORD-Field in Verb-Second Clauses

PARORD is an alternative field to KOORD for verb-second clauses only. Typical PARORD expressions are *denn, weil*¹:



6.1.4 Resumptive Constructions: The LV-Field

Resumptive constructions are analysed as suggested by Höhle (1986) and Kathol (1995), by using the field LV (*Linksversetzung*) which is located to the left of VF. In general, the LV-field is not restricted to one constituent. The typical feature of a resumptive construction is that there is a (pronominal) constituent somewhere in the sentence, on the right hand side of the LV-field, which refers back to the expression within the LV-field. Therefore, we use the X-MOD label to indicate this kind of long-distance dependency.



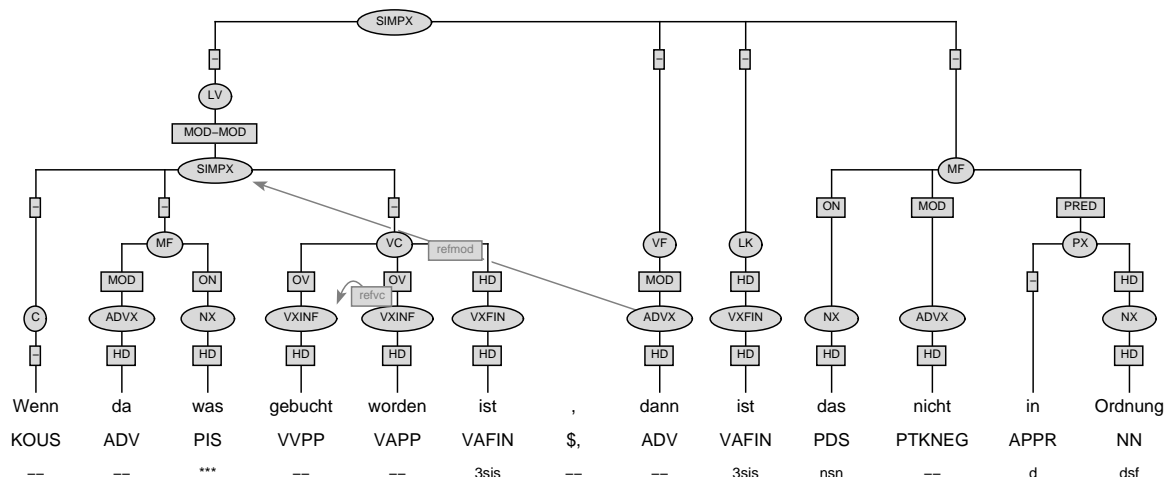
Grammatical functions within a LV-construction are assigned according to the following principle:

- The LV-constituent is licensed by some (pronominal) constituent within the core sentence. The core sentence exceeds from VF to NF. Therefore, the licensing constituent is considered to be modified by the constituent within the LV-field.

¹ *weil* can occur in verb-second and in verb-final clauses. In the first case, it is in the PARORD-field, in the latter case, it belongs to the C-field.

For instance, ON-MOD is licensed by ON like in the example above, which is also in strong accordance with the assumption that the original position of the subject in verb-second clauses is VF.

In constructions with *wenn ... dann ...*, the *wenn*-clause, which is semantically a precondition to the *dann*-clause, is in the LV-field in correlation with *dann*. Therefore, *dann* (MOD) refers back to the *wenn*-clause (MOD-MOD):



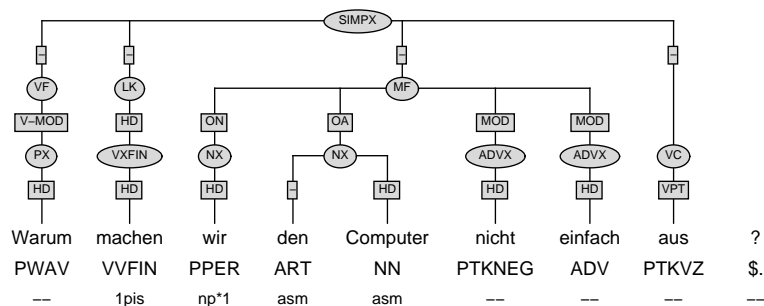
If *dann* is not present in the matrix clause, the *wenn*-clause occurs in VF. In this case, the *wenn*-clause is labelled as MOD because there is no explicit correlating constituent. It rather refers to the whole matrix clause, e.g. (*Wenn da was gebucht worden ist* (MOD), *ist das nicht in Ordnung*.)

6.2 Questions

6.2.1 W-Questions

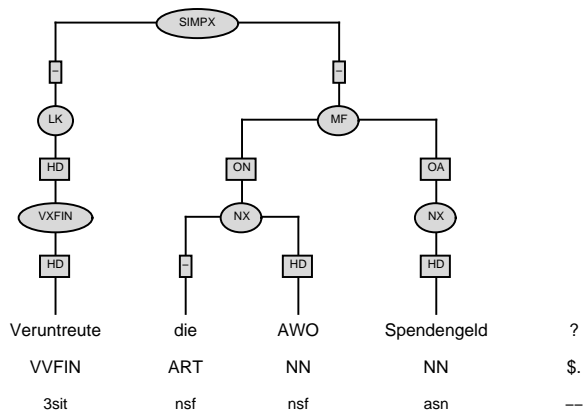
In general, w-questions are verb-second clauses with interrogative pronouns in VF. The problem here is to decide on the syntactic category of the interrogative phrase.

We follow the strategy to assign PX to all PWAVs, which compositionally comprise a preposition such as *wobei*, *wofür*, *wogegen*, *woher*, *womit*, *woran*, *worauf*, *wovon*, *wozu* and also to causal PWAVs such as *warum*, *wieso*, *weshalb*. The (non-compositional) PWAVs *wann*, *wo* are analysed as ADVX.

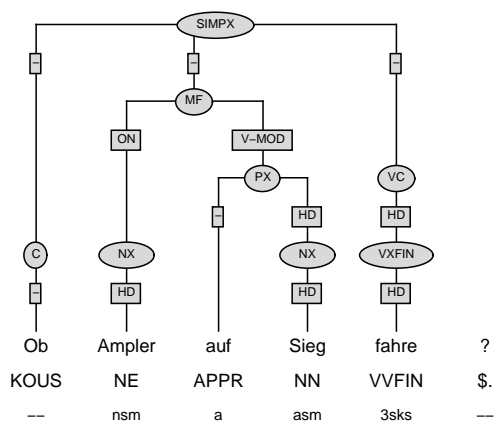
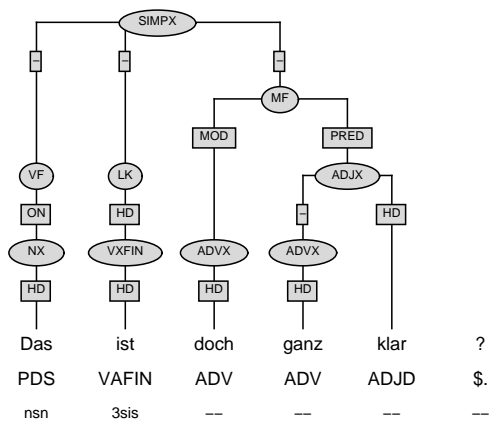


6.2.2 Yes - No Questions

Yes - no questions may occur in various forms, but the most typical form is the verb-first clause:

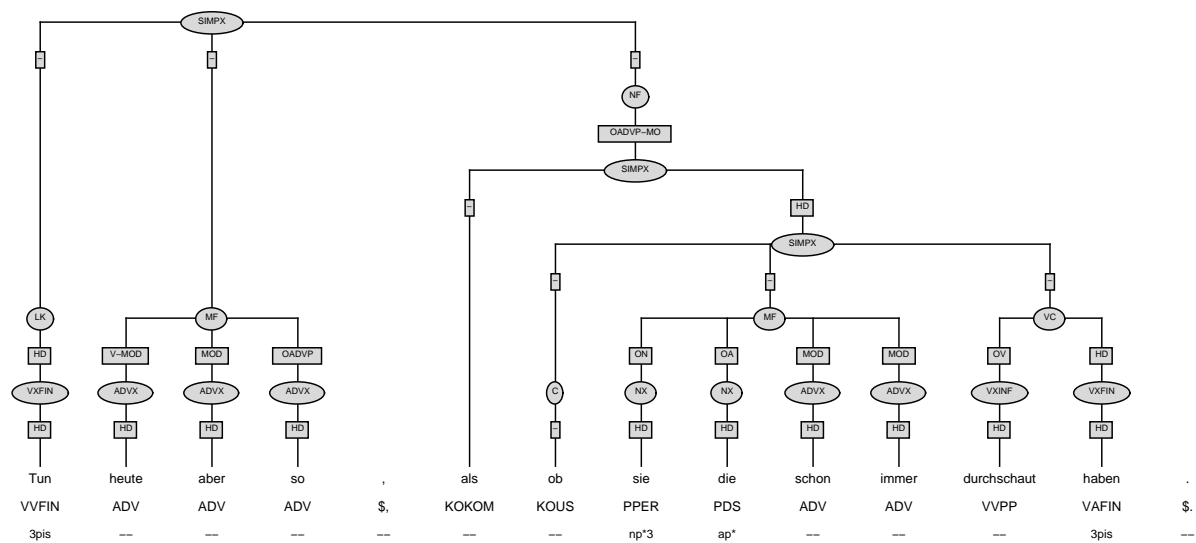
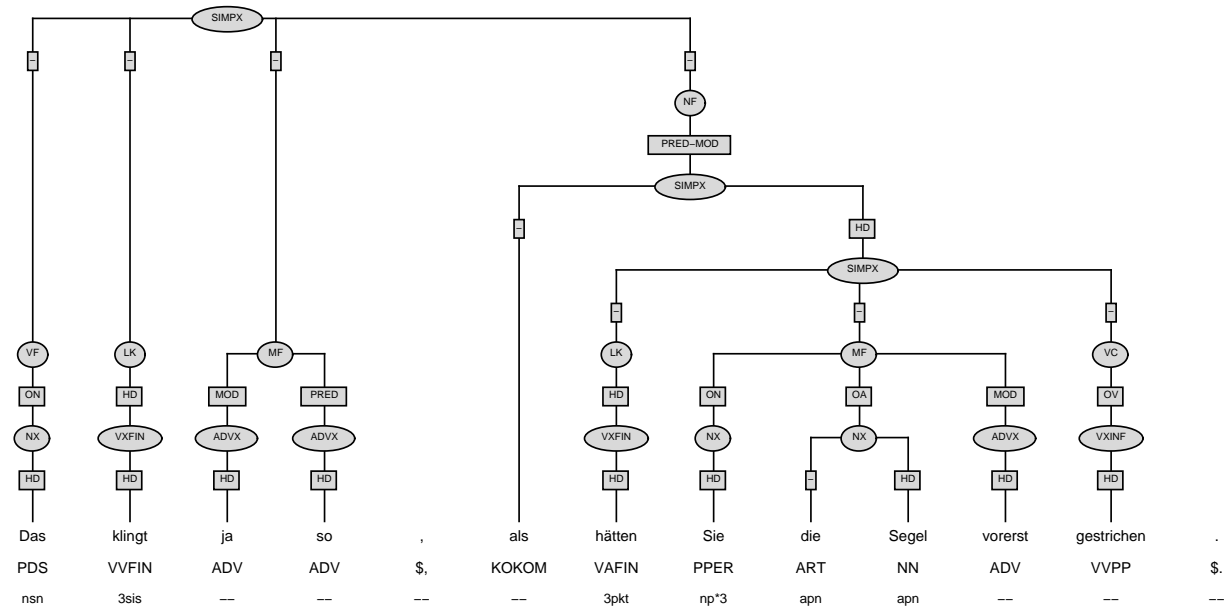


Otherwise, a question mark at the end of a verb-second or verb-final clause indicates that it is actually meant as a question:

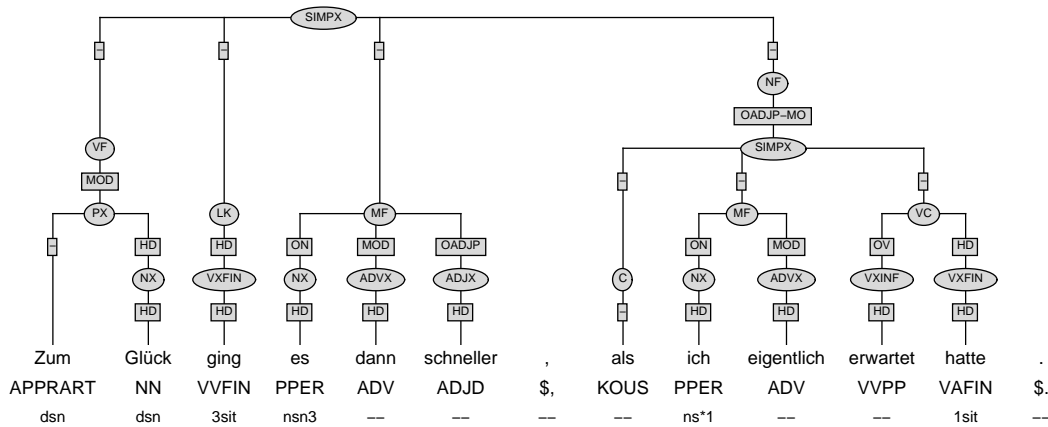


6.3 Clauses of Comparison

Clauses of comparison with *als* and *wie* which are semantically equated with a constituent (e.g. an adverb or an adjective) in the main clause, i.e. the comparison expresses an identity, are annotated with *als* and *wie* as particle of comparison (KOKOM). The sub clause can either be a verb-initial clause (e.g. ..., *als wäre* ...; ..., *als hätte* ...) or a verb-final clause (e.g. ..., *als ob* ...; ..., *wie wenn* ...; ..., *als daß* ...).

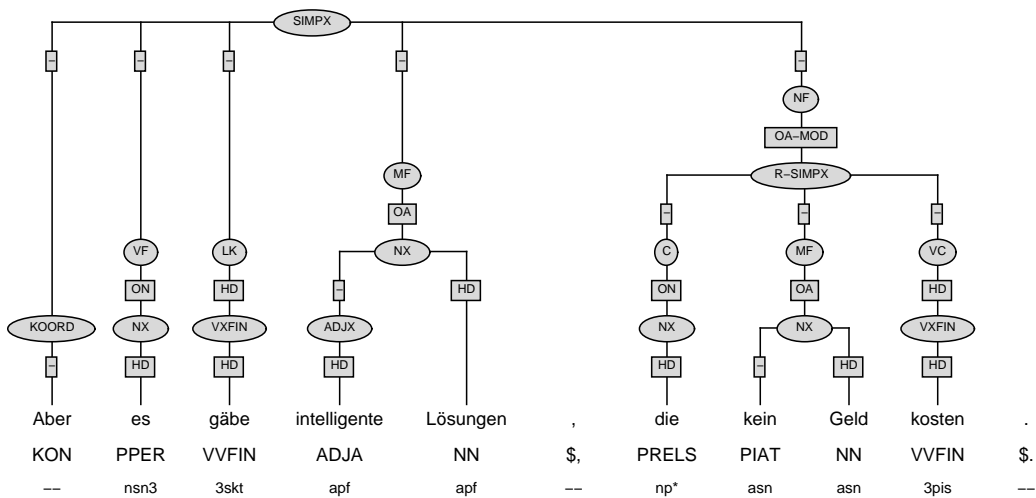


In contrast, clauses of comparison expressing a difference are annotated with the subordinating conjunction *als* (KOUS):

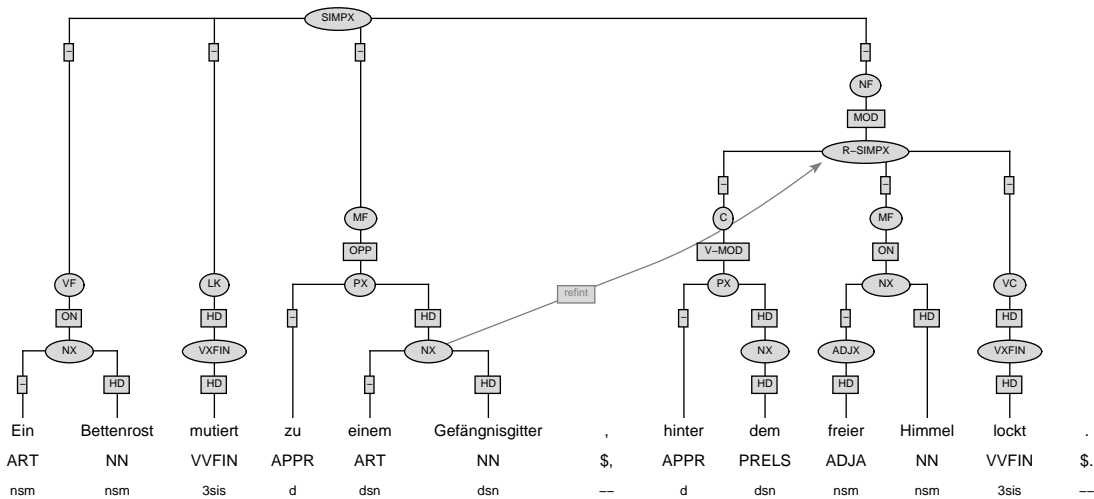


6.4 Relative Clauses

Considering relative clauses (R-SIMPX), the relative pronoun occurs in the C-field. It is first projected to the phrase level before it is attached to the C node. The relative clause itself is located in NF like in the following example if no other constituent follows. Its edge label shows to which constituent of the matrix clause it is related. OA-MOD, for example, suggests that the relative clause refers to OA:

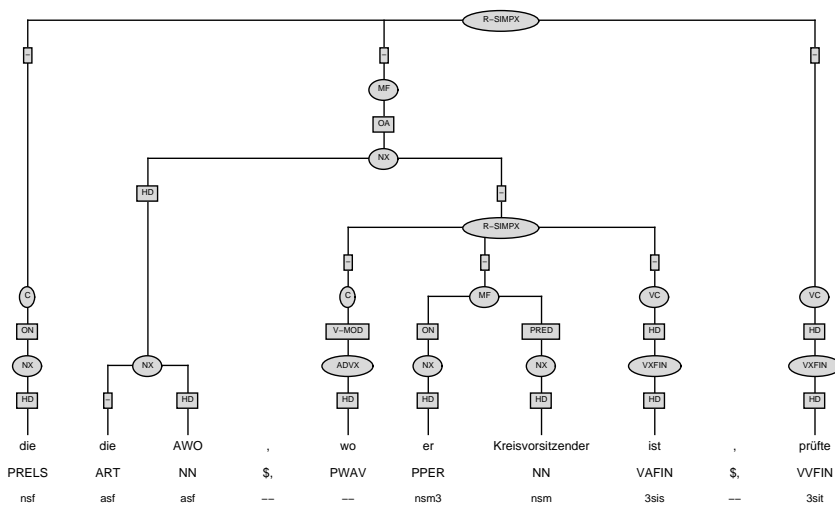


If the head noun phrase of the relative clause is the noun phrase of a prepositional phrase or a postmodifier within a complex phrase, the relative clause is labelled as MOD. Additionally, there is a secondary edge label named *refint* (cf. 3.4.6) from the head noun NX to the relative clause:



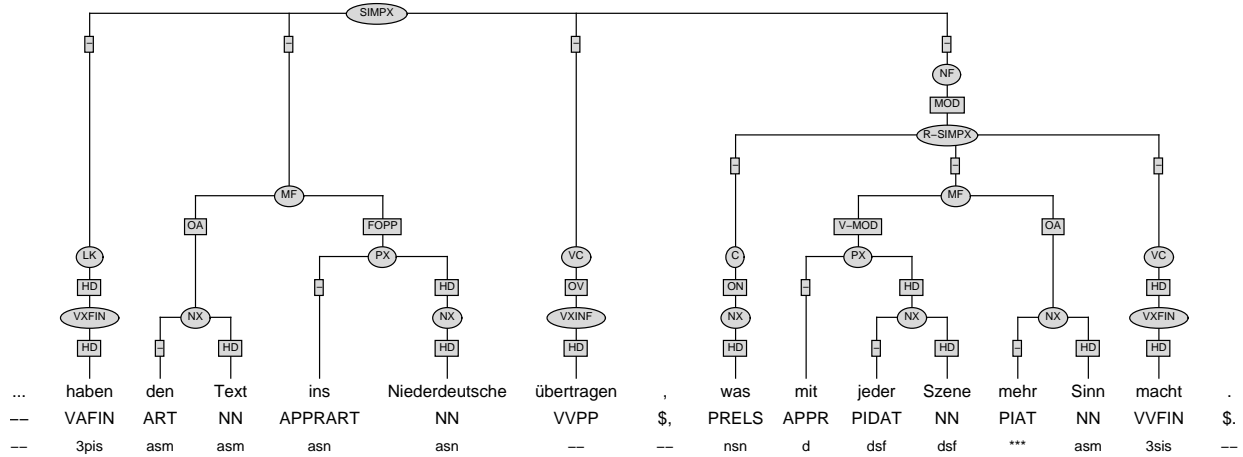
The position of the relative clause in NF is justified by the fact that it does not necessarily occur as an immediate constituent located to the right of the noun phrase to which it refers. For example, a verb complex can occur between the noun phrase and the relative clause (*Der Bettenrost ist zu einem Gefängnisgitter mutiert, hinter dem freier Himmel lockt.*). In sentences like this, the complexity of the noun phrase (NP + relative clause) is important. This so called *heavyness* follows Behaghel's first physical law (Behaghel 1932): complex noun phrases tend to find a position at the end of the sentence even if they deviate from their basic order. If the relative clause does not follow the noun phrase immediately, its unmarked position is in NF. Unless there is strong evidence for a position in MF, the relative clause is located in NF.

If the relative clause and its head noun phrase are adjacent constituents in VF or MF, the relative clause modifies the noun phrase directly as a postmodifier.



6.4.1 Event-modifying Relative Clauses

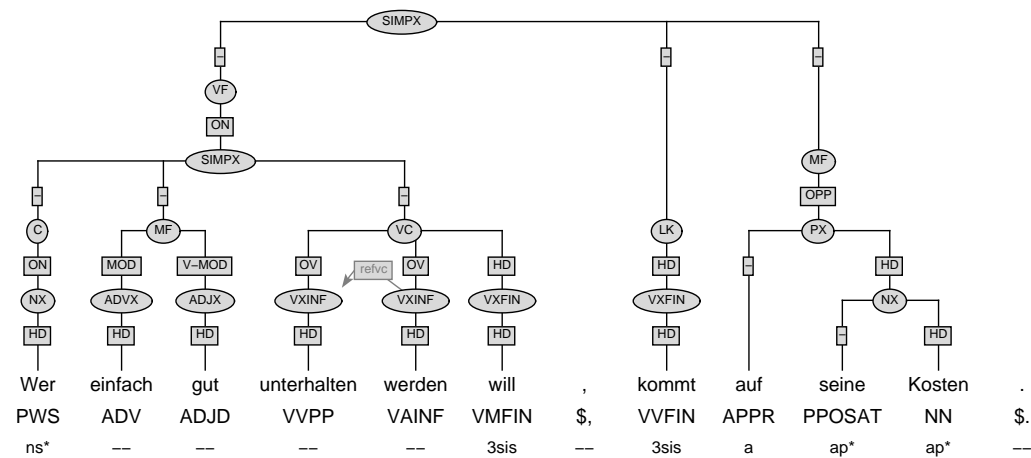
Relative clauses that modify an event which is not expressed by a nominal expression are also annotated as R-SIMPX.

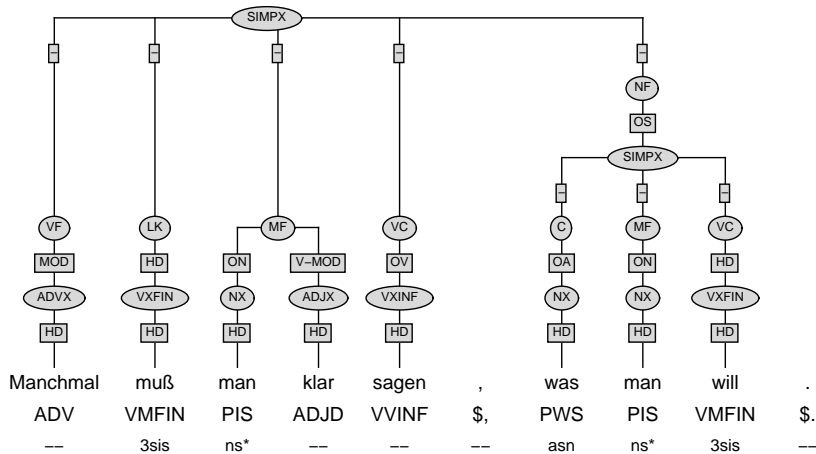


6.4.2 Independent Relative Clauses

Independent relative clauses (also 'nominal relative clauses', in German 'Freie Relativsätze') do not modify a head word but substitute an argument or adjunct in the clause. Consequently, they are labelled SIMPX on sentential level (instead of R-SIMPX) and they function as (sentential) subject (ON) or sentential object (OS). The latter is not uncontroversial since they are distributed like non-sentential, nominal arguments with respect to subcategorization restrictions.

The relative pronoun used in independent relative clauses normally belongs to the *w*-class of relative pronouns such as *wer* or *was* and is tagged with the STTS tag PWS.





Independent relative clauses introduced by *wie* are currently annotated in a different manner. *Wie* is analysed as subordinating conjunction (KOUS). This type of structure is to be revised in a subsequent release.

6.5 Coordination

Coordination is a syntactic phenomenon that occurs on the following annotation levels: phrase level, field level, and sentence level. Within coordinations, the conjuncts are first projected to their phrase, field, or clause level. In a second step, they are attached to their mother node which is n-ary branching (conjunctions between the conjuncts). This scheme is the same for all syntactic categories.

The edge labels between the mother node and the conjuncts of the coordination are labelled as KONJ. This edge label supports the distinction between conjuncts, modifiers, and conjunctions within complex conjunctions (cf. 6.5.3), as well as the distinction between coordinations and elliptical constructions (cf. 6.6).

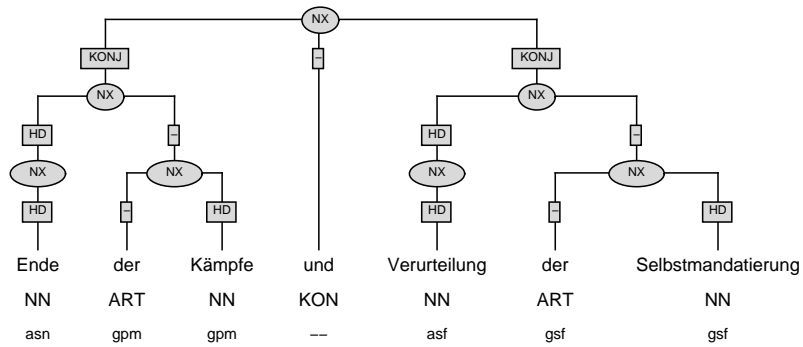
In contrast to coordinating conjunctions in the KOORD-field, coordinating conjunctions in coordinations (*und*, *oder*, etc.) are directly attached to the mother node of the conjuncts. The class of coordinating conjunctions consists of single, e.g. *und*, *oder*, *aber*, *als*, as well as of complex conjunctions, e.g. *entweder oder*, *weder noch*, *sowohl als*. Generally, coordinating conjunctions may coordinate constituents of any category. Moreover, they can form asymmetric coordinations in which the conjuncts belong to different syntactic categories (cf. 6.5.2).² In order to distinguish conjunctions from conjuncts within a coordination, their edge labels are empty.

In the following, coordination on all annotation levels as well as specific cases of coordination, e.g. split coordinations, will be demonstrated.

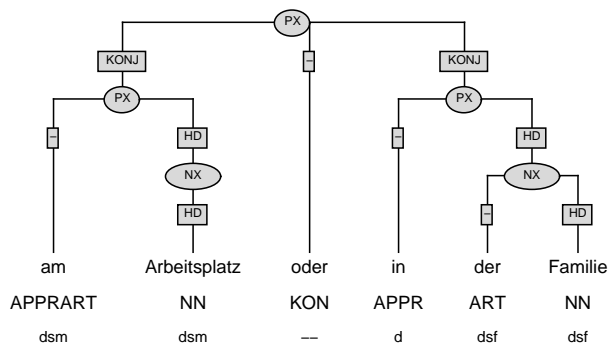
²If *bis* is used as a conjunction like in *10.000 bis (KON) 20.000 koreanischen Daewoo PKW* it is tagged as KON. But remember that *von ... bis ...* phrases are treated differently (cf. 4.4.1).

6.5.1 Coordination of Phrases

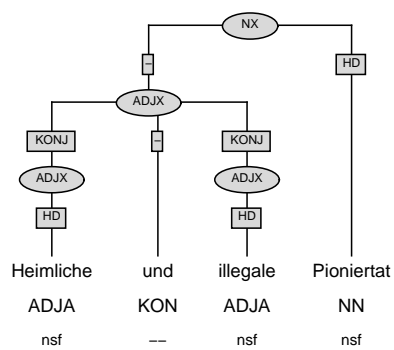
Noun Phrases

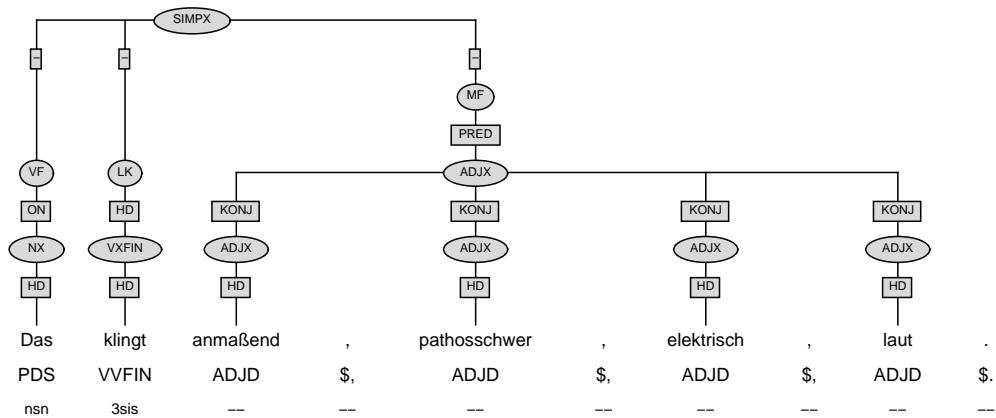


Prepositional Phrases

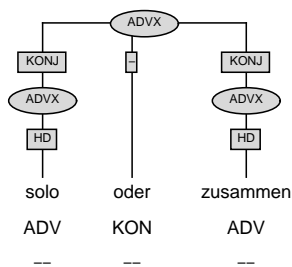


Adjectival Phrases



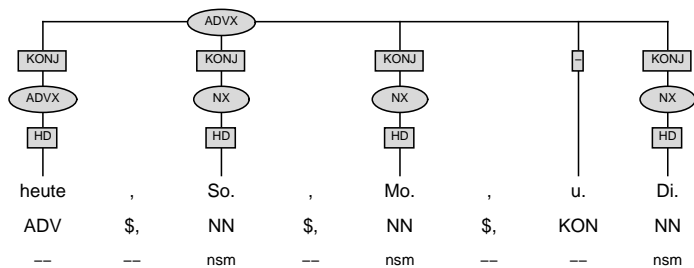


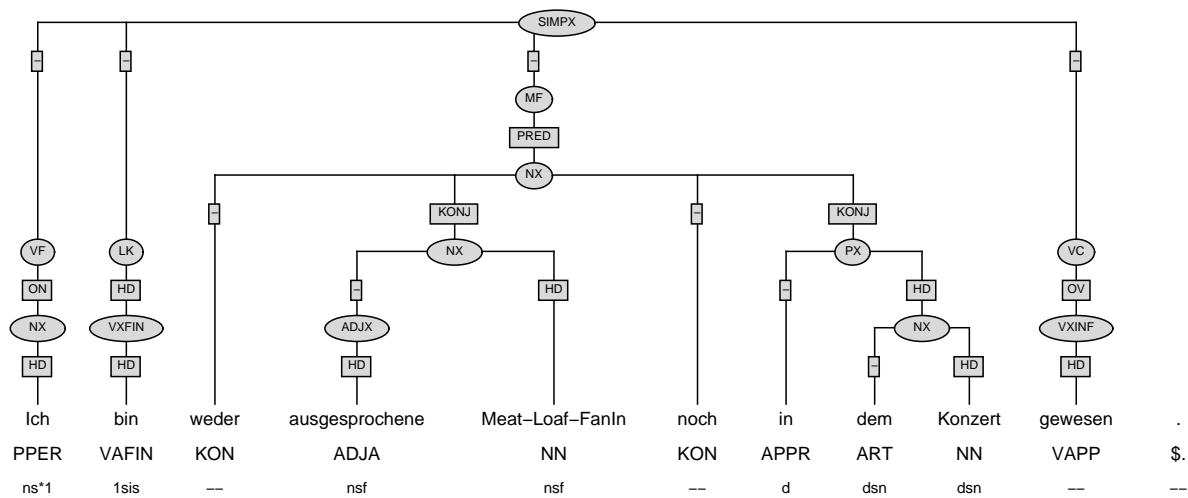
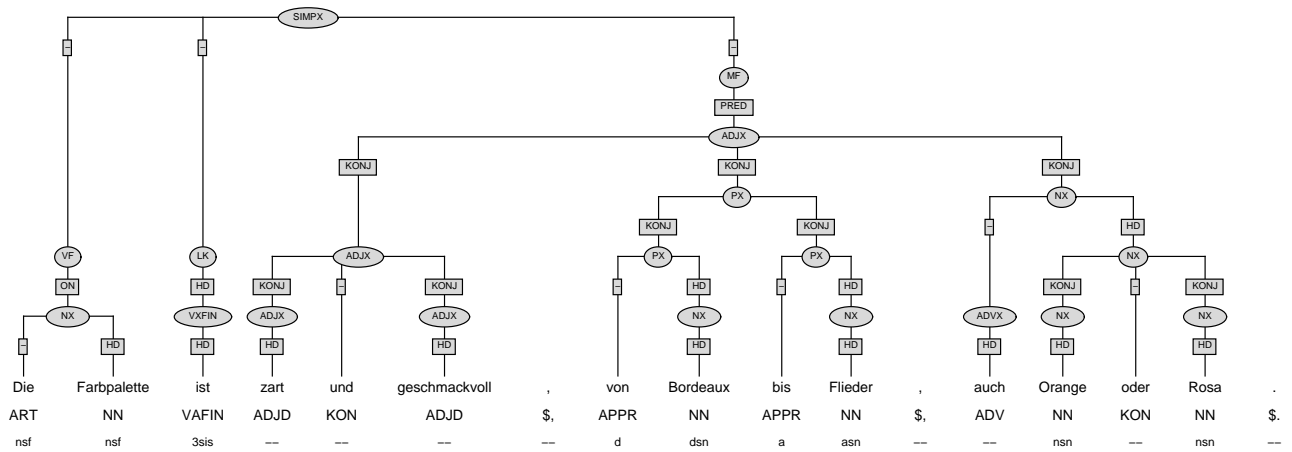
Adverbial Phrases



6.5.2 Asymmetric Coordination

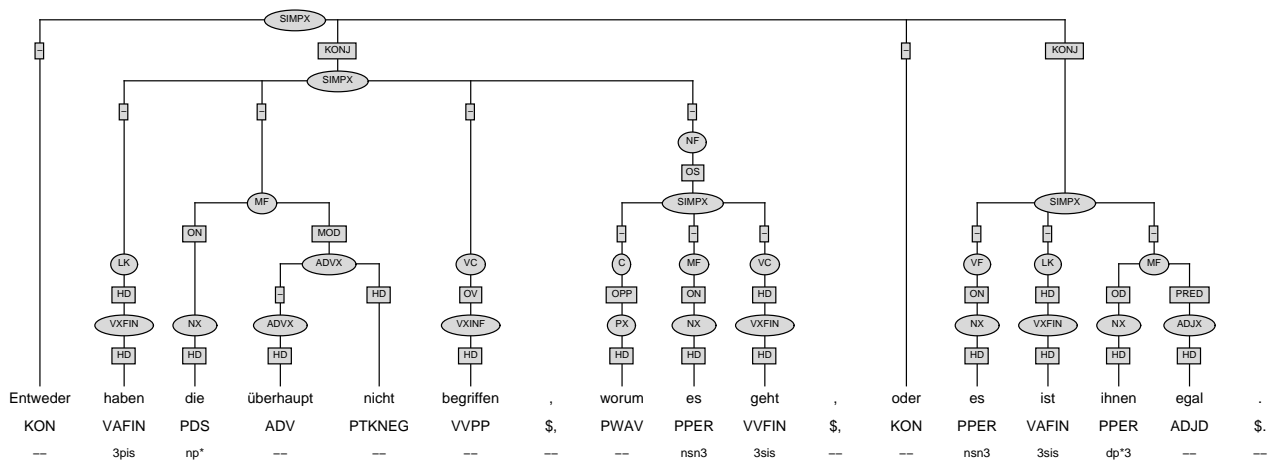
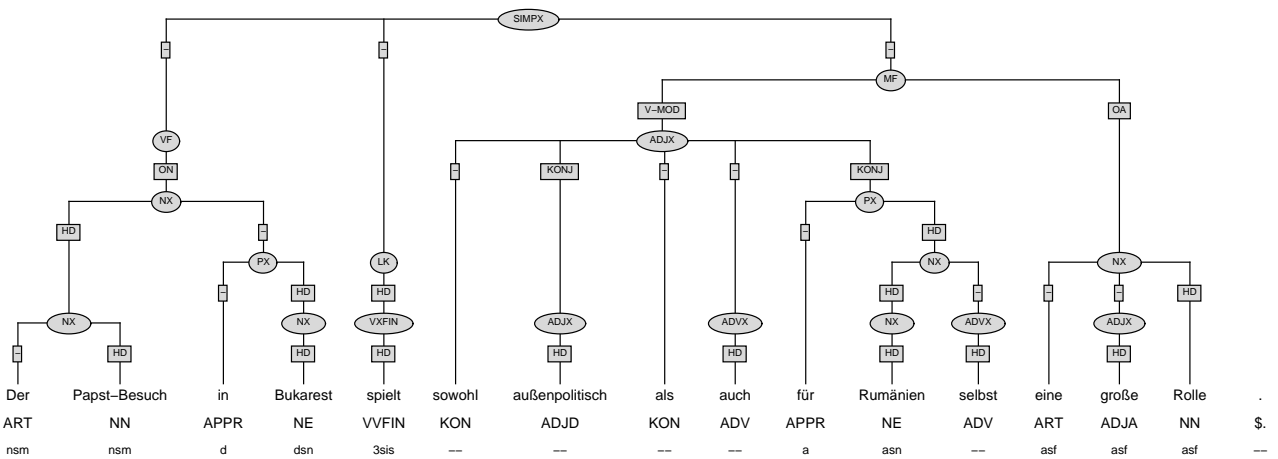
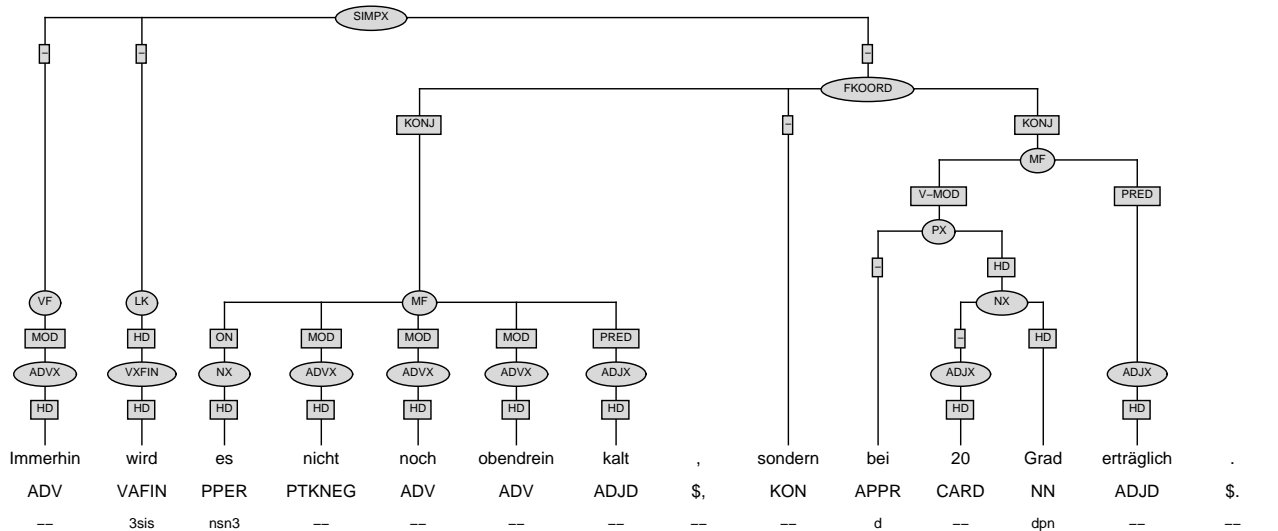
Since constituents of different syntactic categories can be coordinated, it has to be decided on a label for the mother node of the coordination. In this case, the default strategy has been adopted to choose the syntactic category of the left-most conjunct as the category of the entire coordination:





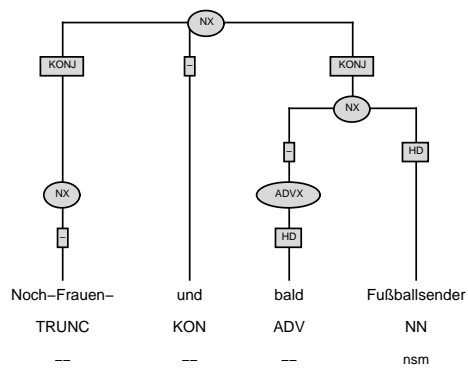
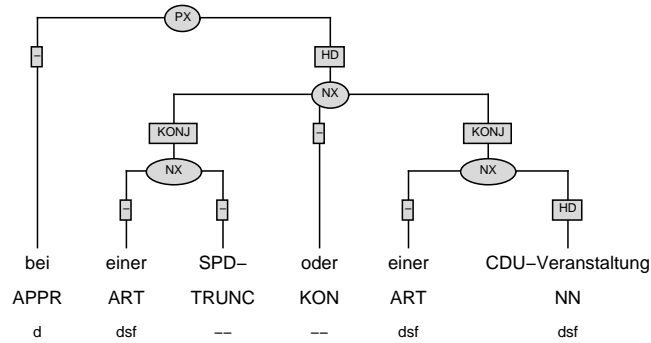
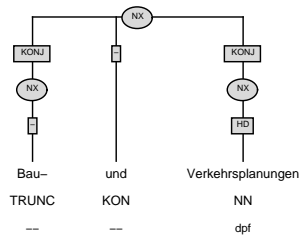
6.5.3 Coordinations with Complex Conjunctions

The conjuncts and conjunctions of a coordination with complex conjunctions are also attached on the same level following the above mentioned rules for coordination. Both parts of complex conjunctions like *entweder oder* and *sowohl als* are tagged as KON. The latter one usually occurs together with the adverb *auch*, which is tagged as ADV, projected to the phrase level, and then attached to the mother node of the coordination. The same applies for *nicht* in coordinations with *sondern*. *Sondern* is tagged as KON, whereas *nicht* is always tagged as PTKNEG:

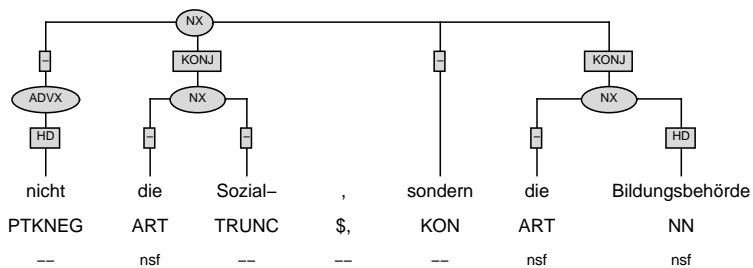
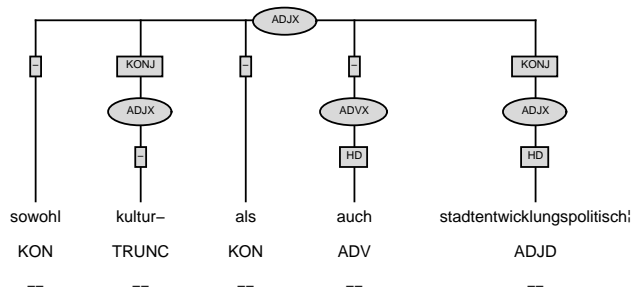


6.5.4 Coordinations with Truncated Words

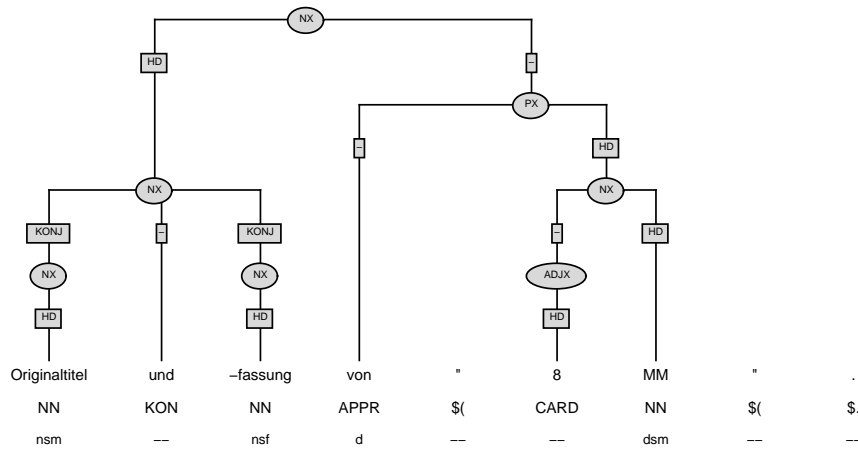
Truncated words are projected to the phrase level. Their edge labels are empty. The phrases of both conjuncts are coordinated. The truncated words do not receive morphological annotation.



In the case of complex conjunctions, the conjuncts are annotated in the same way.

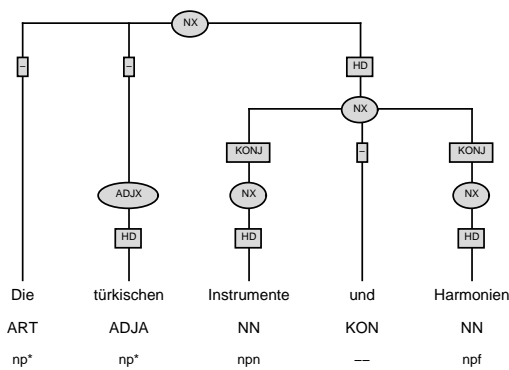
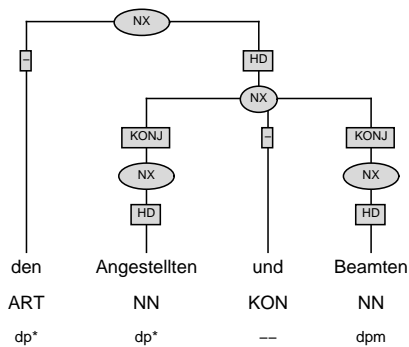


Word initial TRUNCs are different from truncated words which include the second part of a word. The latter ones are treated like complete lexical heads, because they comprise the head morpheme of the complex word.



6.5.5 Attachment Principles of Coordination within Phrases

If two or more nominal conjuncts occur together with a common determiner and/or adjectival phrase, first the conjuncts are projected to their phrase level and then the determiner or the adjectival phrase is attached to the coordination on a higher level according to the *high attachment principle*. Thus, the modification scope comprises the entire coordination. The coordinated part is assigned the head function.

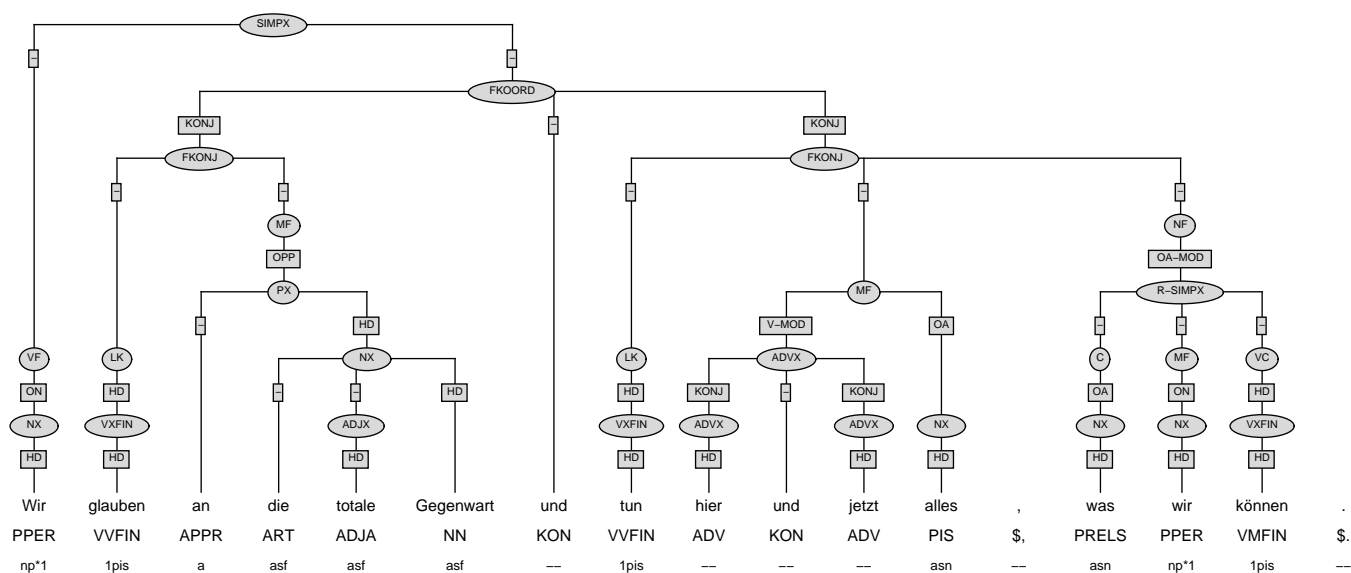


6.5.6 Coordination of Topological Fields

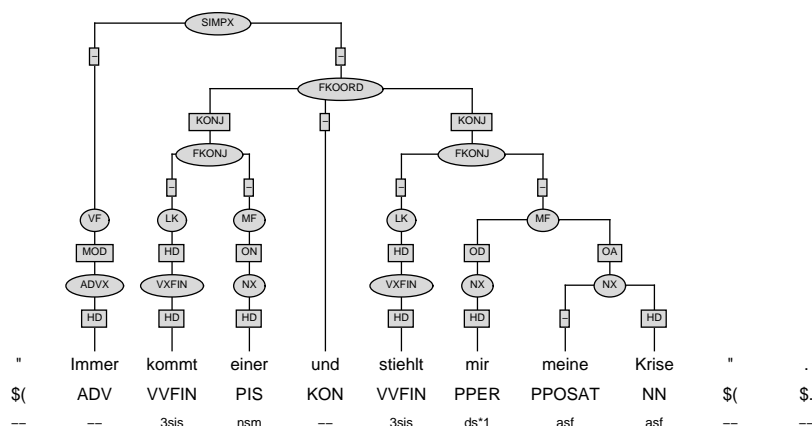
The conjuncts of a coordination of topological fields are either single fields (cf. 6.5.4) or a combination of fields. Possible combinations are, for instance, (MF + VC), (LK + MF), (LK + MF + VC). The node label for these conjuncts is FKONJ (conjunct consisting of fields) and the mother node of a coordination of conjuncts of fields is FKOORD.

In a coordination of conjuncts of fields, the following annotation steps are involved:

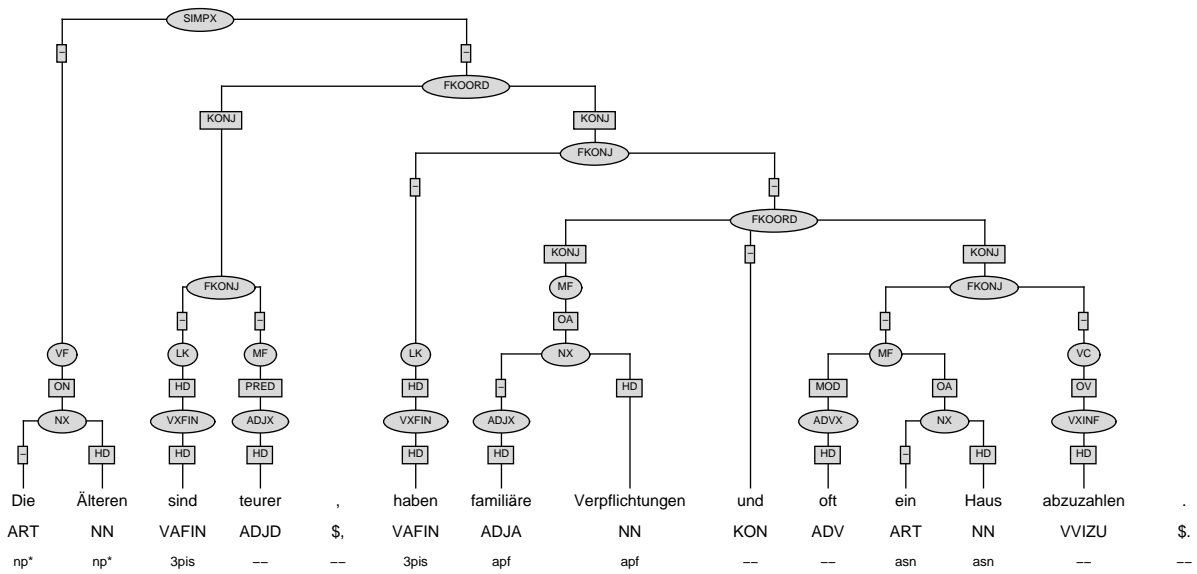
1. The constituents are attached to the fields in which they occur in (MF, VC, NF, etc.).
2. Each conjunct (concatenation of fields or single field) is labelled as FKONJ.
3. The conjuncts are attached to the general coordination field FKOORD.



Oftentimes, the subject of the sentence occurs only in the left field conjunct:

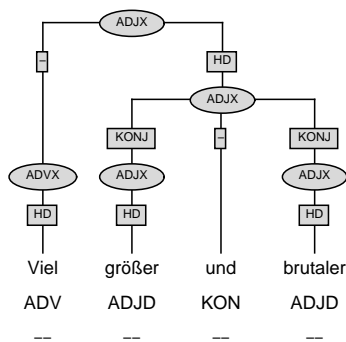


A coordination of fields may also be an embedded structure. In this case, FKOORD functions also as conjunct label:

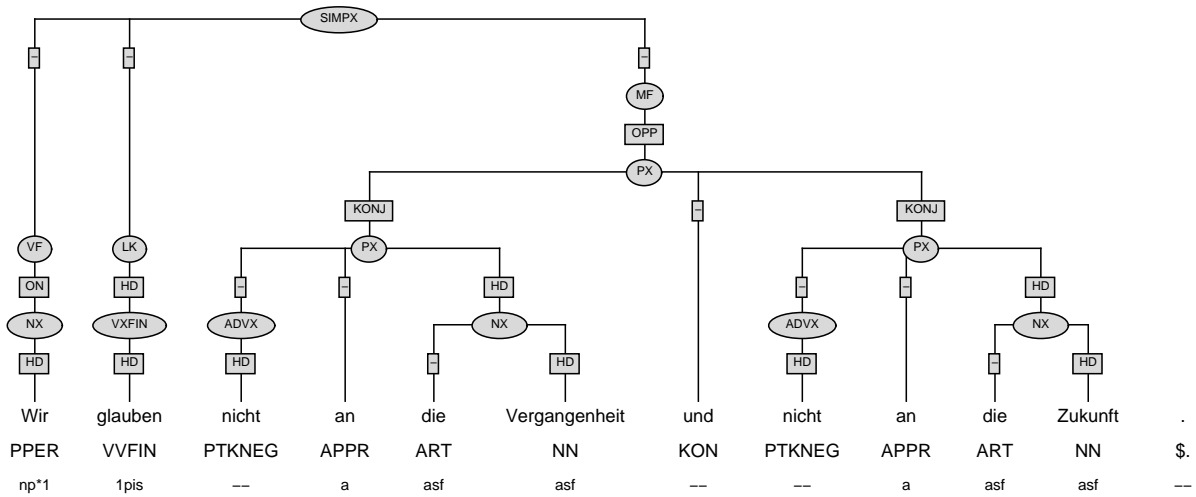


6.5.7 Attachment of Ambiguous Modifiers in Coordination

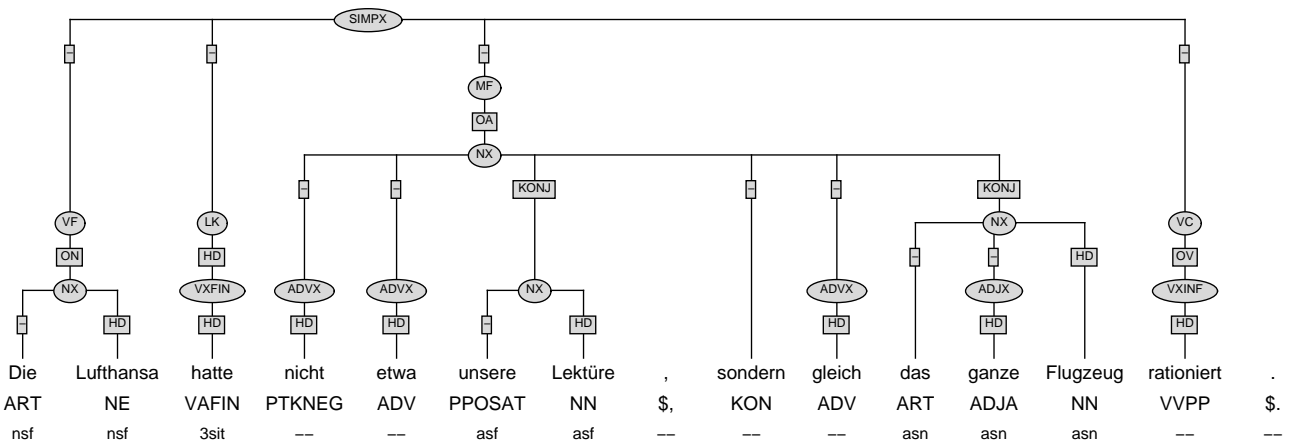
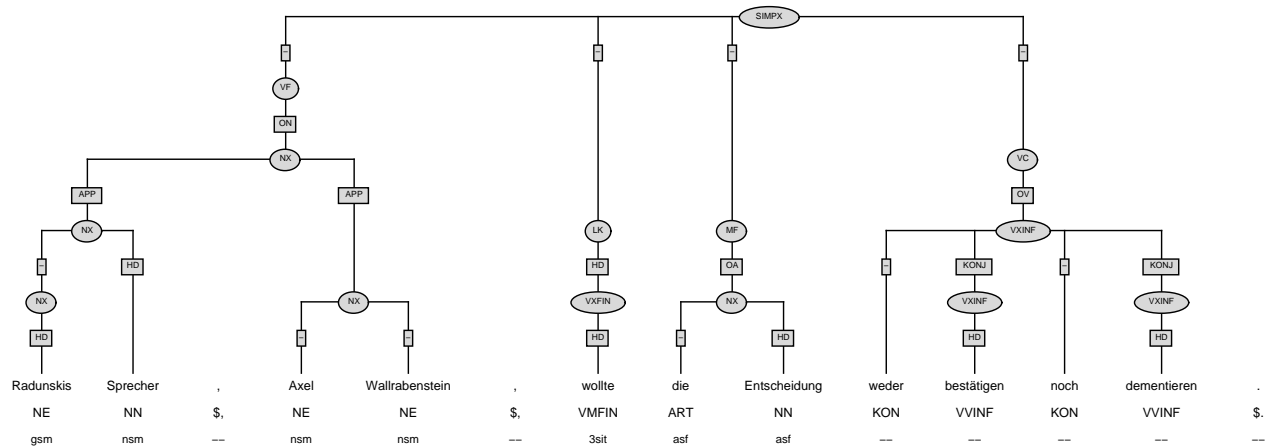
Within phrases, the modification scope of a premodifier can be ambiguous. Therefore, high attachment is applied to preserve ambiguity. In the following example, the adverb modifies the coordination of adjectives rather than only the first adjective:



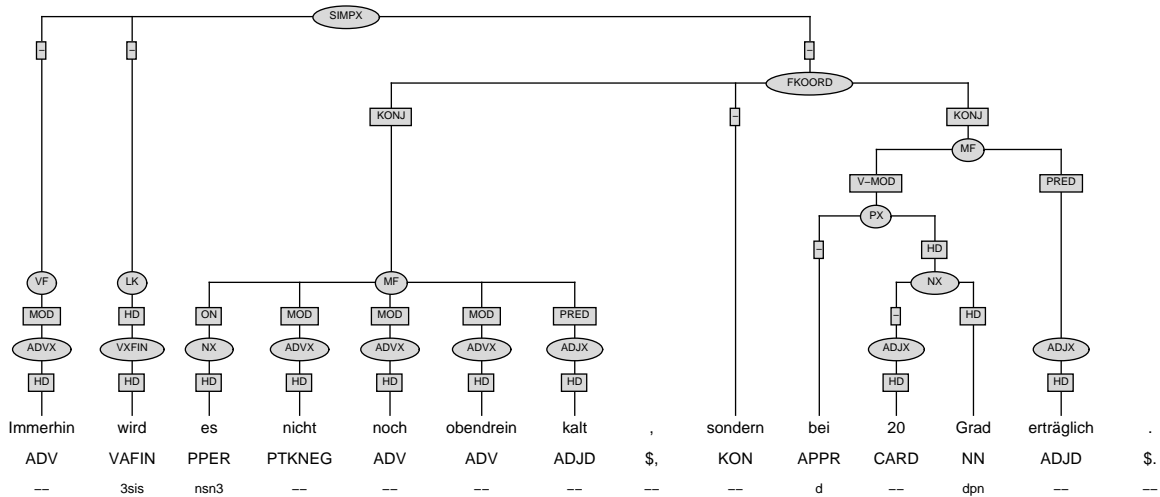
Modifying constituents are attached to a conjunct rather than to a field if their modification scope is limited to the conjunct.



Also in coordinations with complex conjunctions, attachment on the phrase level is applied if possible.

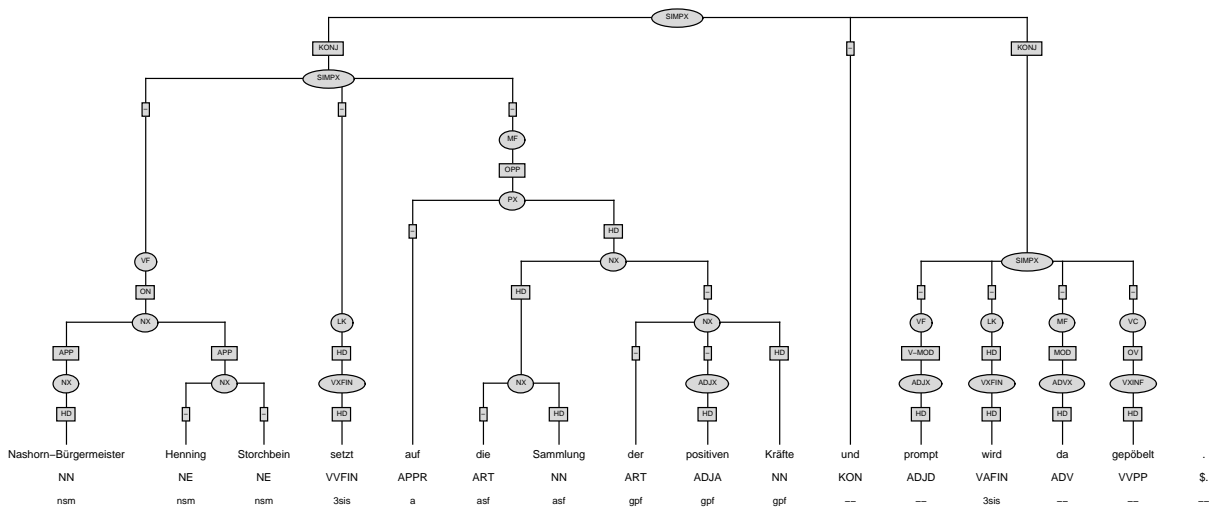


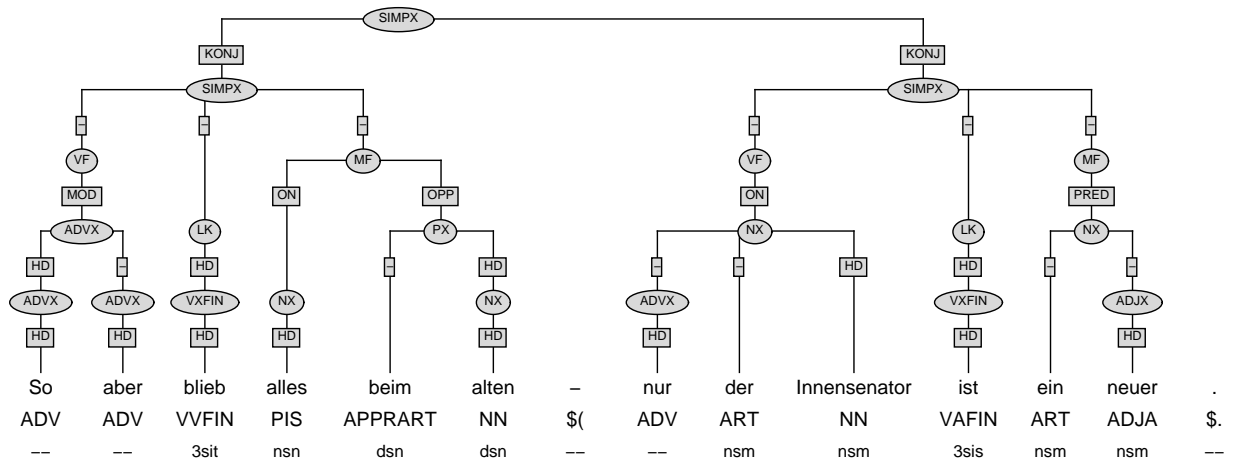
If there is more than one constituent within a conjunct, each with its own grammatical function, these constituents are first attached to the respective field node. Then, the fields are coordinated:



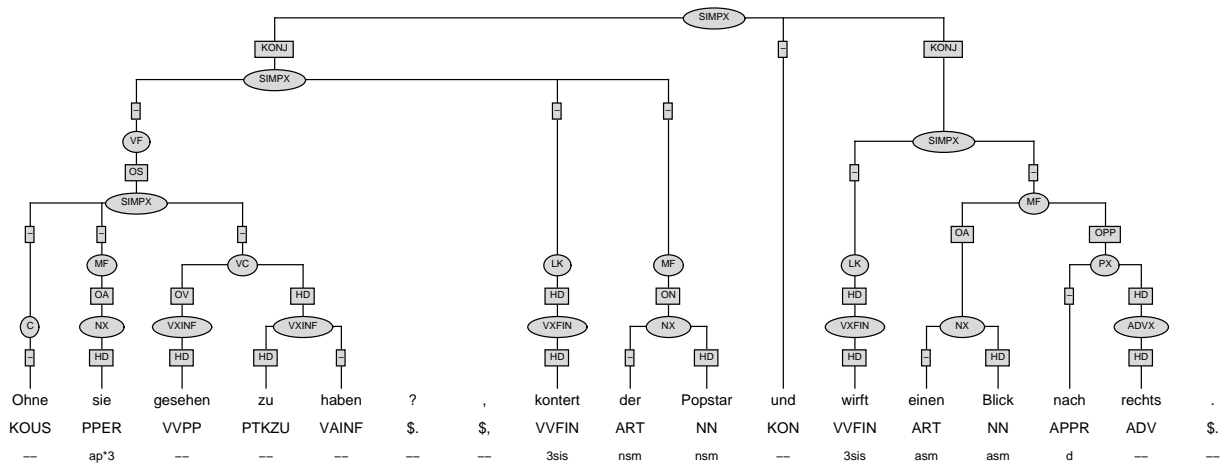
6.5.8 Coordination of Sentences

In accordance with the *longest match principle*, complete sentences are coordinated as paratactic constructions when they belong to the same syntactic unit (cf. 3.4.3), i.e., they are coordinated by a conjunction, a comma, or a dash:

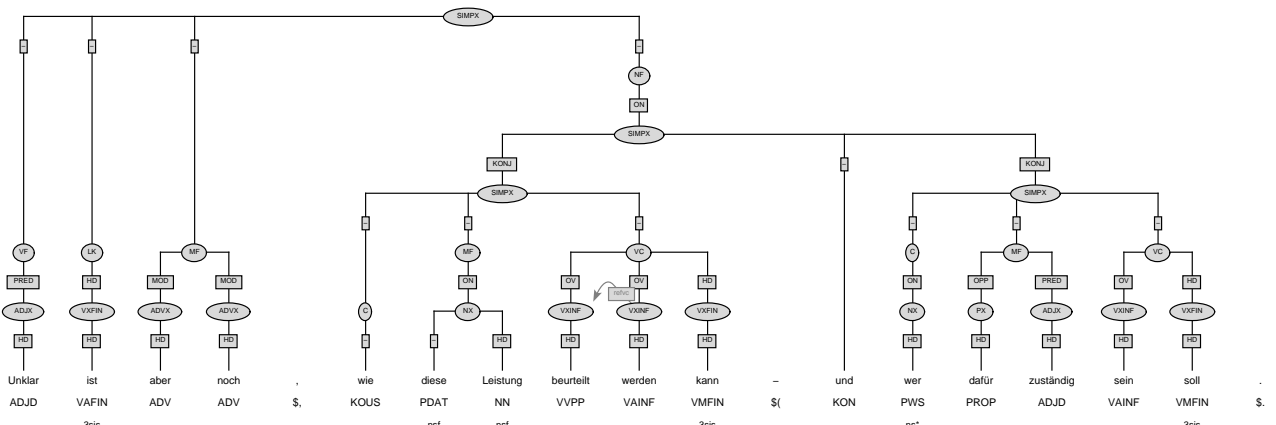




A coordination may also consist of two sentences with the subject of the whole construction only occurring in the left conjunct of the coordination.

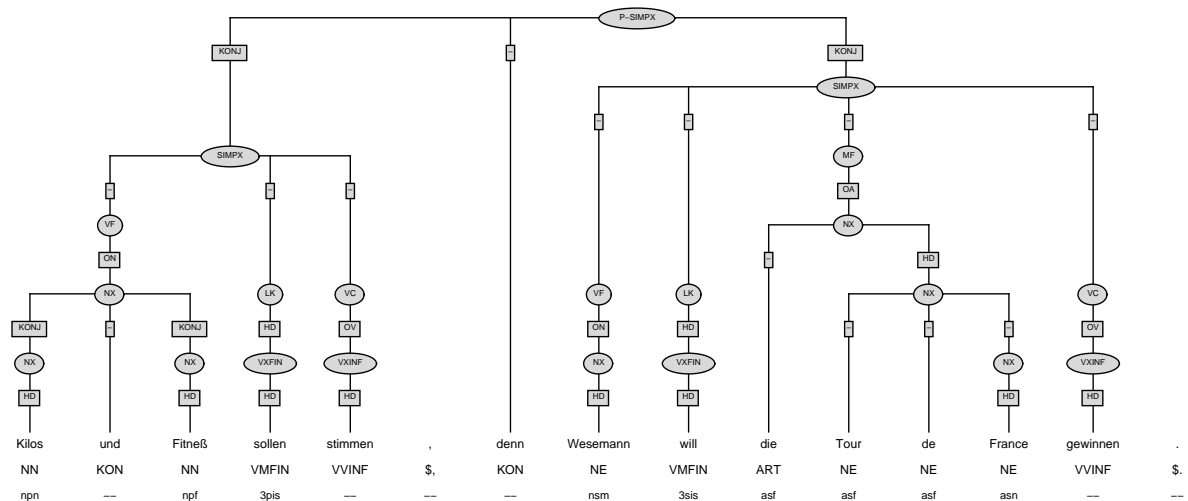


Subclauses (either in VF or in NF) with or even without a conjunction can also be coordinated.

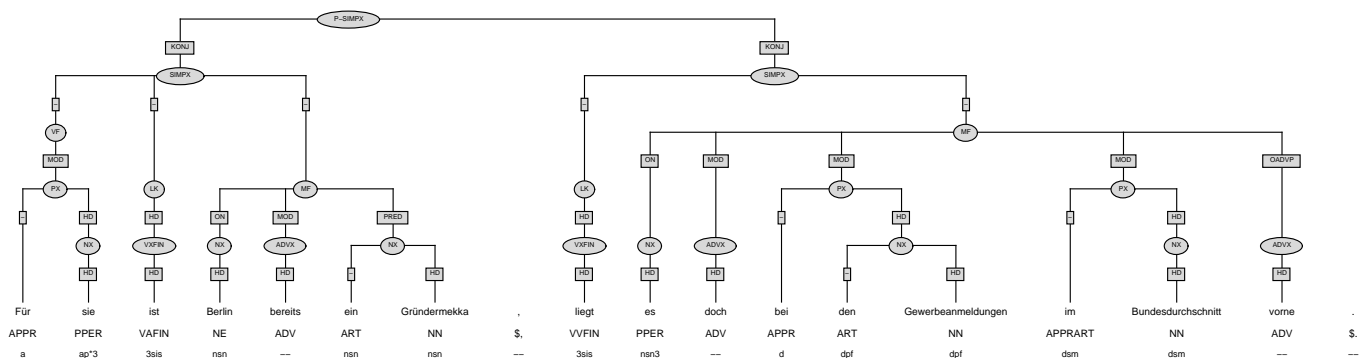


6.5.9 Paratactic Constructions

Paratactic constructions consisting of verb-second clauses conjoined by the conjunctions *denn* and *weil*, which also occur in the PARORD-field in the beginning of a sentence, are treated as syntactically equivalent conjuncts (verb-second instead of verb-final in *weil*-clause). In order to distinguish coordination of sentences with a conjunct of the PARORD field from the above mentioned coordinations of sentences, these paratactic constructions are labelled as P-SIMPX instead of SIMPX.

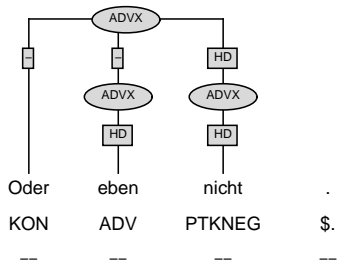
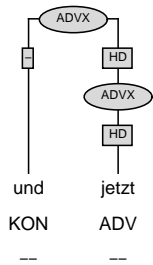


Syntactically coordinated but semantically subordinated main clauses with adversative conjunctive adverbs (e.g. *doch*) are also annotated as paratactic constructions:

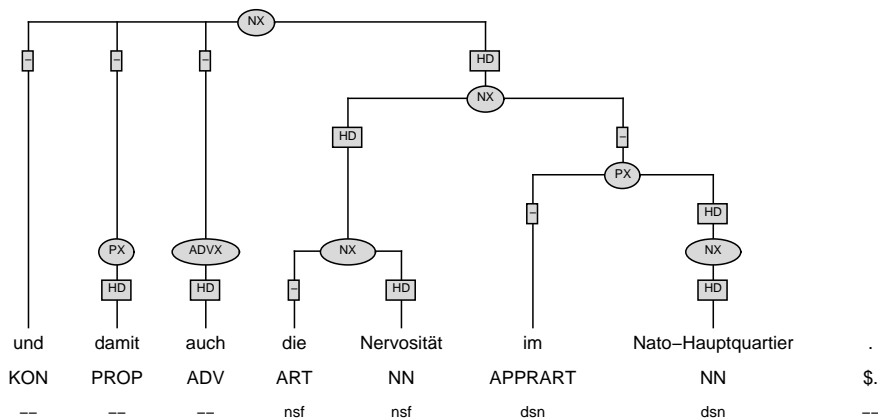
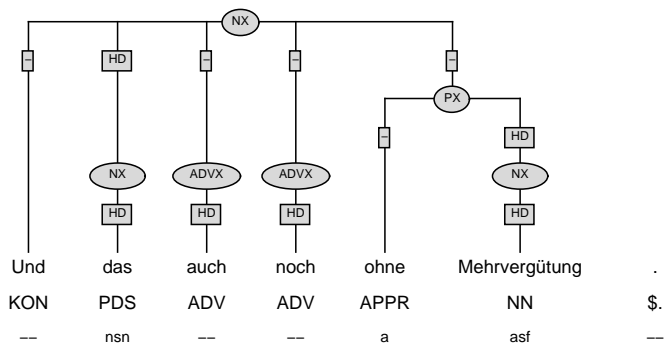


6.5.10 Conjunctions Occurring with Isolated Phrases

If a conjunct occurs isolated with a conjunction, high attachment is applied like in complete coordinations. But for isolated conjuncts, the conjunct is annotated as the head of the construction (HD instead of KONJ).

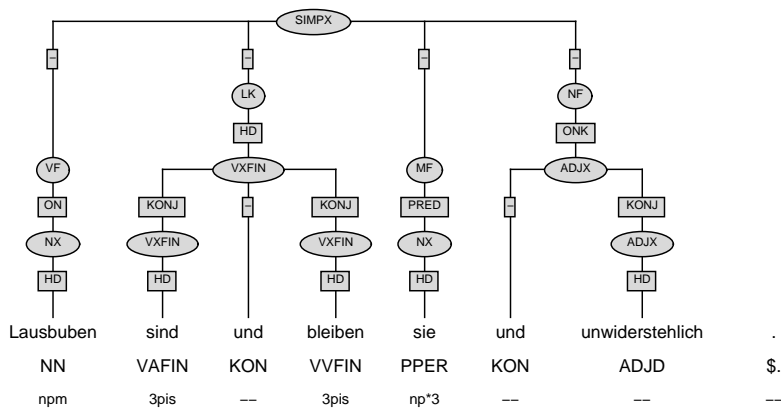
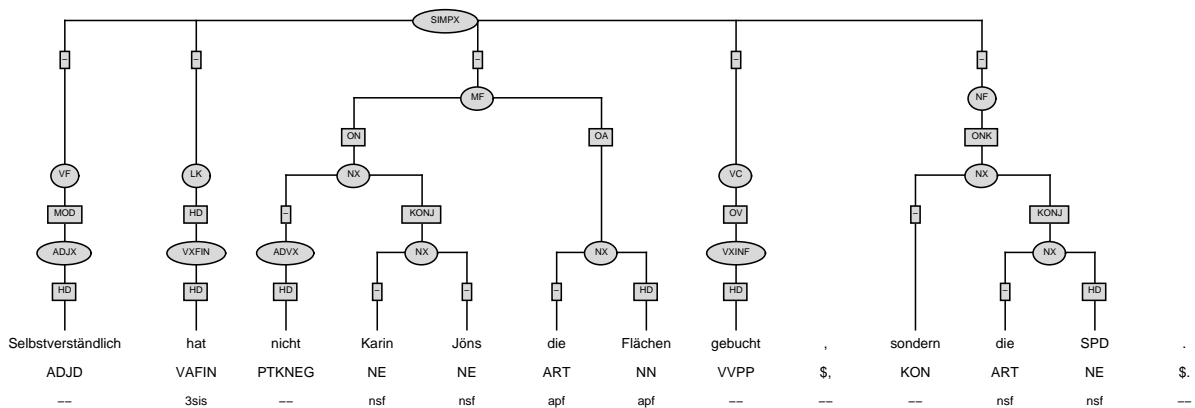
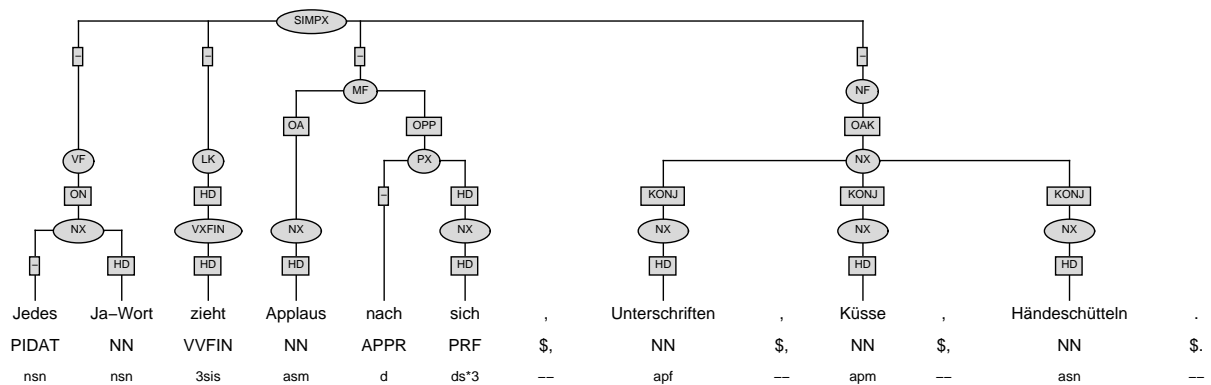


If there are modifiers which do not modify the conjunct itself because they are ambiguous or might modify something else rather than the conjunct, they are attached on the same (high) level as the conjunction:



6.5.11 Split Coordinations

Closely related to isolated conjuncts are *split* coordinations. Generally, the left conjunct of a split coordination is located in MF, in rare cases in VF, and the right conjunct occurs in NF. In order to express the relation between them, the left conjunct carries the label of its grammatical function (ON, OA, OD, etc.) whereas the right conjunct carries a label that denotes that it is the conjunct of this grammatical function (e.g. ONK, OAK, ODK, etc.). In asymmetric coordination, the syntactic category of the second split conjunct determines the syntactic category one level higher up:

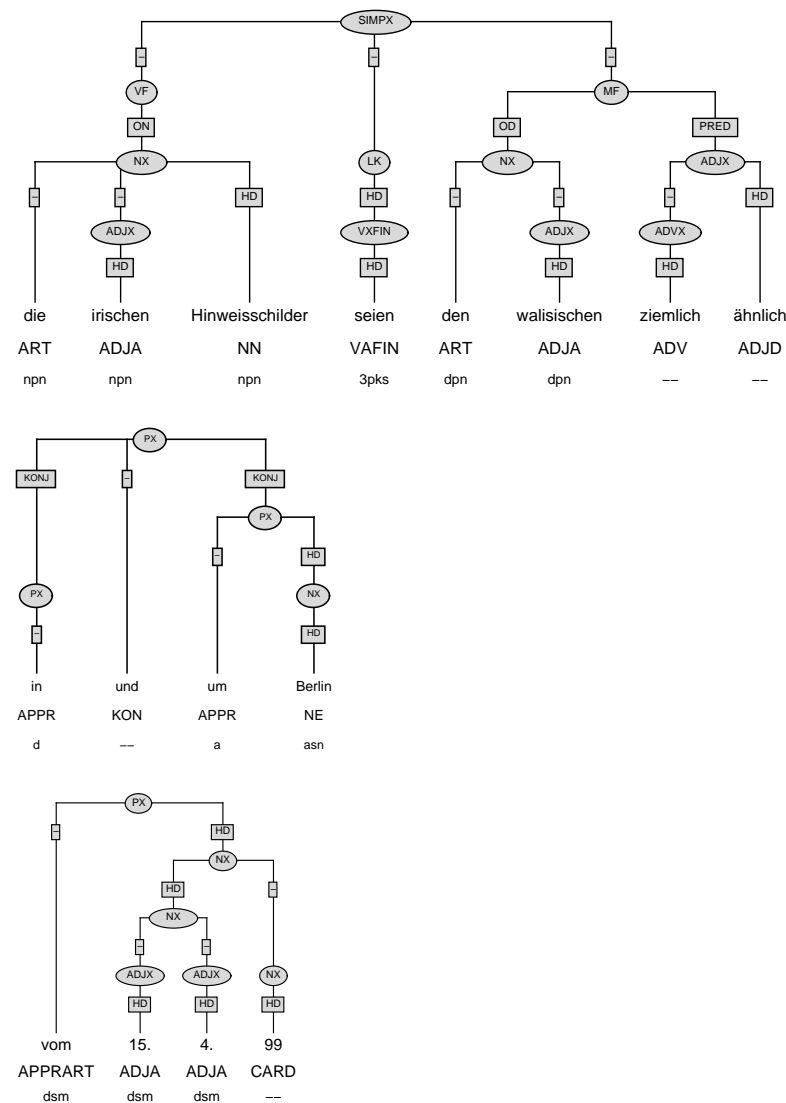


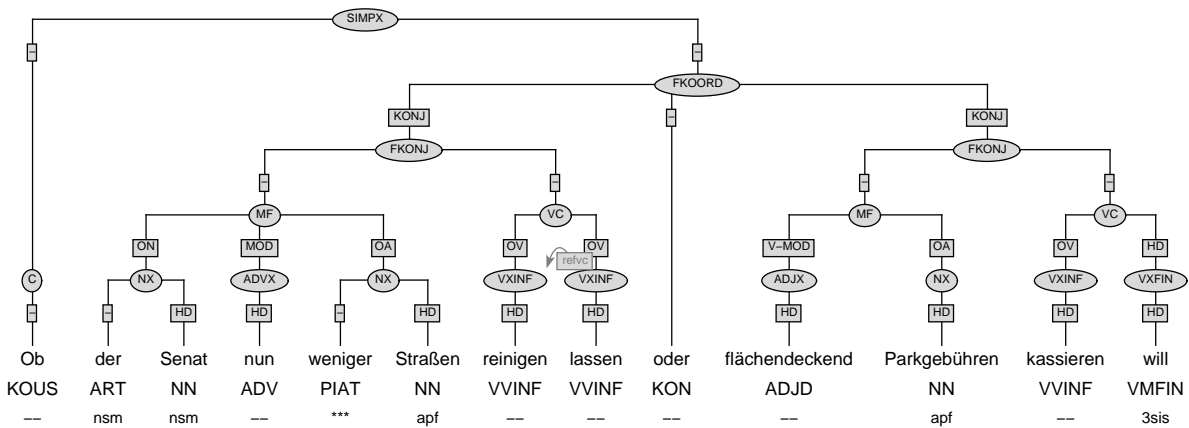
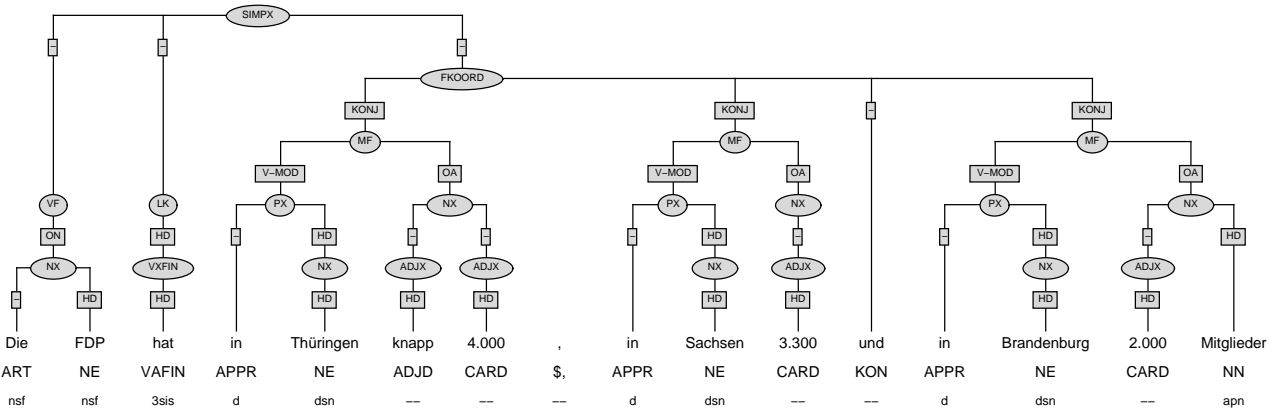
6.6 Elliptical Constructions

In elliptical constructions, syntactically necessary linguistic elements are missing which can be reconstructed from the context or the speech situation. Elliptical constructions appear on the phrase level as well as on the sentence level.

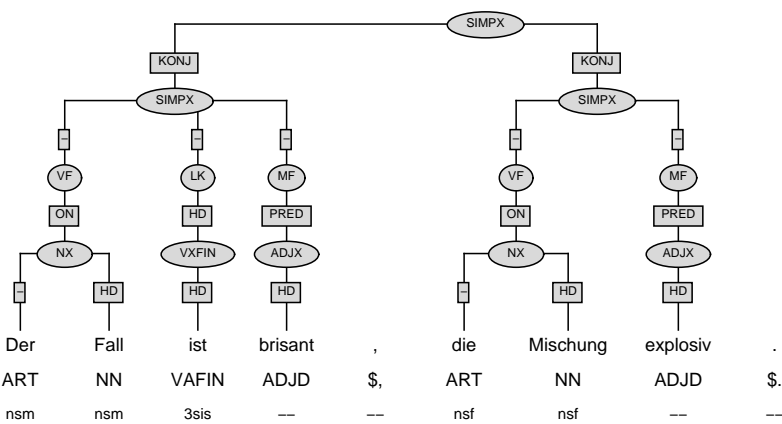
The model of topological fields does not make any assumptions about dependency relations, but it allows that topological fields may be left empty. For the description of elliptical sentence constructions, the scheme of topological fields is an appropriate model because neither crossing branches nor traces have to be used to annotate the surface structure of a sentence.

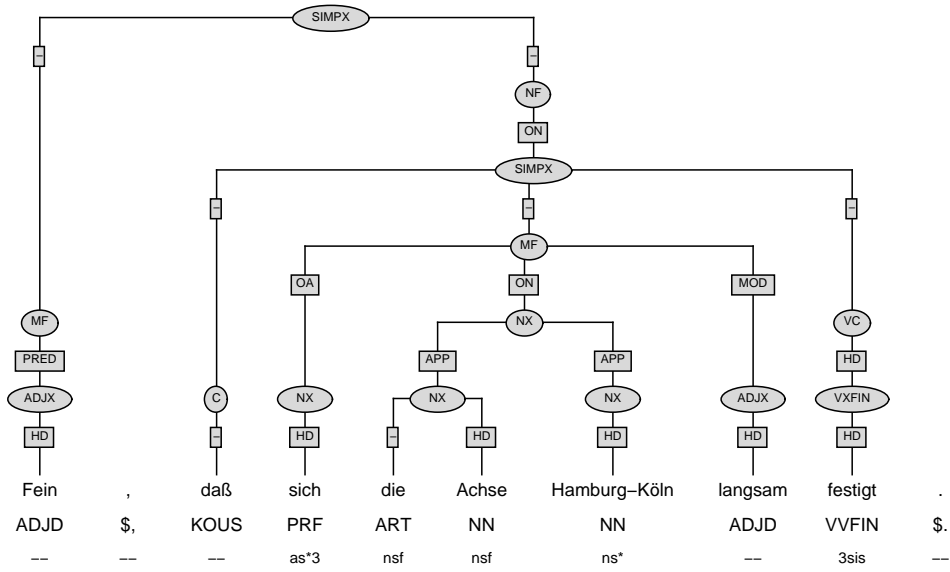
In elliptical phrases, the head word is missing. They are annotated like phrases without a head. Therefore, the edge labels of an elliptical phrase are empty:





In elliptical sentence constructions, specific topological fields are not occupied. All constituents are attached to the appropriate field. In the first example, LK in the second conjunct is missing. In the second example, the subject is in NF and the main clause is lacking a verbal constituent:





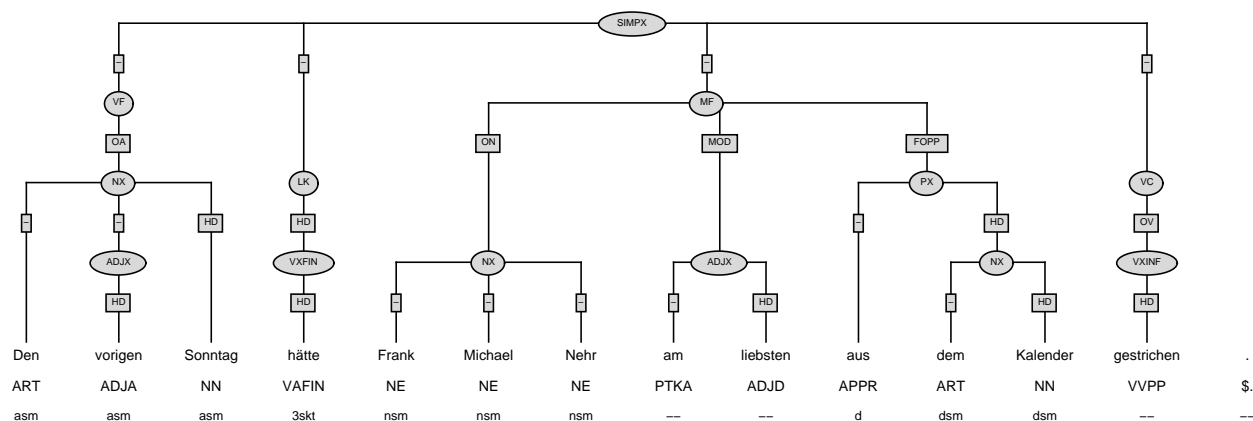
Chapter 7

The Annotation of Specific Syntactic Phenomena

7.1 Superlative and Comparative Forms

7.1.1 Superlative Forms

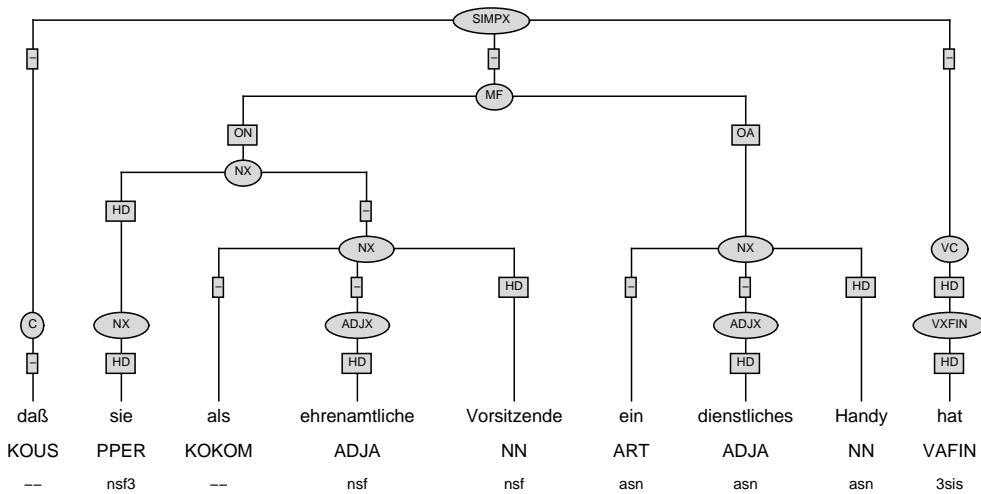
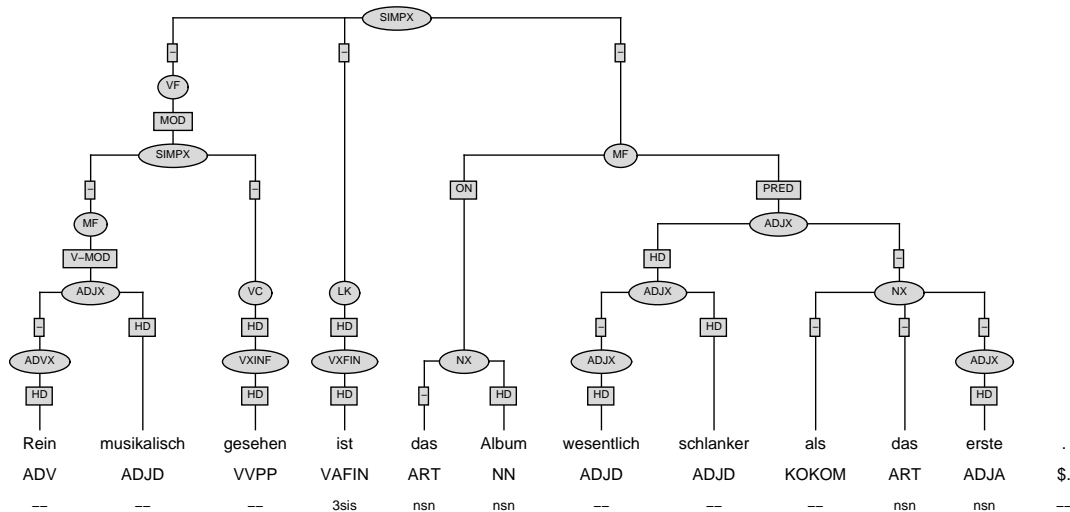
The particle *am*, which occurs as a particle with an adjective or an adverb in superlative constructions, is tagged as PTKA. Both, the particle and the adjective/adverb are attached on the same level forming an adverbial/adjectival phrase:



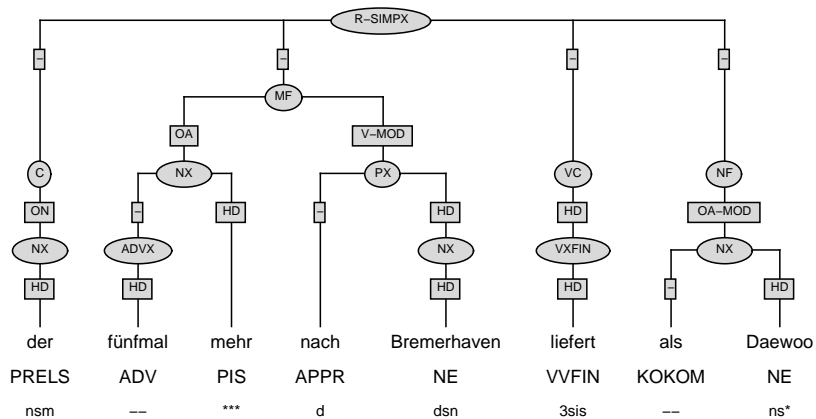
7.1.2 The Comparative Particles *wie* and *als*

Comparative particles in German are *als* and *wie*, in rare cases also *denn* (e.g. *Die werden dort seliger schlummern denn je.*). These particles are tagged as KOKOM and occur with all types of syntactic phrases (NX, ADVX, PX, etc.). They are directly attached to an adjacent comparative phrase. In case of a comparative phrase with a postmodifier, they are directly attached to the highest node of the complex phrase.

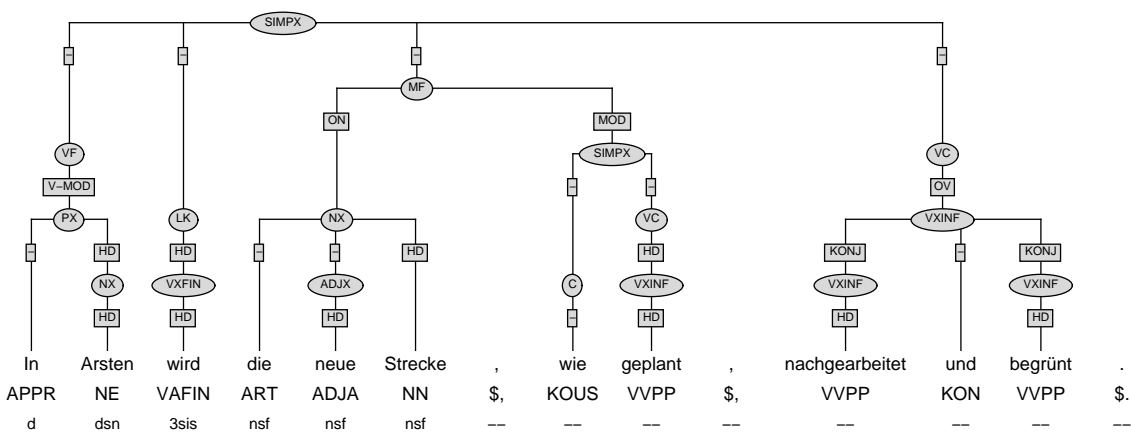
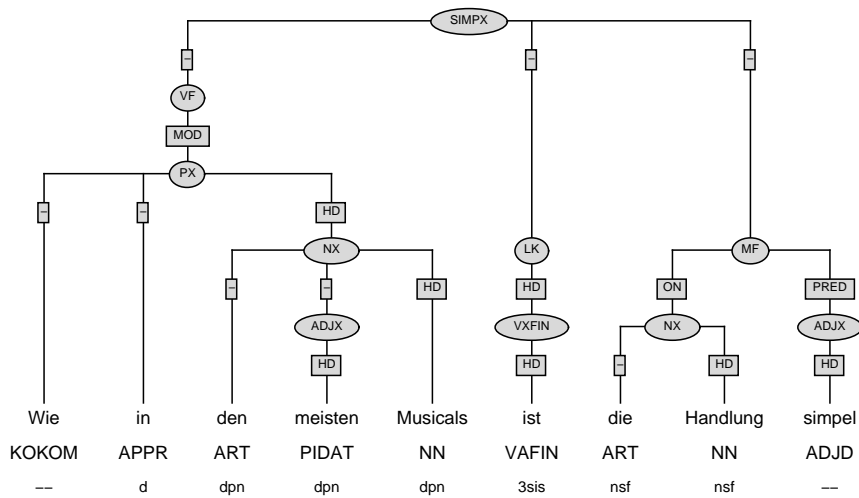
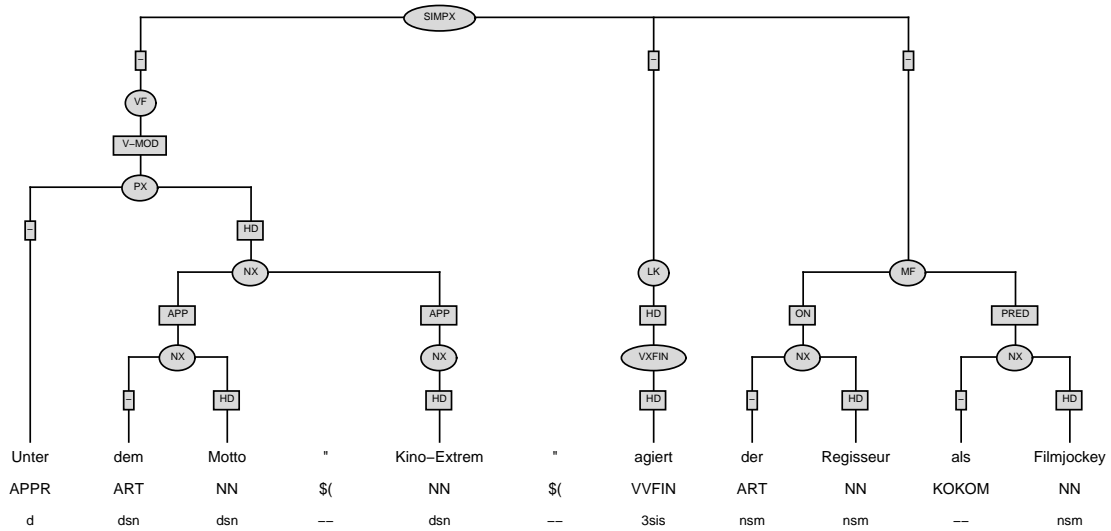
A comparative phrase can occur as an adjacent postmodifier of the head phrase:



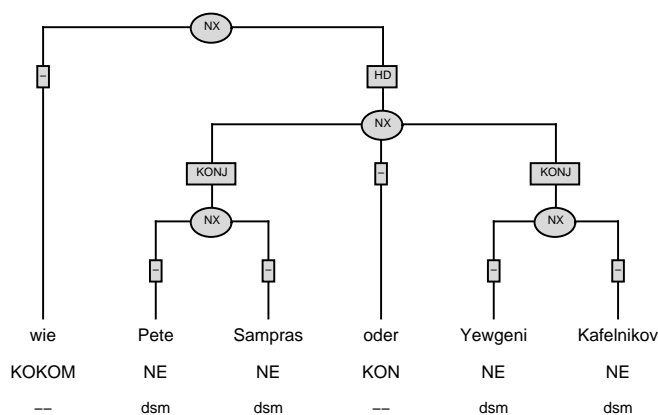
If there is a long-distance dependency between the comparative phrase and the head phrase, the dependency relation is denoted with the respective X-MOD label.



In case of a long-distance dependency between the comparative phrase and the main verb (cf. 4.7.9), the comparative phrase is either a complement (e.g. PRED) or an ambiguous or unambiguous modifier of the main verb (MOD or V-MOD).



The *high attachment principle* applies when the comparative particle has scope over a coordination of phrases (cf. 6.5.5). In this case, the two conjuncts are coordinated first. Then the particle is attached on a higher level.



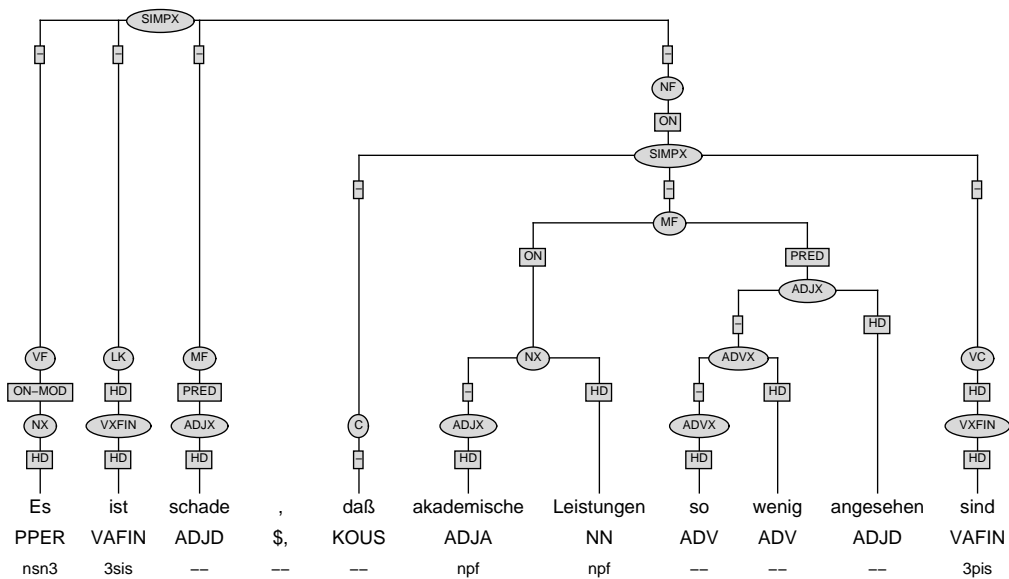
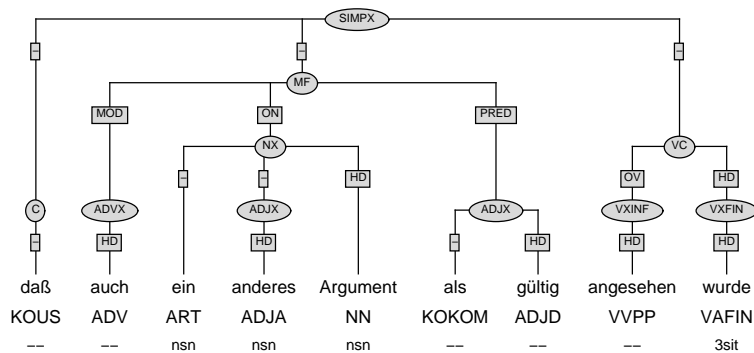
7.2 Verbal and Adjectival Use of Participles

In German, verbal participles which are passive verb forms (*Der Mensch wird angesehen*) can be used as adjectives: it can either function as an attribute adjective (*der angesehene Mensch*) or - depending on the context - also as a predicative adjective (*der Mensch ist angesehen.*). In contrast to the auxiliary *werden* in verbal passives, the auxiliary *sein* is used in constructions with adjectival passives. Concerning the problematic distinction between verbal and adjectival passives, we adapted the criteria in the Stuttgart-Tübingen tag set (STTS) (Schiller et al. 1995).¹

1. Can the sentence be transformed into active form keeping the same semantics? If yes → VVPP
2. Is there a *von*-PP or an equivalent PP that gives evidence for verb semantics? If yes → VVPP
3. Is it possible to substitute the word in questions by a semantically similar adjective? If yes → ADJD

The following two tree structures show the annotation of the verbal and adjectival passives of the verbal participle *angesehen*. In the first example, the verbal participle is analysed as a VVPP in VC. In the second example, the verbal participle has an adjectival reading and is annotated as an ADJD in MF.

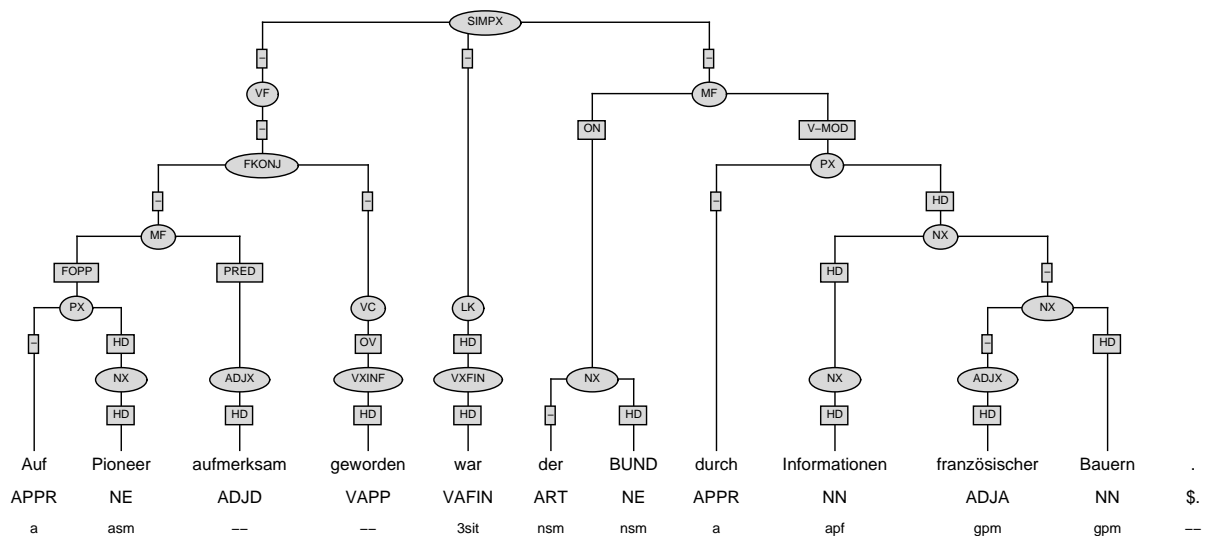
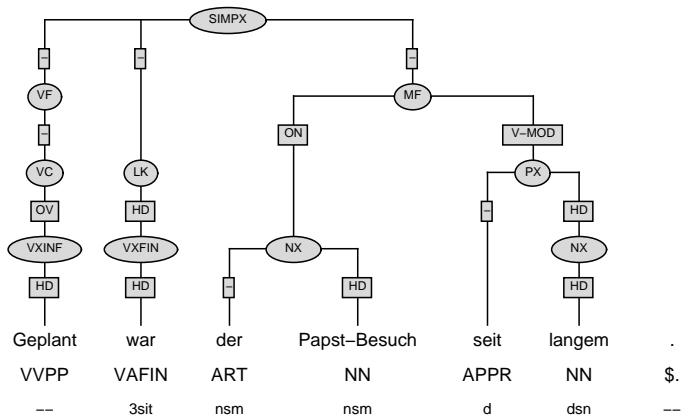
¹Concerning the differences between verbal and adjectival passives in English cf. Bresnan (1995).



7.3 Topicalization

Topicalization is almost exclusively found in verb-second clauses. Consequently, the subject is not in the first position of the clause. Topicalized constructions bring about word order phenomena which differ from those occurring in MF, e.g., non-finite parts of VC are not allowed in MF.

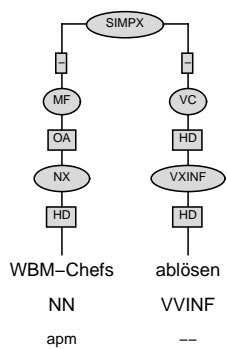
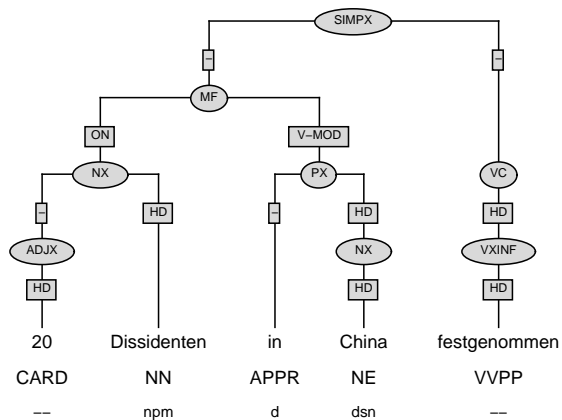
Our annotation principles demand to analyse the topicalized verb complex and its non-finite parts as VC in the first position of the clause. VC is then attached to VF. If a part of MF is topicalized along with VC, first MF and VC are combined to form FKONJ before they are attached to VF:



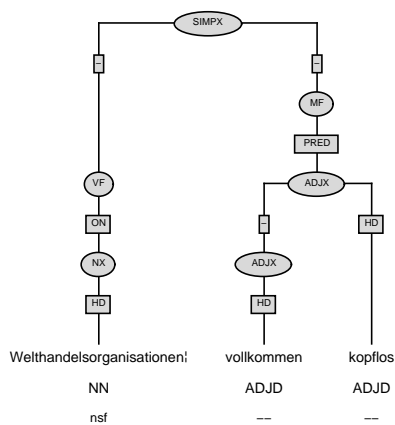
7.4 Headlines

The syntax of headlines differs from other syntactic constructions in so far as headlines² often lack the finite verb or a verb at all. If a headline has only an infinitive, the case assignment follows the preference principle formulated in 5.2. Therefore, we assume in general the more plausible grammatical function in each case: a passive constructions with ON in MF if the verb in VC is a past participle and an active construction with OA in MF if the verb in VC is an infinitive.

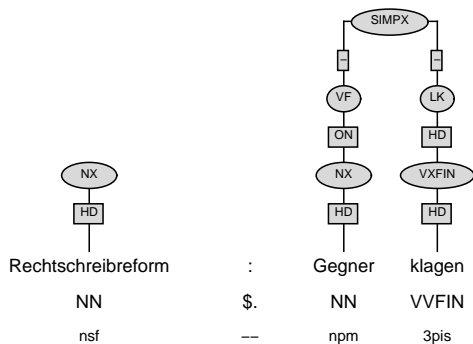
²The identifier "HEADLINE" is inserted into the comment line above the sentence for each syntactic unit which is marked as a headline in the original data.



A headline can also consist of an elliptical sentence (cf. 6.6):



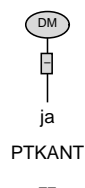
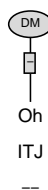
Headlines can also consist of more than one syntactic structure, for instance, separated by a colon or a dash (cf. 4.7.2 and 5.2):



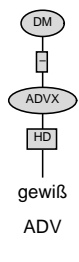
7.5 Discourse Markers

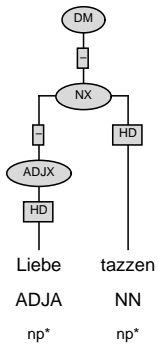
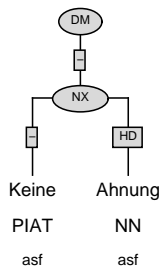
Generally, discourse markers are expressions or phrases of greeting, apologizing, thanking, short emotional utterances, and interjections. Their node label is DM. The edge label of a discourse marker is empty, i.e., it does not have a head. Typical discourse markers are: *ja, nein, hallo, oh, aha, pst, nunja, gewiß, toll, nun ja*, etc.

In most cases, discourse markers occur as isolated expressions. Interjections, tagged as ITJ, are directly projected to DM without internal structure. The same applies for answer particles (PTKANT):

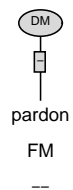
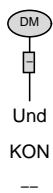


Phrases which function as discourse markers are first projected to their phrase level before they are assigned the node label DM.

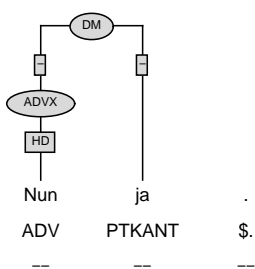




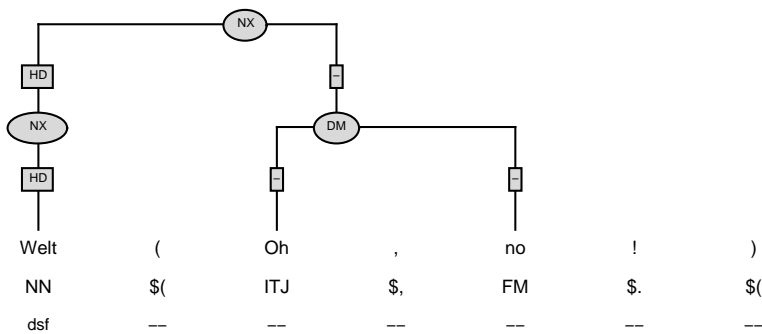
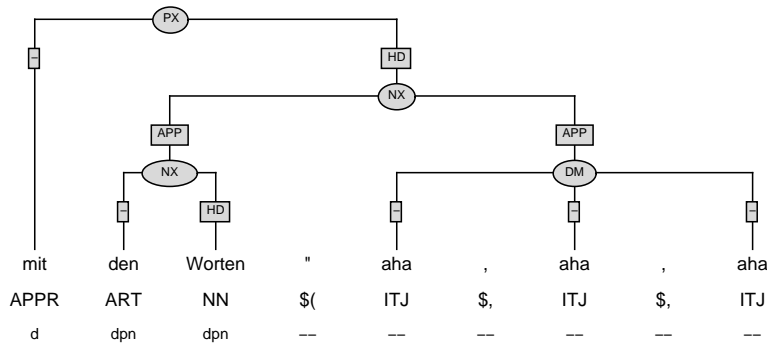
Isolated conjunctions and foreign language discourse markers are tagged according to their part of speech (KON and FM) and are projected to DM:



Discourse markers may also consist of an interjection or an answer particle and a phrase:

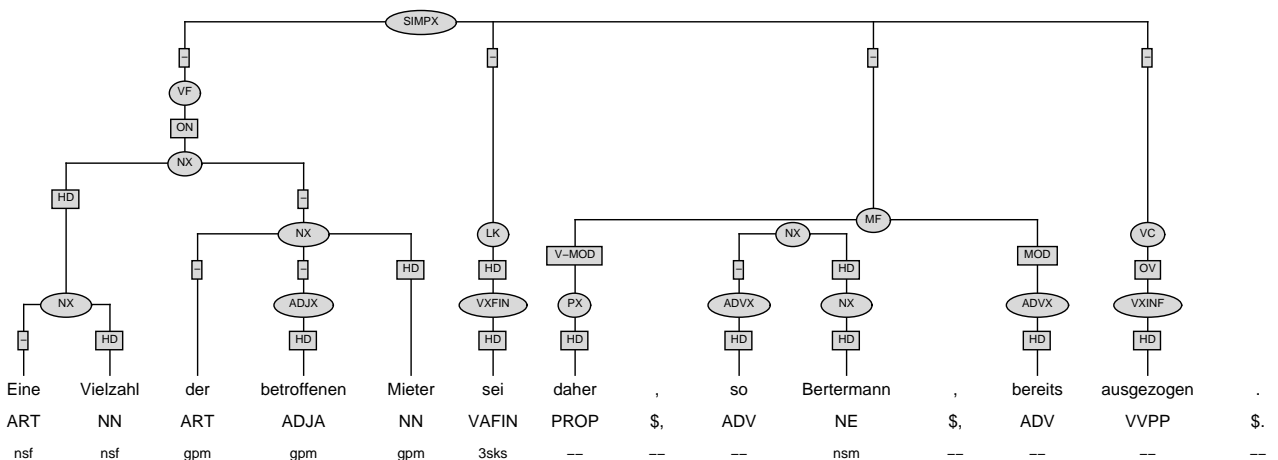


In some cases, discourse markers have a grammatical function within a phrase or a clause. Therefore, they are attached to the syntactic structure:

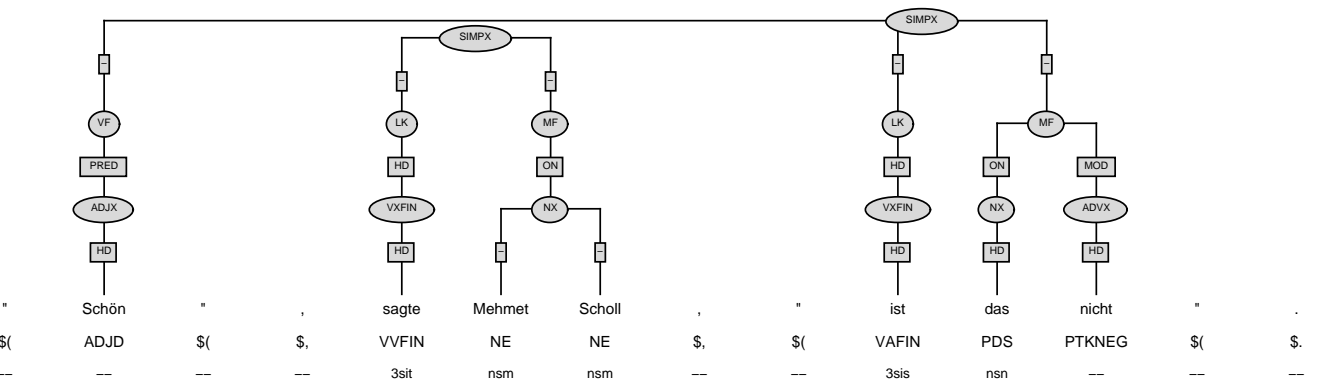
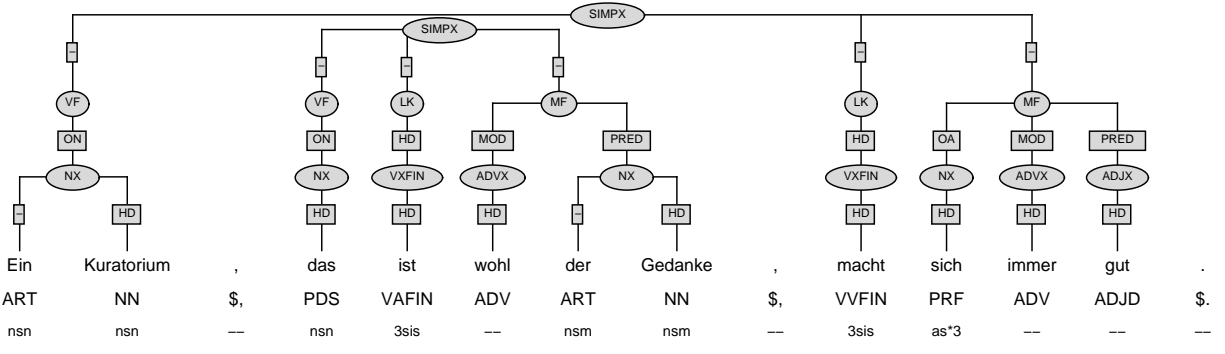


7.6 Parentheses

Parentheses occur as interjective utterances within a sentence. Since there is no dependency relation between the parenthesis and the rest of the construction, the parenthesis is not attached to the surrounding constituents. Often parentheses occur as SIMPX-clauses. Insertions like *sagte Mehmet Scholl* into direct speech are also annotated as parenthesis.³



³On the TüBa-D/Z web page (<http://www.sfs.uni-tuebingen.de/en/ascl/resources/corpora/tueba-dz.html>), the treebank is also available in the Penn Treebank formats version 1 and 2. In version 1, parentheses are attached to the tree structure with the edge label PAR. For further details about the Penn Treebank formats cf. 9.



Chapter 8

Criteria for the Distinction of Grammatical Functions

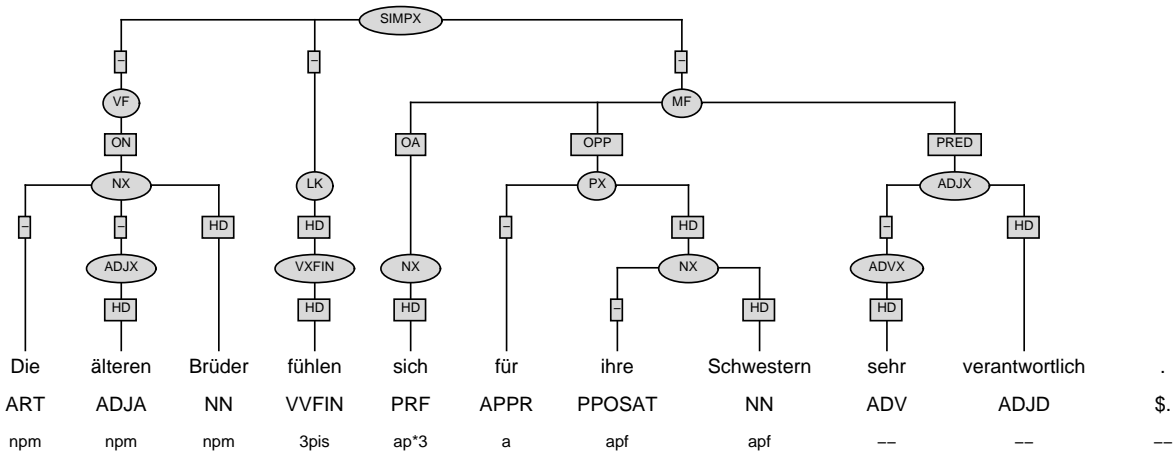
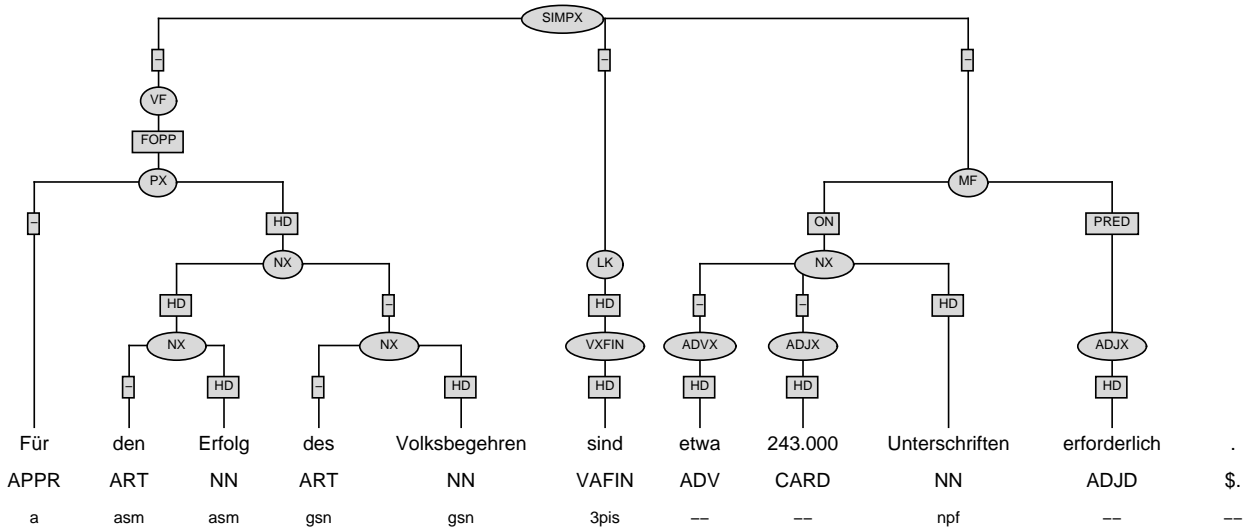
8.1 Subcategorization of Verbs

The *TüBa-D/Z-Verblast* document¹ lists all verbs occurring in the treebank with their specific subcategorization frames. This reference list guarantees the consistent annotation of grammatical functions. For a detailed description of constructing the verb list see (Hinrichs and Telljohann 2009).

8.2 Subcategorization of PREDs

Since constituents which predicates subcategorize for have grammatical function within a sentence, they are neither marked as PRED-MOD nor attached to the predicate itself. These constituents are attached to a field and assigned the respective grammatical function like the constituent which is marked as FOPP and OPP in the following examples:

¹In case of interest, please refer to web page (<http://www.sfs.uni-tuebingen.de/en/ascl/resources/corpora/tueba-dz.html>) for contact information.



8.3 Distinction of FOPP, OPP, and V-MOD

One of the major problems is to distinguish, whether a given PP is an obligatory (OPP) or an optional (FOPP) complement of a specific verb in a specific reading, or whether it is a free adjunct (V-MOD) of that verb.

The *TüBa-D/Z-Verblst* is intended as a reference for these problematic cases.

In the following, we will briefly describe what criteria have been used in order to decide about the subcategorization with respect to PP complements/modifiers:

1. A PP is called **OPP** within a sentence if the sentence were ungrammatical without the OPP (or if there was at least a very noticeable change of meaning). For instance, *Sie gehen [OPP gegen die Faschisten] vor./ Das Gesetz ist [OPP in Kraft] getreten.*
2. A PP is called **FOPP** if it can be left out of this specific sentence without causing ungrammaticality (or a very noticeable change of meaning) and if its preposition is selected by this specific verb. For instance, *Insgesamt berichtet die Polizei [FOPP von 19 Festnahmen und 98 Ingewahrsamnahmen]./Später würden wir [FOPP über*

Auswandern] nachdenken. Here, the prepositions select these specific verbs and the PPs cannot be added to any arbitrary verb (which is possible for free adjuncts). In addition, in passive clauses, the subject of the original active clause, which has the form of a prepositional phrase, is marked as FOPP (*Sie wurden [FOPP von Autonomen] umringt.*).

3. A PP is called **V-MOD** if its preposition is not selected by this specific verb, i.e., it can be exchanged by any other modifying PP, and similarly, this PP can occur with arbitrary verbs (*Nur [V-MOD im griechischen Lager] gab es Probleme*). Typical V-MODs are temporal or local adjuncts specifying time and location of the action, event, or state expressed by the verb.

8.4 Distinction of MOD, MOD-MOD, and V-MOD

A typical case of modification of modifiers is a temporal expression (V-MOD) that further specifies another temporal expression (MOD-MOD) in the same clause:

1. *[V-MOD am Samstag] finden [MOD-MOD ab 16 Uhr] Führungen statt.*
[MOD-MOD Wann] finden [V-MOD am Samstag] Führungen statt?

2. *[MOD da] finden [V-MOD am Samstag] Führungen statt.*
[V-MOD wann] finden [MOD da] Führungen statt?

[MOD dann] finden [V-MOD am Samstag] Führungen statt.

da, dann, etc. can be either temporal, causal, consequential, or local expressions. Thus, one cannot make sure whether the following time expression *am Samstag* really refers to them. The only obvious observation is that the *time* expression is a V-MOD in any case.

For resumptive constructions (LV), there is also a clear criterion concerning the modification relations. Within a verb-second clause, a modifier occurring in VF is MOD/X-MOD, whereas the modifier in LV is MOD-MOD, not vice versa, because the modifier in VF occurs within the core of the sentence, whereas the modifier in LV has to be licensed by some other constituent in the core sentence, e.g. *Wenn da was gebucht worden ist, dann ist das nicht in Ordnung.* (cf. 6.1.4).

8.5 Distinction of ON, PRED, ON-MOD, and PRED-MOD

It is not always trivial to distinguish which constituent is ON, PRED, or ON-MOD for predicative verbs. For this reason, a few criteria and examples are listed here that can be of help. Here are some properties of ON and PRED:

1. Typically, PRED occurs in MF, whereas ON occurs in VF of verb-second clauses. This should be considered for annotation, if no other criterion (as described below) applies.

2. Subject-verb agreement always has to be taken into account. For instance, if the verb is in plural form, the subject has to be plural as well.
3. If there is a suitable NP that could serve as subject, then this NP is annotated as subject rather than any other constituent with a different syntactic category (PP, ADVP, etc.).

For verb-second clauses, it is important to follow these two steps in exactly this order to stick to the distributional criterion that has been chosen for the PRED/ON distinction:

1. Have a look at the constituent in VF. If it is an NP which might serve as subject and if it agrees with the verb, annotate it as ON.
2. If it does not agree with the verb, annotate it as PRED (ADJP, ADVP, PP, etc.).

Examples:

1. *[ON neue Wortschöpfungen] sind [PRED es] nur.
[PRED es] sind nur [ON neue Wortschöpfungen].
oder sind [PRED es] nur [ON neue Wortschöpfungen].
[PRED das] sind ohnehin [ON die schwächsten Partner].
[ON die schwächsten Partner] sind [PRED das] ohnehin.
oder sind [PRED das] ohnehin [ON die schwächsten Partner].*

Subject-verb agreement suggests that *neue Wortschöpfungen* und *die schwächsten Partner* are the subject, because of their plural form regardless in which field they occur.

2. *[ON die Ursache] war [PRED unklar].
[PRED unklar] war [ON die Ursache]
[ON Candan Ercettin] ist [PRED überall].
[PRED überall] ist [ON Candan Ercettin].*

ADJPs and ADVPs typically have PRED function when occurring together with predicative verbs and NP subjects.

3. *[PRED aus den Trauernden] wird [ON ein wütender Mob].*

ein wütender Mob is considered the subject, because it is a noun phrase. Therefore, the prepositional phrase is PRED.

4. *[ON das] ist [PRED eine einmalige Chance].
[ON eine einmalige Chance] ist [PRED das].
[ON es] ist [PRED der erste Besuch eines Papstes].
[ON der erste Besuch eines Papstes] ist [PRED es].*

[ON Hauptauftraggeber] ist [PRED die Bremer Verwaltung].
[ON die Bremer Verwaltung] ist [PRED Hauptauftraggeber].

The NP in VF position agrees with the verb and therefore has subject priority. As a consequence, the constituent in MF is PRED.

5. *[PRED wer] bin [ON ich].*
[PRED was] ist [ON das].

In w-questions, the interrogative pronoun is always PRED because here also the agreement rule applies.

6. *[ON-MOD es] sei [PRED wichtig], [ON daß man ...].*
[ON Aufgabe des Festspielhauses] sei [PRED-MOD es], [PRED das Haus spielfertig zu halten].

If a sentential subject or a sentential predicate occurs with an expletive *es*, the expletive *es* is either ON-MOD or PRED-MOD (cf. 4.2.10).

Chapter 9

The TüBa-D/Z Data Formats

The TüBa-D/Z treebank is released in four different data formats :

1. the NEGRA Export format
2. the Penn Treebank format version 1 and version 2
3. the Export-XML format (incl. anaphora and coreference relations)
4. the CoNLL format 2006, 2010 and 2011/12

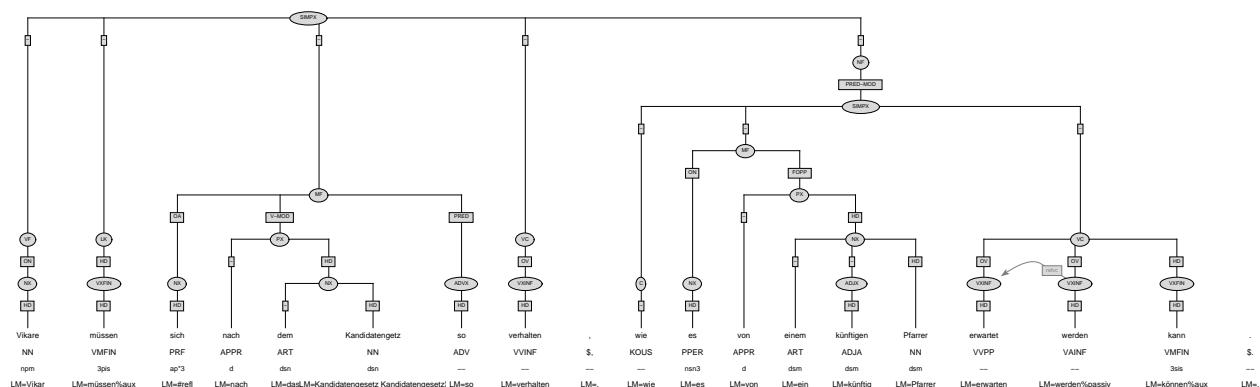
9.1 The NEGRA Export Format

This format is provided by the annotation tool Annotate (Brants and Skut 1998), it is created automatically from the database underlying the annotation process in Annotate. The NEGRA Export format is a line-oriented pointer-based representation of the syntactic annotation. It is also the most complete data format since it preserves all the information available during the manual annotation. A more complete description of the negra Export format can be found in (Brants 1997).

There are two versions of this format; NEGRA Export format 3 contains 3 layers of token information (token, POS tags, morphology), lemmata are displayed in the ‘comment’ column; NEGRA Export format 4 contains a 4th layer of lemma annotation.

An example of the NEGRA Export format 4 is given below, combined with the graphical representation of the syntactic annotation for the sentence ”Vikare müssen sich nach dem Kandidatengesetz so verhalten, wie es von einem künftigen Pfarrer erwartet werden kann”.

Graphical representation (print out of the annotate tool):



The first line of the sentence representation (marked as 'begin of sentence' (BOS)) includes the sentence id (here: 24539), the identity of the last annotator (here the one with id 2), the time of the last modification (in UNIX format, i.e. seconds since 1/1/1970) and the id of the origin of the file (1146 points to article 155 of the edition of 11/7/1992).

In the right column, secondary edges (here: 'refvc' pointing from node # 518 to node # 517, a dependency within the verbal complex) as well as corrections of misspellings (here: 'Kandidatengesetz') are also represented. Optionally, there is a version of NEGRA Export format 4 that contains anaphoric relations (here: 'es' is marked as expletive).

Export format 4:

```

#BOS 24539 2 1134150923 1146
Vikare          Vikar          NN      npm          HD      500
müssen         müssen%aux    VMFIN   3pis         HD      502
sich           #refl        PRF     ap*3         HD      504
nach           nach          APPR    d            -       506
dem            das           ART     dsn          -       505
Kandidatengetz Kandidatengesetz NN      dsn          HD      505      %% Kandidatengesetz
so             so            ADV     --           HD      507
verhalten      verhalten     VVINF   --           HD      509
,              ,             $,      --           --      0
wie            wie           KOUS    --           -       511
es             es            PPER    nsn3         HD      512
von            von           APPR    d            -       515
einem         ein           ART     dsm          -       514
künftigen     künftig      ADJA    dsm          HD      513
Pfarrer       Pfarrer      NN      dsm          HD      514
erwartet      erwarten     VVPP    --           HD      517
werden        werden%passiv VAINF   --           HD      518
kann          können%aux    VMFIN   3sis         HD      519
.             .            $.      --           --      0
#500          --           NX      --           ON      501
#501          --           VF      --           -       523
#502          --           VXFIN   --           HD      503
#503          --           LK      --           -       523
#504          --           NX      --           OA      508
#505          --           NX      --           HD      506
#506          --           PX      --           V-MOD    508
#507          --           ADVX    --           PRED     508
#508          --           MF      --           -       523
#509          --           VXINF   --           OV       510
#510          --           VC      --           -       523
#511          --           C       --           -       521
#512          --           NX      --           ON      516      %% R=expletive
#513          --           ADJX    --           -       514
#514          --           NX      --           HD      515
#515          --           PX      --           FOPP     516
#516          --           MF      --           -       521
#517          --           VXINF   --           OV       520
#518          --           VXINF   --           OV       520      refvc   517
#519          --           VXFIN   --           HD      520
#520          --           VC      --           -       521
#521          --           SIMPX   --           PRED-MOD 522
#522          --           NF      --           -       523
#523          --           SIMPX   --           --      0
#EOS 24539

```

The only deviation from context-freeness which the annotation scheme allows concerns the annotation of parentheses. Parentheses are annotated as separate trees with no attachment to surrounding trees. The following tree gives an example for such a phenomenon (for a more complete description of the annotation cf. 7.6).

9.2 The Penn Treebank Format

There exist two versions of the Penn Treebank format which will be introduced in the following sections.

9.2.1 The Penn Treebank Format Version 1

Version 1 is based on the format of the Penn Treebank (Mitchell et al. 1993). The attachment of constituents is shown via bracketing and indentation. Thus, all constituents which show the same level of indentation are attached on the same level. In the Penn Treebank format, grammatical functions, which are shown in the NEGRA Export format in the column "edge label", are attached to the syntactic label via a colon. Thus, the label "NX:OA" means that the constituent is a noun phrase with the grammatical function accusative object.

The Penn Treebank format is a representation that combines the linear representation of words with their attachment to higher constituents. For this reason, this format is restricted to completely context-free tree structures, i.e. it cannot adequately represent the annotation of parentheses in TüBa-D/Z. In order to capture the original syntactic annotation as well as the original word order in the sentence, it was decided to introduce a new edge label to mark such cases: PAR. Thus, the sentence "So etwas , sagen die Abgeordneten , hätten sie auch noch nicht erlebt .", as shown above is represented in the Penn Treebank format by the following bracketed structure:

Comments are preceded by a double '%' sign. The comment behind the structure is intended to help the reader locate the beginning of the parenthesis and it is not part of the actual data.

%% sent. no. 7307

```
(
(SIMPX
  (VF
    (NX:OA
      (ADVX
        (ADV:HD So)
      )
    )
    (PIS:HD etwas)
  )
)
)
($, ,)
(SIMPX:PAR                                %% here starts the parenthesis!
  (LK
    (VXFIN:HD
      (VVFIN:HD sagen)
    )
  )
  (MF
    (NX:ON
      (ART die)
      (NN:HD Abgeordneten)
    )
  )
)
)
($, ,)
(LK
  (VXFIN:HD
    (VAFIN:HD hätten)
  )
)
(MF
  (NX:ON
    (PPER:HD sie)
  )
  (ADVX:MOD
    (ADV:HD auch)
  )
  (ADVX:MOD
    (ADVX
      (ADV:HD noch)
    )
    (PTKNEG:HD nicht)
  )
)
)
)
)
)
($. .)
)
```


Commas, which are not attached to the tree, are indented on the highest level although they are included in the bracketing of the constituent surrounding them. In the sentence below, e.g., the first comma is grouped into the noun phrase NX via word order. The indentation, however, signals that the comma cannot necessarily be attached to this node. It is also conceivable that it may be attached to one of the lower nodes, NX or R-SIMPX. In the case of the second comma, there are even more possible attachment sites.

%% fragment of sent. no. 33

```
(
  (R-SIMPX
    (C
      (NX:ON
        (PRELS:HD die)
      )
    )
  )
  (MF
    (NX:OA
      (NX=ORG:HD
        (ART:-NE die)
        (NN:HD AWO)
      )
    )
  ($, ,)
  (R-SIMPX
    (C
      (PX:V-MOD
        (PWAV:HD wo)
      )
    )
    (MF
      (NX:ON
        (PPER:HD er)
      )
      (NX:PRED
        (NN:HD Kreisvorsitzender)
      )
    )
    (VC
      (VXFIN:HD
        (VAFIN:HD ist)
      )
    )
  )
  )
  )
  )
  ($, ,)
  (VC
    (VXFIN:HD
      (VVFIN:HD prüfte)
    )
  )
  )
  )
  ($ . .))
```

9.2.2 The Penn Treebank Format Version 2

Version 2 of the Penn Treebank format has no unattached phrases. For comparison, the same sentences used to demonstrate the Penn Treebank version 1 format are repeated here in version 2 format. Please note that in the data each sentence is on one line, and that the multi-line formatting here is for the human reader only:

```
%% sent. no. 7307
(VROOT:--
  (SIMPX:--
    (VF:-
      (NX:OA
        (ADVX:-
          (ADV:HD So)
        )
        (PIS:HD etwas)
      )
    )
    ($,:-- ,)
    (SIMPX:--
      (LK:-
        (VXFIN:HD
          (VVFIN:HD sagen)
        )
      )
      (MF:-
        (NX:ON
          (ART:- die)
          (NN:HD Abgeordneten)
        )
      )
    )
    ($,:-- ,)
    (LK:-
      (VXFIN:HD
        (VAFIN:HD hätten)
      )
    )
    (MF:-
      (NX:ON
        (PPER:HD sie)
      )
      (ADVX:MOD
        (ADV:HD auch)
      )
      (ADVX:MOD
        (ADVX:-
          (ADV:HD noch)
        )
        (PTKNEG:HD nicht)
      )
    )
    (VC:-
      (VXINF:OV
        (VVPP:HD erlebt)
      )
    )
  )
  ($,:-- .)
)
```

Sentence #33 fragment

```
(R-SIMPX:MOD
  (C:-
    (NX:ON
      (PRELS:HD die)
    )
  )
  (MF:-
    (NX:OA
      (NX=ORG:HD
        (ART:-NE die)
        (NN:HD AWO)
      )
      ($, :-- ,)
      (R-SIMPX:-
        (C:-
          (ADVX:V-MOD
            (PWAV:HD wo)
          )
        )
        (MF:-
          (NX:ON
            (PPER:HD er)
          )
          (NX:PRED
            (NN:HD Kreisvorsitzender)
          )
        )
        (VC:-
          (VXFIN:HD
            (VAFIN:HD ist)
          )
        )
      )
    )
  )
  ($, :-- ,)
  (VC:-
    (VXFIN:HD
      (VVFIN:HD prüfte)
    )
  )
)
```

9.3 The Export-XML Format

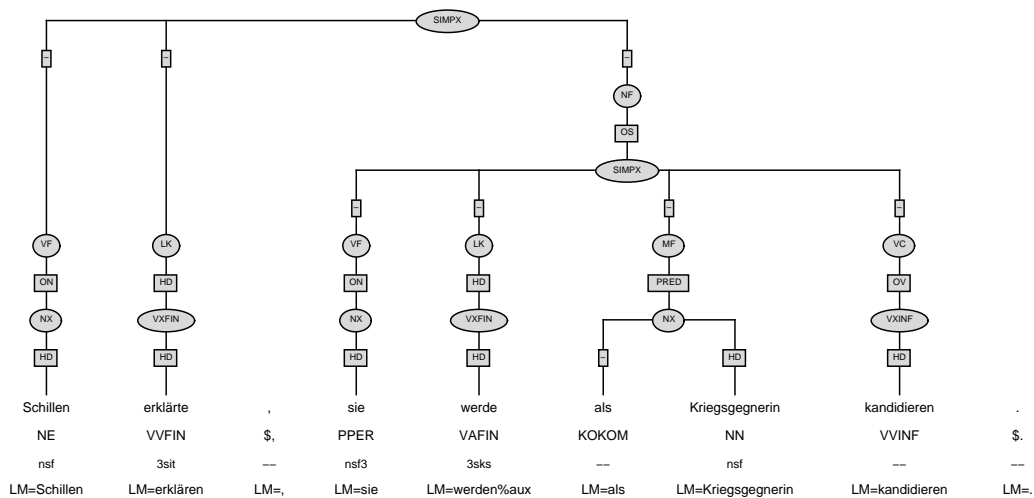
The XML format is a custom-made XML format that follows the NEGRA Export file format. It is designed to accommodate all original information provided in the Export format, including e.g. comments . Dominance relations between nodes are represented directly within the XML tree structure. Nodes without parent node (e.g. the sentence node SIMPX and punctuation marks do not have any "parent" attribute. They have the edge label "-" that can be linked to an implicit root node. Thus, it is possible to represent parentheses without the use of additional labels.

Anaphora is expressed by a link between two related nodes. Coreference sets therefore are represented implicitly by chains of nodes that are part of a referential relation.

The following example shows the XML structure for the sentence "Schillen erklärte, sie werde als Kriegsgegnerin kandidieren". The personal pronoun "sie" is anaphoric to the antecedent noun phrase "Schillen". In the XML document, a <relation> tag is added below each node that is part of a referential relation. It encodes the type of referential relation and the node ID of the antecedent node. In our example, the antecedent is the node with ID s1723_500, that is the NX dominating the named entity "Schillen". This NX in turn is in a coreferential relationship with node s1721_11 (word number 11 in sentence 1721), thus part of a coreference chain.

For extensive documentation of the ExportXML format as well as a Java API for reading the format, please visit the webpage at: <http://www.sfs.uni-tuebingen.de/en/ascl/resources/corpora/export-format.html>

Graphical representation of the tree without annotation of the referential relation:



XML format including the referential relation:

```
<sentence xml:id="s1723">
  <node xml:id="s1723_514" cat="SIMPX" func="--">
    <node xml:id="s1723_501" cat="VF" func="-" parent="s1723_514">
      <node xml:id="s1723_500" cat="NX" func="ON" parent="s1723_501">
        <relation type="coreferential" target="s1721_11"/>
        <ne xml:id="ne_29507" type="PER">
          <word xml:id="s1723_1" form="Schillen" pos="NE" morph="nsf" lemma="Schillen" func="HD" parent="s1723_500"/>
        </ne>
      </node>
    </node>
  </node>
  <node xml:id="s1723_503" cat="LK" func="-" parent="s1723_514">
    <node xml:id="s1723_502" cat="VXFIN" func="HD" parent="s1723_503">
      <word xml:id="s1723_2" form="erklärte" pos="VVFIN" morph="3sit" lemma="erklären" func="HD" parent="s1723_502"/>
    </node>
  </node>
  <word xml:id="s1723_3" form="," pos="$," lemma="," func="--"/>
  <node xml:id="s1723_513" cat="NF" func="-" parent="s1723_514">
    <node xml:id="s1723_512" cat="SIMPX" func="OS" parent="s1723_513">
      <node xml:id="s1723_505" cat="VF" func="-" parent="s1723_512">
        <node xml:id="s1723_504" cat="NX" func="ON" parent="s1723_505">
          <relation type="anaphoric" target="s1723_500"/>
          <word xml:id="s1723_4" form="sie" pos="PPER" morph="nsf3" lemma="sie" func="HD" parent="s1723_504"/>
        </node>
      </node>
    </node>
    <node xml:id="s1723_507" cat="LK" func="-" parent="s1723_512">
      <node xml:id="s1723_506" cat="VXFIN" func="HD" parent="s1723_507">
        <word xml:id="s1723_5" form="werde" pos="VAFIN" morph="3sks" lemma="werden%aux" func="HD" parent="s1723_506"/>
      </node>
    </node>
    <node xml:id="s1723_509" cat="MF" func="-" parent="s1723_512">
      <node xml:id="s1723_508" cat="NX" func="PRED" parent="s1723_509">
        <word xml:id="s1723_6" form="als" pos="KOKOM" lemma="als" func="-" parent="s1723_508"/>
        <word xml:id="s1723_7" form="Kriegsgegnerin" pos="NN" morph="nsf" lemma="Kriegsgegnerin" func="HD" parent="s1723_508"/>
      </node>
    </node>
  </node>
  <node xml:id="s1723_511" cat="VC" func="-" parent="s1723_512">
    <node xml:id="s1723_510" cat="VXINF" func="OV" parent="s1723_511">
      <word xml:id="s1723_8" form="kandidieren" pos="VVINF" lemma="kandidieren" func="HD" parent="s1723_510"/>
    </node>
  </node>
</node>
</node>
</node>
</node>
<word xml:id="s1723_9" form="." pos="$." lemma="." func="--"/>
</sentence>
```

9.4 The CoNLL Format (2006, 2010, 2011/2012)

In the following, the different CoNLL formats will be presented.

9.4.1 The CoNLL 2006 Format

The CoNLL format contains a dependency version of TüBa-D/Z in the format of the CoNLL-X shared task. The conversion was done automatically, but is oriented at the annotation guidelines by Foth (2006).

The CoNLL format is a table format containing a series of tabular-separated lines. Each line contains the following information:

1. ID – a sequential ID for each token
2. FORM – the word form (token)
3. LEMMA – the gold standard lemma of the token
4. CPOSTAG – simplified part-of-speech tag
5. POSTAG – part-of-speech tag according to STTS tag set
6. FEATS – tag with morphological information
7. HEAD – regent of the token in the dependency analysis, or “0” for tokens without regent
8. DEPREL – dependency relation between the token and its regent, or “ROOT” for tokens without regent
9. PHEAD – Projective head of current token (always “_”)
10. PDEPREL – Dependency relation to PHEAD (always “_”)

Sentence 1723 in CoNLL 2006 format:

1	Schillen	Schillen	N	NE	nsf	2	SUBJ	-	-
2	erklärte	erklären	V	VVFIN	3sit	0	ROOT	-	-
3	,	,	\$,	\$,	--	2	-PUNCT-	-	-
4	sie	sie	PRO	PPER	nsf3	5	SUBJ	-	-
5	werde	werden%aux	V	VAFIN	3sks	2	S	-	-
6	als	als	KOKOM	KOKOM	--	8	KOM	-	-
7	Kriegsgegnerin	Kriegsgegnerin	N	NN	nsf	6	CJ	-	-
8	kandidieren	kandidieren	V	VVINF	--	5	AUX	-	-
9	.	.	\$.	\$.	--	8	-PUNCT-	-	-

See <http://ilk.uvt.nl/conll/> for more details about the CoNLL 2006 format.

9.4.2 The CoNLL 2010 Format

The CoNLL 2010 format was used in the SemEval-2010 shared task, Coreference Resolution In Multiple Languages, and contains the following columns:

1. ID - word identifiers in the sentence
2. TOKEN - word forms
3. LEMMA - word lemmas (gold standard manual annotation)
4. *PLEMMA - word lemmas predicted by an automatic analyzer
5. POS - coarse part of speech
6. *PPOS - same as 5 but predicted by an automatic analyzer
7. FEAT - morphological features (part of speech type, number, gender, case, tense, aspect, degree of comparison, etc., separated by the character "—")
8. *PFEAT - same as 7 but predicted by an automatic analyzer
9. HEAD - for each word, the ID of the syntactic head ('0' if the word is the root of the tree)
10. *PHEAD - same as 9 but predicted by an automatic analyzer
11. DEPREL - dependency relation labels corresponding to the dependencies described in 9
12. *PDEPREL - same as 11 but predicted by an automatic analyzer
13. NE - named entities
14. *PNE - same as 13 but predicted by a named entity recognizer
15. *PRED - predicates are marked and annotated with a semantic class label
16. *PPRED - Same as 13 but predicted by an automatic analyzer
17. COREF - coreference annotation in open-close notation

Columns marked with "*" are always filled with "_", since they are either predicted values or the information is not available in the TüBa-D/Z. These columns are included to conform to the format.

1	Schillen	Schillen	-	NE	-	nsf	-	2	-	SUBJ	-	(PER)	-	-	-	(0)
2	erklärte	erklären	-	VVFIN	-	3sit	-	0	-	ROOT	-	*	-	-	-	-
3	,	,	-	\$,	-	-	-	0	-	ROOT	-	*	-	-	-	-
4	sie	sie	-	PPER	-	nsf3	-	5	-	SUBJ	-	*	-	-	-	(0)
5	werde	werden%aux	-	VAFIN	-	3sks	-	2	-	S	-	*	-	-	-	-
6	als	als	-	KOKOM	-	-	-	8	-	KOM	-	*	-	-	-	-
7	Kriegsgegnerin	Kriegsgegnerin	-	NN	-	nsf	-	6	-	CJ	-	*	-	-	-	-
8	kandidieren	kandidieren	-	VVINF	-	-	-	5	-	AUX	-	*	-	-	-	-
9	.	.	-	\$.	-	-	-	0	-	ROOT	-	*	-	-	-	-

See <http://stel.uib.edu/semEval2010-coref/datasets> for more details about the CoNLL 2010 format.

9.4.3 The CoNLL 2011/2012 Format

The CoNLL 2011/2012 format contains the following columns:

1. Document ID - the newspaper article id in the form TYYMMDD.articleNumber
2. Part Number - the GLOBAL sentence ID. Numbering does not restart within each document, therefore the part number corresponds the sentence ID in the treebank. Thus a document should be solely identified by the doc ID.
3. Word number
4. Word itself
5. Part-of-Speech
6. Parse bit - represents parse in a bracketed structure
7. Predicate lemma - the lemma of every token is represented
8. *Predicate Frameset ID - PropBank frameset ID of the predicate in column 7
9. Word sense - GermanNet ID of the word sense
10. *Speaker
11. Named Entities
12. *Predicate Arguments
13. Coreference - coreference chain information

Columns marked with "*" are always filled with "-", since they are either predicted values or the information is not available in the TüBa-D/Z. These columns are included to conform to the format.

T990507.136	1723	1	Schillen	NE	(VROOT:--(SIMPX:--(VF:-(NX=PER:ON*)))	Schillen	-	-	-	(PER)	-	(0)
T990507.136	1723	2	erklärte	VVFIN	(LK:-(VXFIN:HD*))	erklären	-	-	-	*	-	-
T990507.136	1723	3	,	\$,	*	,	-	-	-	*	-	-
T990507.136	1723	4	sie	PPER	(NF:-(SIMPX:OS(VF:-(NX:ON*)))	sie	-	-	-	*	-	(0)
T990507.136	1723	5	werde	VAFIN	(LK:-(VXFIN:HD*))	werden%aux	-	-	-	*	-	-
T990507.136	1723	6	als	KOKOM	(MF:-(NX:PRED*	als	-	-	-	*	-	-
T990507.136	1723	7	Kriegsgegnerin	NN	*)	Kriegsgegnerin	-	-	-	*	-	-
T990507.136	1723	8	kandidieren	VVINFIN	(VC:-(VXINF:OV*))))	kandidieren	-	-	-	*	-	-
T990507.136	1723	9	.	\$.	*	.	-	-	-	*	-	-

See <http://conll.cemantix.org/2012/data.html> for more details about the CoNLL 2011/2012 format.

References

- Bech, G. 1955–57. *Studien über das deutsche Verbum infinitum*. Kopenhagen. 2 Bände. 2. unveränderte Auflage 1983 mit einem Vorwort von Catharine Fabricius-Hansen. Tübingen: Max Niemeyer.
- Behaghel, O. 1932. *Deutsche Syntax (Eine geschichtliche Darstellung), Band 4*. Heidelberg: Carl Winter.
- Brants, T., and W. Skut. 1998. Automation of treebank annotation. In *Proceedings of the Conference on New Methods in Language Processing (NeMLaP-3/CoNLL98), January 14-17, 1998, Sydney, Australia*, 49–57. Sydney.
- Brants, T. 1997. *The NeGra Export Format for Annotated Corpora*. Universität des Saarlandes, Computational Linguistics, Saarbrücken, Germany.
- Brants, T. 1998. *TnT – A Statistical Part-of-Speech Tagger*. Universität des Saarlandes, Computational Linguistics, Saarbrücken, Germany.
- Bresnan, J. 1995. Lexicality and Argument Structure. In *Invited Paper given at the Paris Syntax and Semantics Conference*. Paris. October 12-14, 1995. URL: <http://www-csli.stanford.edu/bresnan/download.html>.
- Drach, E. 1937. *Grundgedanken der Deutschen Satzlehre*. Frankfurt/Main.
- Drosdowski, G. (Ed.). 1995. *Duden "Die Grammatik der deutschen Gegenwartssprache"*. Mannheim, Leipzig, Wien, Zürich: Dudenverlag.
- Eisenberg, P. 1999–2001. *Grundriß der deutschen Grammatik, Band 2: Der Satz*. Stuttgart, Weimar: J.B. Metzler.
- Engel, U. 1996. *Deutsche Grammatik*. Heidelberg: Julius Groos Verlag.
- Erdmann, O. 1886. *Grundzüge der deutschen Syntax nach ihrer geschichtlichen Entwicklung dargestellt*. Stuttgart: Cotta. Erste Abteilung.
- Foth, K. A. 2006. Eine umfassende Constraint-Dependenz-Grammatik des Deutschen. Technical report, Fachbereich Informatik der Universität Hamburg.
- Grewendorf, G. 1991. *Aspekte der deutschen Syntax*. Vol. 33 of *Studien zur deutschen Grammatik*. Tübingen: Gunter Narr Verlag.
- Helbig, G., and J. Buscha. 1998. *Deutsche Grammatik. Ein Handbuch für den Ausländerunterricht*. Leipzig: Verl. Enzyklopädie. 18 edition.
- Herling, S. H. A. 1821. Über die Topik der deutschen Sprache. In *Abhandlungen des frankfurterischen Gelehrtenvereins für deutsche Sprache*, 296–362, 394. Frankfurt/Main. Drittes Stück.

- Hinrichs, E. W., and H. Telljohann. 2009. Constructing a Valence Lexicon for a Treebank of German. In *Proceedings of the Seventh International Workshop on Treebanks and Linguistic Theories (TLT 7): January 23-24, 2009, Groningen, The Netherlands*. URL: <http://www.let.rug.nl/tlt/>.
- Hinrichs, E. W., J. Bartels, Y. Kawata, V. Kordoni, and H. Telljohann. 2000. The Tübingen Treebanks for Spoken German, English, and Japanese. In W. Wahlster (Ed.), *Verbmobil: Foundations of Speech-to-Speech Translation*. Berlin: Springer.
- Höhle, T. N. 1986. Der Begriff ‘Mittelfeld’. Anmerkungen über die Theorie der topologischen Felder. In A. Schöne (Ed.), *Kontroversen alte und neue. Akten des 7. Internationalen Germanistenkongresses Göttingen*, 329–340. Tübingen: Niemeyer.
- Kathol, A. 1995. *Linearization-Based German Syntax*. PhD thesis, Ohio State University.
- Kiss, T. 1995. *Infinitive Komplementation. Neue Studien zum deutschen Verbum infinitum*. Tübingen: Max Niemeyer.
- Kübler, S., and H. Telljohann. 2002. Towards a dependency-based evaluation for partial parsing. In *Beyond PARSEVAL – Towards Improved Evaluation Measures for Parsing Systems – (LREC 2002 Workshop), Las Palmas, Gran Canaria, June 2002*.
- Mitchell, M., B. Santorini, and M. A. Marcinkiewicz. 1993. Building a Large Annotated Corpus of English: The Penn Treebank. *Computational Linguistics* 19(2):313–330.
- Naumann, K., and V. Möller. 2007. *Manual for the Annotation of in-document Referential Relations*. University of Tübingen, May 2007.
- Plaehn, O. 1998. *Annotate – Bedienungsanleitung*. FR 8.7 Computerlinguistik, Projekt C3 Nebenläufige Grammatische Verarbeitung, Sonderforschungsbereich 378, Ressourcennadaptive Kognitive Prozesse, 13. April 1998. URL: <http://www.coli.uni-saarland.de/projects/sfb378/negra-corpus/annotate.html>.
- Pütz, H. 1986. *Über die Syntax der Pronominalform ‘es’ im modernen Deutsch*. Tübingen: Stauffenburg. 2 edition.
- Schiller, A., S. Teufel, and C. Thielen. 1995. Guidelines für das Tagging deutscher Textcorpora mit STTS. Technical report, Universities of Stuttgart and Tübingen. URL: <http://www.sfs.nphil.uni-tuebingen.de/ELWIS/stts/stts.html>.
- Schnorr, V. 1991. Problems of Lemmatization in the Bilingual Dictionary. In F. J. Hausmann, O. Reichmann, H. E. Wiegand, and L. Zgusta (Eds.), *Wörterbücher Ein internationales Handbuch zur Lexikographie, Dritter Teilband*, 2813–2817. Walter de Gruyter, Berlin, New York.
- Stegmann, R., H. Telljohann, and E. W. Hinrichs. 2000. Stylebook for the German Treebank in VERBMOBIL. Verbmobil-report 239, University of Tübingen.
- Telljohann, H., E. W. Hinrichs, S. Kübler, H. Zinsmeister, and K. Beck. 2012. *Stylebook for the Tübingen Treebank of Written German (TüBa-D/Z)*. University of Tübingen, January 2012.

- Trushkina, J. 2004. *Morpho-syntactic annotation and dependency parsing of German*. PhD thesis, University of Tübingen. URL: <http://w210.ub.uni-tuebingen.de/dbt/volltexte/2004/1523>).
- Versley, Y., K. Beck, E. W. Hinrichs, and H. Telljohann. 2010. A Syntax-first Approach to High-quality Morphological Analysis and Lemma Disambiguation for the TüBa-D/Z Treebank. In *Proceedings of the Ninth International Workshop on Treebanks and Linguistic Theories (TLT 9): December 3-4, 2010, Tartu, Estonia*. URL: <http://www.math.ut.ee/tlt9/>.

Index

- accusative object, double, 89
- AcI, 88
- adverbial adjective, 21, 74, 76
- adverbial phrase, 25, 79
- ambiguity, 18, 31, 40, 96, 98–101, 120
- apposition, 26, 49
- attributive adjective, 21, 32, 42, 43, 74, 77

- C-field, 16, 17, 25, 102, 105, 109
- cardinal numbers, 21, 47, 64, 65
- circumposition, 21, 74
- coherency, 87, 88
- comparatives, 13, 130
- CoNLL format, 146
- context-freeness, 12
- coordination, 13, 15, 20, 25, 44, 72, 77, 112, 114–116, 118–124, 126, 132

- Dependency Grammar, 32, 71
- determiner phrase, 25, 70
- discourse marker, 11, 13, 22, 25, 32, 137, 138

- edge labels, 12, 17, 19, 20, 26–28, 31, 32, 97, 112
- elliptical construction, 13, 19, 22, 112, 127, 128, 136
- Ersatzinfinitiv, 16, 25, 83, 84
- expletive, 26, 67, 68, 145
- export-XML format, 146

- flat clustering principle, 18, 40, 81, 102
- foreign language material, 21, 53, 63

- headline, 11, 13, 22, 83, 96, 97, 135
- high attachment principle, 18, 40, 118, 132

- imperative, 22, 89
- incoherency, 88
- infinitives with *zu*, 85, 86, 88
- initial field, 14, 16, 25, 99
- isolated phrase, 28, 29, 100, 101, 124

- KOORD-field, 16, 17, 25, 104, 112

- lassen, 88
- lemma, 19, 33–38, 146
- lemmatization, 33–38
- levels of annotation, 19
- long-distance dependency, 12, 31, 105, 131
- longest match principle, 18, 22, 30, 122

- modal verbs, 22, 94

- named entities, 12, 20, 27, 48, 56–60, 63
- Negra export format, 146
- Negra treebank, 10
- node labels, 12, 19, 20, 25, 27, 53, 83, 101, 119, 137
- nominalized adjective, 76
- non-ambiguity, 31, 97
- non-words, 22, 66

- ordinal numbers, 64

- paratactic construction, 25, 122, 124
- parenthesis, 11, 13, 22, 139
- PARORD-field, 16, 17, 25, 104, 105, 124
- part-of-speech tags, 12, 19, 28, 66, 72
- particle verb, 91
- Penn Treebank format, 146
- postmodifier, 18, 40, 46, 48, 50, 55, 60, 62, 79, 99, 109, 110, 130
- postnominal modifier, 46, 47

postposition, 21, 74
 predicate, 91, 92, 141, 145
 predicate-argument structure, 11, 17
 predicative adjective, 21, 74, 75, 133
 preference principle, 97, 135
 premodifier, 40, 43, 60–62, 64, 75, 76, 78,
 79, 99, 120
 pronominal modifier, 42, 44, 46
 preposition, 21, 32, 41, 71, 72, 142, 143
 proper noun, 21, 27, 40, 41, 44, 52, 53,
 56, 63
 punctuation marks, 22, 28, 49, 99, 100

 relative clause, 16, 25, 98, 109, 110
 relative clause, event-modifying, 111
 relative clause, independent, 111
 resumptive construction, 15, 16, 25, 105,
 143
 reusability, 10, 11

 secondary edge label, 12, 20, 26, 31, 81,
 82, 88, 109
 split coordination, 26, 112, 126
 superlative forms, 130
 syntactic dependencies, 12
 syntactic-semantic node labels, 20, 57, 58

 TüBa-D/S treebank, 10
 TüBa-D/Z data formats, 10, 146
 theory-neutrality, 10, 11
 TIGER treebank, 10
 topicalization, 13, 134
 topological fields, 12–17, 19, 20, 31, 32,
 42, 92, 96, 102, 119, 127, 128
 truncated word, 22, 116

 verb complex, 12, 14–16, 25, 81–83, 85,
 87, 88, 91, 110, 134
 verb particle, 22, 90
 VERBMOBIL treebank, 10, 28