

## Non-thematic part

*Ralf Hauser / Angelika Storrer*

### Dictionary Entry Parsing Using the LEXPARSE System

1. Introduction
2. Dictionary Entry Parsing: Problems and Requirements
  - 2.1. Dictionary Entry Parsing
  - 2.2. Structural Properties of Dictionary Texts
  - 2.3. Dictionary Entry Parsing vs. Sentence Parsing
  - 2.4. Requirements for a Dictionary Entry Parser
3. LEXPARSE System Design and Concepts
  - 3.1. The System's Architecture
  - 3.2. LEXPARSE Concepts
4. The LEXPARSE Grammar Formalism
  - 4.1. Formal Description of the LEXPARSE Grammar
  - 4.2. Annotations to the Formalism
  - 4.3. Applying the LEXPARSE Grammar Formalism
5. The Implementation
  - 5.1. Module LEXPARSE
  - 5.2. Module GRAMMAR
  - 5.3. Module PPROCESS
  - 5.4. Module PREP
  - 5.5. Module SCAN
  - 5.6. Module PARSE
  - 5.7. Module OUTPUT
  - 5.8. Module ERROR
  - 5.9. Module GREP
  - 5.10. Module INI
  - 5.11. Module TOOLS
6. Conclusion and Outlook
7. References
- Appendixes
  - A List of All Available XCodes
  - B A Sample '.INI' Configuration File for LEXPARSE
  - C A Sample LEXPARSE Grammar
  - D Sample Sessions
    - D.1. **Treffen**, DUDEN-STILWÖRTERBUCH, Page 703
    - D.2. **trennen**, DUDEN-STILWÖRTERBUCH, Page 704
    - D.3. **überhaupt**, DUDEN-STILWÖRTERBUCH, Page 716

#### 1. Introduction

Dictionary entry parsing, i.e., the automatic conversion of a typesetting representation of dictionary text into a format which explicitly represents the hierarchical structure of dictionary

entries, has become an increasingly important topic in Computational Lexicography. Whereas earlier programs were tailored to the conversion of one specific dictionary, the growing interest in machine readable dictionaries led to the development of general-purpose dictionary entry parsers, which may be used in converting any type of typesetting representation of dictionary text (cf. NEFF & BOGURAEV 1990 and BLÄSI & KOCH 1991).

Dictionary entry parsers analyze an input stream coming from typesetting tapes using a dictionary entry grammar and then generate a representation of the dictionary text in which the hierarchical structure of the entries is explicitly shown so that each lexicographic item can be individually accessed by a software system. Such a representation of both dictionary form and dictionary subject is generally called *Lexical Database* and is used in different fields:

### Computational Lexicography:

Methods and techniques are developed and tested to transform the lexical information stored in lexical databases into a format which can be processed by a natural language processing (NLP) system (cf. BOGURAEV & BRISCOE 1989).

Since the manual construction of NLP lexicons is a time and labor consuming task, one hopes that the automatic exploitation of machine readable dictionaries (MRDs) using these methods will reduce the effort required in the development of NLP lexicon components (cf. BOGURAEV 1991).

### Traditional Lexicography:

Lexical databases are fundamental components of lexicographic workbenches, i.e., software systems for computer-aided dictionary publishing. Dictionary entries stored in the database can be edited and updated directly thereby speeding up the process of revising existing dictionaries and making new dictionaries. Furthermore, the multiple ways in which lexical databases can be queried allow for a more consistent and reflected description of the lexical items in the dictionaries (cf. BLÄSER & WERMKE 1990).

Electronic editions of dictionaries based on lexical databases can easily be developed offering a faster and more flexible access to the dictionary text using techniques such as hierarchical, relational or object-oriented databases or hypertext systems.

Dictionary entry parsers automatically assign structure to the dictionary text represented on typesetting tapes using a dictionary entry grammar. Dictionary entry grammars define conditions for well-formedness of dictionary entries and specify partitive and precedence relations between the constituents of the dictionary entry structure. Dictionary entry structures not licensed by the dictionary entry grammar will be marked as non-wellformed. As a consequence, the process of dictionary entry parsing – aside from its main goal of converting the typesetting tape into a lexical database – has the side-effect of detecting errors and inconsistencies in the structural encoding of the dictionary. Thus, not only can complete and correct dictionary entry grammars be used in the conversion step, but they may be reused for checking structural well-formedness of revised entries in future editions of the dictionary in question as well. A prerequisite for this use of dictionary entry grammars is that the grammar formalism is simple and easy to learn, so that changes in the dictionary entry structure, introduced in connection with the revision of the dictionary, can be adapted without difficulty.

Our paper deals with the problems of dictionary entry parsing and presents the main features of the LEXPARSE system, a dictionary entry parser developed by Ralf Hauser within the framework of the *ELWIS* project.<sup>1</sup>

LEXPARSE supports a context-free grammar formalism supplemented by special operators

<sup>1</sup> *ELWIS* is a research project on the corpus-based development of lexical knowledge bases (the acronym stands for *Korpusunterstützte Entwicklung lexikalischer Wissensbasen*) carried out at the University of Tübingen since 1992 and funded by the Ministry of Science and Research Baden-Württemberg (cf. STORRER, FELDWEG & HINRICHS 1993).

which are needed for the particular requirements of dictionary entry parsing. It has been tested with parts of the DUDEN-STILWÖRTERBUCH (DUDEN-2) and the DUDEN-BEDEUTUNGSWÖRTERBUCH (DUDEN-10).<sup>2</sup> and is currently being used for the parsing of the DEUTSCHES RECHTSWÖRTERBUCH (Akademie der Wissenschaften in Heidelberg/Germany) and the FRÜHNEUHOCHDEUTSCHES WÖRTERBUCH (German language department of the University of Heidelberg) and will be applied to a bilingual dictionary in an EC funded R&D project. In most cases, our examples will refer to the dictionary entry grammar which was developed for the DUDEN-STILWÖRTERBUCH (cf. ENGELKE 1994). Although the system was designed for the parsing of dictionary text, the system can also be used for the parsing of text types with similar structural properties such as bibliographies and encyclopedias.

The development of the system and of the above mentioned dictionary entry grammars was based on the theory on lexicographic texts by H.E. WIEGAND. In the following section we will briefly introduce some basic notions of this theory needed in describing the key problems of dictionary entry parsing. Most of these problems can be traced back to the fact that intellectual segmentation of dictionary text into meaningful lexicographic text segments presupposes an understanding of its semantic content, whereas automatic segmentation can only rely on formal properties of the text segments. We will show that the problems of dictionary entry parsing differ considerably from those related to the parsing of natural languages and we will compile a list of specific requirements which a parsing system has to meet in order to cope with these problems. We will then present a detailed description of the design and of the general concepts of LEXPARSE explaining how these concepts together with the grammar formalism lead to a powerful tool for dictionary entry parsing which will allow a straightforward treatment of these problems. Finally, we will describe the modular architecture of the system implemented in the object-oriented programming language C++ and outline how the modules of the system interact.

## 2. Dictionary Entry Parsing: Problems and Requirements

### 2.1. Dictionary Entry Parsing

A Dictionary entry parser is a software system which segments dictionary text stored on typesetting tapes into functional text segments and generates an explicit representation of the dictionary structure by consulting a user-supplied dictionary entry grammar. The development of dictionary entry grammars must be based on a theory which provides methods and categories for the analysis of the dictionary form and the dictionary subject. In recent metalexicographic research on printed dictionaries, the idea that dictionaries can be observed as a special text type to be described with textlinguistic categories has succeeded in being accepted. A formalized theory on lexicographic texts worked out by H.E. WIEGAND provides:

- a sound theoretical framework for the development of dictionary entry grammars by formally describing the various types of text structures to be found in standardized dictionaries;
- methods to intellectually analyze and segment entries of standardized monolingual dictionaries;
- a substantial terminological framework to elaborate the problems of dictionary entry parsing as well as the different strategies for their solution.

In the following we will informally introduce those basic notions of the theory which are

<sup>2</sup> We would like to thank Dr. MATTHIAS WERMKE and the DUDEN editorial for supplying us with this material.

needed to understand the basic concepts and strategies of the LEXPARSE approach to dictionary entry parsing.<sup>3</sup>

## 2.2. Structural Properties of Dictionary Texts

The main part of a general monolingual dictionary is the word list, consisting of lexicographic descriptions for lemma signs, i.e., those lexical signs which are represented as *lemmata* (headwords) in the dictionary. The basic textual units of a word list are the *dictionary entries*. Dictionary entries can be further segmented into functional text segments which are either *structure indicators* carrying information on the dictionary form, or *lexicographic items* carrying information on the dictionary subject. The *microstructure* and the *macrostructure* are both order structures which are characteristic for the word list of a dictionary and have to be recognized by a dictionary entry parser.<sup>4</sup>

### 2.2.1. The Macrostructure

The macrostructure defines the ordering of the lemmata in the word list by means of a specific *access alphabet*, i.e., an alphabet in which all kinds of diacritics are integrated. Strictly-alphabetically ordered dictionaries arrange their lemmata strictly according to this alphabet, whereby two kinds of arrangements have to be distinguished: *Niching* dictionaries cluster sub-entries, i.e., each sub-lemma immediately follows the preceding sub-entry without beginning a new line, whereas *straight-alphabetical* dictionaries have no lexicographic grouping. *Nested* dictionaries reveal the same type of grouping as niched ones, but are not strictly alphabetical, because the alphabetic order can be interrupted in order to exhibit morphosemantic relationships between the (sub)-lemmata.<sup>5</sup>

### 2.2.2. The Microstructure

The microstructure of a dictionary is determined by the classes of lexicographic items (*item classes*) used in the dictionary and specific relations which hold between them.<sup>6</sup> Two types of relations are relevant for the development of dictionary entry grammars to be used by a dictionary entry parser:

- *Partitive relations* are defined on the set of classes of items and determine the partitive microstructure.
- *Precedence Relations* are defined on the set of classes of items which cannot be further segmented into smaller functional text segments (*classes of elementary items*) thereby determining the precedential microstructure.

Both relations together determine the *hierarchical microstructure* of dictionary entries which can be represented by tree diagrams.<sup>7</sup> Classes of elementary items are the terminal nodes, classes of non-elementary items, i.e., complex categories such as example groups or lexico-

<sup>3</sup> Detailed and formal descriptions of different parts of the theory are given in different publications, e.g., in WIEGAND 1989a; 1989b; 1989c; an overview is given in WIEGAND 1991.

<sup>4</sup> Other types of structures described in the above mentioned theory on lexicographic texts, e.g., addressing structures, cohesion structures, theme-rheme structures, mediostructures have not yet been dealt with within the context of dictionary entry parsing.

<sup>5</sup> For a detailed description of the different kinds of access structures and macrostructures see WIEGAND 1989a.

<sup>6</sup> In our discussion we use the term 'microstructure' to refer to the 'abstract microstructure' in the sense of WIEGAND 1989b without further differentiating in abstract und concrete microstructures.

<sup>7</sup> A detailed description of different kinds of microstructures in general monolingual dictionaries can be found in WIEGAND 1989b; 1989c.

graphic comments, are the non-terminal nodes of such a tree representation. Hierarchical microstructures can be described as a context-free dictionary grammar  $DG \langle CEI, CNI, R, WA \rangle$  where:

- *CEI* (the set of *Classes of Elementary Items*) is the terminal alphabet of  $DG$ ;
- *CNI* (the set of *Classes Non-elementary Items*) is the set of non-terminal symbols of  $DG$ ;
- *R* (*Rules*) is a set of context-free rewrite rules;
- *WA* (*Wörterbuchartikel* = dictionary entry) is an element of *CNI* and the initial symbol of  $DG$ .

### 2.2.3. Structure indicators

Structure indicators (*SI*s) carry information about the form of the dictionary and have the genuine function of supporting the user's perception of the macrostructure and the microstructure. Two types of *SI*s must be distinguished:

- Symbols like punctuation marks, brackets or other special characters are used as *nontypographic structure indicators* ( $SI_{nyp}$ ) which give structural information, e.g., separate lexicographic items from each other. Along with the lexicographic items, they constitute the main types of lexicographic text segments.
- The fonts and typefaces in which text segments appear in the dictionary have the function of *typographic structure indicators* ( $SI_{typ}$ ). With respect to the printed dictionary,  $SI_{typ}$  are not segments of the dictionary text, but rather attributes of segments such as items and  $SI_{nyp}$ . They belong – from a semiotic point of view – to a secondary system of signs, which can only fulfill their sign function with the help of a primary system of signs.<sup>8</sup> This, however, does not hold true for the way in which the dictionary text is represented on typesetting tapes: In this representation, typographic information is encoded directly in the form of control codes denoting different kinds of fonts and typefaces.

Together  $SI_{typ}$  and  $SI_{nyp}$  both help human users to reconstruct the textual structure of a dictionary. An elaborate system of *SI*s allows for rapid access to the desired information and can therefore considerably speed up the process of dictionary look-up. *SI*s play a central role in dictionary entry parsing because the automatic segmentation and analyses of dictionary text is mainly controlled by the interpretation of *SI*s. Whereas human users of dictionaries can rely on their linguistic and metalexicographic knowledge, dictionary entry parsers have to cope without this knowledge. The quality of the parsing results are substantially impaired if the system of *SI*s is not carefully considered or if it is used in an inconsistent way.

### 2.3. Dictionary Entry Parsing vs. Sentence Parsing

On a very high level of abstraction, the parsing of syntactic structures of natural language sentences and the parsing of dictionary entry structures are similar tasks, namely the automatic assignment of non-linear structure to a linear input stream of text segments by consulting a grammar. A closer look, however, reveals significant differences between the problems related to sentence parsing and those related to dictionary entry parsing.<sup>9</sup>

- A sentence parser has to cope with recursion, whereas a dictionary entry parser has to adequately handle the iteration of an arbitrary number of items belonging to the same item class, all being direct constituents of the same item group (e.g. a number of example sentences all of which are direct constituents of an example group item or a number of translations belonging to the same translation group). Recursion is only needed for dictionaries with niched or nested entries and only in those cases where the structure of the sub-entries is identical to the structure of the main entry.
- A sentence parser generally consults a lexicon in order to assign the words of the sentence to syntactic categories. A dictionary entry parser, however, when classifying the lexicographic text segments as be-

<sup>8</sup> Cf. WIEGAND 1991:26f.

<sup>9</sup> On this subject see also NEFF, BIRD & RIZK 1988.

longing to specific item classes, has to manage without a lexicon: It is not feasible to list all items of unrestricted item classes, like example or definition sentences in a lexicon, since this would mean a listing of an infinite number of phrases and sentences which can, in principle, be used in order to describe or exemplify the meaning of the lemmata in the dictionary.<sup>10</sup> Lists of permissible values can only be specified for 'closed' item classes (e.g. items giving the gender of a noun or items giving information about inflection) – the parser can then test a lexicographic text segment in question against these lists and assign it to the appropriate item class.

- An important issue of sentence parsing is the recognition and resolution of various types of structural and semantic ambiguities. These phenomena do not have a counterpart in dictionary entry parsing where the hierarchical microstructure is unambiguously specified. However, dictionary entry parsing has to cope with two types of functional ambiguities both related to *SI*s, i.e., those text segments which play a crucial role for the identification and interpretation of the dictionary entry structure:
  - One type of ambiguity occurs if the same *SI* has potentially different functions: The slash in the DUDEN-STILWÖRTERBUCH for example is used, on the one hand, to separate lexical variants from each other (e.g. two variants of a phraseme component as in *den starken Mann mimen/markieren*) and, on the other hand, it indicates the beginning and end of an item giving pragmatic-semantic information (as in *grüner Star /eine Augenkrankheit/*). The dealing with this type of ambiguity is simple as long as the two functions are – linguistically speaking – in complementary distribution, i.e., if each function is bound to a specific context. Error-prone, with respect to dictionary entry parsing, are those cases where contextual features do not suffice to identify the function of the *SI*.
  - Most dictionary texts contain symbols which – depending on the context – can either be a *SI* giving information on the dictionary form or can be part of an item giving information on the dictionary subject.<sup>11</sup> An example of this second type of ambiguity is the colon in the DUDEN-STILWÖRTERBUCH which can either function as a *SI* separating the comment on form from the rest of the dictionary entry, or can be part of an example group, e.g., in *die Mannschaften trennten sich 0:0* (the teams parted with a score of 0:0). In the following we will use the term *indefinite structure indicator* (*SI<sup>ind</sup>*) for the symbols which reveal this type of ambiguity between dictionary form and dictionary subject in order to distinguish them terminologically from symbols called *definite structure indicators* (*SI<sup>def</sup>*), which may also be functionally ambiguous but which always carry information about the dictionary form (like the slash mentioned above).

Dictionary entry parsers have to provide facilities to handle both types of ambiguities, because the automatic identification of the dictionary entry structure is mainly controlled by the correct interpretation of *SI*s. The wrong interpretation of *SI<sup>ind</sup>* will inevitably result in bad segmentation results and in inadequate structural analyses. In the following section we will show that the LEXPARSE systems provides straightforward solutions to these problems based on a differentiated handling of *SI<sup>def</sup>* and *SI<sup>ind</sup>*.

## 2.4. Requirements for a Dictionary Entry Parser

In order to cope with the specific problems related to the automatic conversion of typesetting tapes into lexical databases, a dictionary entry parser must meet the following requirements:

- The segmentation of the typesetting tape into separate dictionary entries and the structural analyses of these entries must be based primarily upon the automatic recognition and the correct interpretation of *SI*s.
- Ambiguities due to *SI<sup>ind</sup>* have to be recognized and handled adequately.
- The grammar formalism supported by the dictionary entry parser has to provide operators and mechanisms to cope with phenomena – such as the handling of fonts and typefaces and the management of counters – that cannot be expressed by means of a context-free grammar formalism.
- Iteration of an arbitrary number of items belonging to the same item class is typical for hierarchical microstructures of dictionary entries and has to be supported by the grammar formalism. Recursive rules are not appropriate, because they result in a misleading representation of iterative structures.

<sup>10</sup> The permissible values of these item classes can only be specified by means of very general data types.

<sup>11</sup> For the distinction between dictionary form and dictionary subject see WIEGAND 1991:18f.

In addition, a dictionary entry parser must meet some general requirements:

- A simple and easy to learn yet powerful grammar formalism should assist the user in developing dictionary entry grammars.
- In order to meet specific user needs, the format of the parse trees should be user-configurable. The system should at least be capable of displaying parse trees as hierarchical attribute-value structures for an easy-to-read view, and as annotated SGML structures for a convenient conversion of the data into other formats.
- The parser should provide special functions and protocol files in order to facilitate the development and the debugging of dictionary entry grammars.
- Since the process of parsing is usually carried out offline, speed of execution plays a minor role for dictionary entry parsing. However, brief response times are an advantage, especially during the time period in which the dictionary entry grammar is being developed and the user heavily interacts with the system.
- The program should be capable of permanently saving all settings and options in order to facilitate subsequent parsing sessions with identical settings and options.
- The implementation of the program should be independent of a specific host architecture or operating system so as to obtain a maximum portability.

In the following we will explain how the LEXPARSE parsing system meets these requirements. In contrast to former approaches to dictionary entry parsing, the LEXPARSE system was not built upon an existing natural language parser but was developed from scratch, taking into account the specific problems related to the parsing of dictionary text stored on typesetting tapes.

### 3. LEXPARSE System Design and Concepts

#### 3.1. The System's Architecture

The LEXPARSE system consists of two files: the configuration file and the program executable comprising the following modules:

- The *main function* initializes all data structures, error handlers and uses default values for the program's settings. It also contains the main loop for dictionary entry parsing.
- The *parsing engine* itself is split into three different components:
  - the preprocessor
  - the scanner
  - the parser
- The *output function* displays the resulting parse tree in different user-definable styles and performs textual conversions to the terminal nodes of the parse tree to meet user-specific needs.

Figure 1 illustrates the architecture of the LEXPARSE system.

The *configuration file* is a pure ASCII text file, maintained by a standard text editor, which contains all settings, commands and the dictionary entry grammar. The configuration file may be split into several files if a separation between general settings controlling the general parse process and dictionary-specific settings (the dictionary entry grammar and control of the scanner and the preprocessor) is desired.

At runtime, LEXPARSE reads and interprets the configuration file, verifies the grammar and checks all settings for completeness. If no errors did occur so far, the program opens the specified input file for parsing.

The *preprocessor* reads the input file on a line-by-line basis and converts tape specific format sequences and codes, i.e., removes superfluous control codes and changes escape sequences into readable characters.

The *scanner* splits the lines being read by the preprocessor into atomic items, the so-called *tokens*. These tokens may represent textual elements (e.g. words, numbers, punctuation

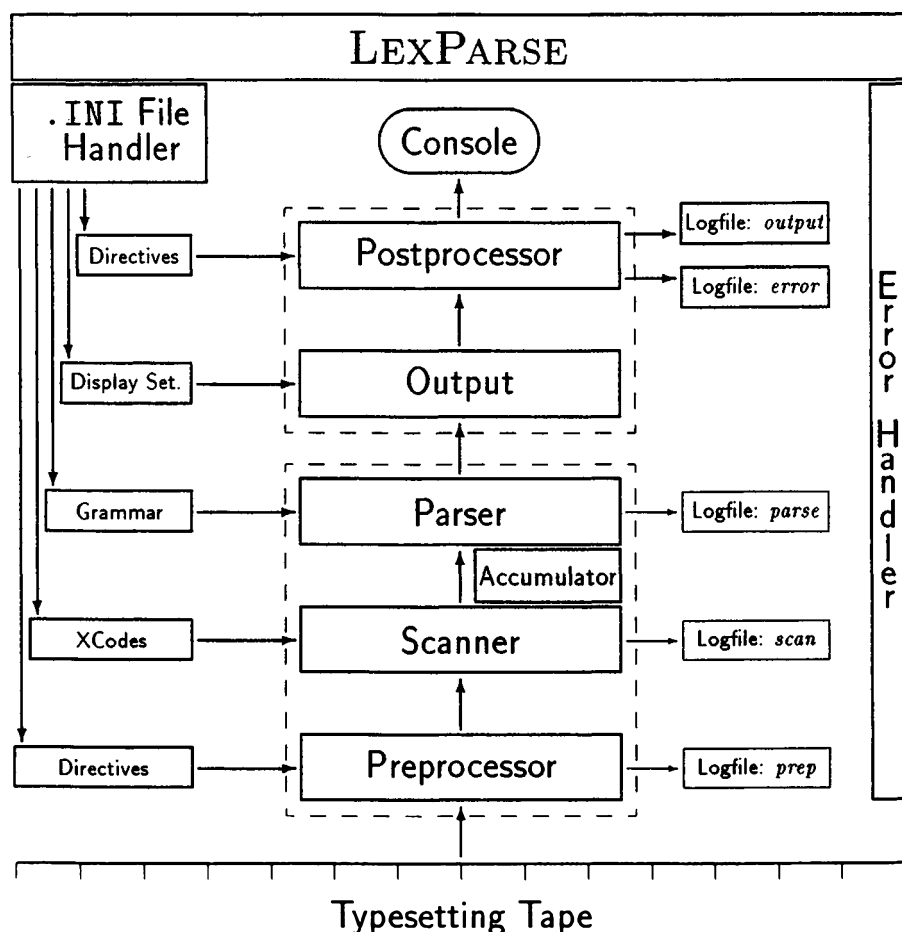


Figure 1: The Design of the LEXPARSE Program

characters) or structure indicators (which may be defined as XCodes; cf. the next section for details).

The *parser* matches these tokens with the rules of the dictionary entry grammar and creates a parse tree by expanding the starting symbol of the grammar.

The *postprocessor* applies textual conversions to the terminal nodes of the parse tree, i.e., the representation of the lexicographic items. For instance German 'Umlaute' may be converted into the respective T<sub>E</sub>X or L<sup>A</sup>T<sub>E</sub>X representation.

The program displays the parse tree according to user-defined style settings. Entries being marked as non-well-formed are displayed with a maximum parse tree so that the source for the error is almost immediately apparent (cf. example D2 in the appendix).

All actions performed by the system are logged in different *protocol files* in order to facilitate debugging. The user can specify whether the protocol shall comprise all entries or only those which were marked as non-well-formed. A final *status report* informs the user about the total number of correct and incorrect entries and the time spent on parsing.



### 3.2. LEXPARSE Concepts

After having applied all preprocessor directives, the scanner splits the input into *tokens*, i.e., atomic strings being enclosed by token delimiters. LEXPARSE distinguishes two classes of tokens, the so-called *XCode tokens* and the so-called *literal tokens*. XCodes are predefined symbols which are generally used to represent *Sl*s. A literal is a string containing a word, a number or a punctuation character.

Tokenization is controlled by the following principles:

- Word delimiters always separate consecutive tokens. Word delimiters are:<sup>12</sup>
  - the space character (ASCII 32),
  - the newline/linefeed characters (ASCII 10 resp. 13),<sup>13</sup>
  - the formfeed character (ASCII 12), and
  - the tabulator character (ASCII 9).<sup>14</sup>
- A change from a letter to a digit and vice versa results in separate tokens.<sup>15</sup>
- A change from a letter to a punctuation character and vice versa results in separate tokens.
- An occurrence of an XCode being introduced by a preprocessor directive results in separate tokens.

The difference between XCodes and literals is crucial to the LEXPARSE approach and will be discussed in more detail in the following sections.

#### 3.2.1. XCode Tokens

XCodes are predefined symbols provided by the LEXPARSE system which are generally used to represent *Sl*s. In the typesetting representation of the dictionary text, *Sl*s are represented either as control codes (this applies to all *Sl*<sub>typ</sub>) or as ordinary characters (this applies to most *Sl*<sub>typ</sub>).

In general, the grammar writer assigns XCodes to all patterns denoting *Sl*<sub>def</sub> and, in specific cases, also to the patterns denoting *Sl*<sub>ind</sub>.<sup>16</sup>

The grammar writer who is in charge of this assignment has to consider the predefined semantics of the XCodes:<sup>17</sup> Punctuation characters must be assigned to the appropriate XCodes (e.g. the colon ':' must be assigned to the XCode XPCOL denoting the punctuation character 'colon' (:)) and control codes must be assigned to their appropriate XCodes (e.g. a control code &234&16 indicating the beginning of a new dictionary entry in the typesetting tape must be assigned to the XCode XFLBE denoting the beginning of dictionary entries).

LEXPARSE offers, aside from user-definable XCodes, a wide variety of predefined XCode label, each of which consists of a combination of exactly five letters and belongs to a specific XCode class:

<sup>12</sup> The configuration option [Scanner] AddLetters can be used to restrict the set of delimiters by expanding the set of characters which are regarded as 'letters'.

<sup>13</sup> The configuration switch [Settings] SkipNLChar determines whether these characters act as delimiters or whether they are simply ignored.

<sup>14</sup> The configuration switch [Settings] TabAsXCode determines whether the tabulator character is regarded as an XCode or whether it is simply a delimiter.

<sup>15</sup> By default the set of characters being considered as 'letters' contains all uppercase letters A-Z and all lowercase letters a-z. This set may be further expanded using the configuration option [Scanner] AddLetters.

<sup>16</sup> A detailed discussion of strategies for the assignment of XCodes to *Sl*<sub>ind</sub> can be found in HAUSER 1993.

<sup>17</sup> This does not apply to XCodes of class XU\*\*\*, i.e., user-definable XCodes.

| XCodes       | Class                               |
|--------------|-------------------------------------|
| <b>XB***</b> | <i>Brackets</i>                     |
| <b>XP***</b> | <i>Punctuation characters</i>       |
| <b>XC***</b> | <i>(Other) Characters</i>           |
| <b>XF***</b> | <i>Format</i>                       |
| <b>XFT**</b> | <i>Format: Typefaces (Subclass)</i> |
| <b>XU***</b> | <i>User-definable</i>               |
| <b>X_EOF</b> | <i>End Of File</i>                  |

A complete list of all predefined XCodes and their semantics can be found in the appendix.

Various operators provided by the LEXPARSE grammar formalism allow for a differentiated treatment of XCodes: The skip operators consume XCode tokens whereas the \$ operators block the consumption of XCode tokens.<sup>18</sup> In section 4.2. we will show how this can be applied to the functional disambiguation of *SIs*<sup>ind</sup>.

### 3.2.2. Literal Tokens

All tokens not assigned to XCodes are regarded as literal tokens or *literals*. These literals can be words, numbers or punctuation characters; a word is defined as being any sequence of 'letters'. The set of characters being regarded as 'letters' contains all uppercase letters A–Z, all lowercase letters a–z, and all characters as defined by the configuration option **[Scanner] AddLetters**.

Any sequence of digits (0-9) is regarded as being a number.

All characters which are neither digits nor letters and which are not defined as XCodes are treated as punctuation characters which always result in single literal tokens.

In contrast to XCodes, literals are consumed by both skip operators and \$ operators.

### 3.2.3. Typeface Changes and Typeface States

Typesetting tapes represent typographic information in the form of control codes which indicate transitions from one typeface to another typeface.<sup>19</sup> In most cases, the end of a specific typeface is not explicitly encoded but rather implicitly indicated by various types of events, e.g., the end of a paragraph, the end of an entry or the beginning of another typeface. The scope of a specific typeface thus reaches from the occurrence of the respective control code up to one of these events.

All control codes indicating typographic information are assigned to XCodes of class **XFT\*\***. LEXPARSE provides a large set of XCodes to cover all possible typefaces and supports two alternative approaches to the handling of typographic information:

- *Operation mode 'TypefaceChanges'*: Within the operation mode 'TypefaceChanges'<sup>20</sup>, XCodes denoting typefaces (XCodes of class **XFT\*\***) are treated like all other XCodes: if an XCode is specified as a terminal in a grammar rule, the parser expects (and consumes) the respective XCode token in the input. The rule fails, if the next token available in the input does not match this XCode terminal.

An example grammar which makes use of the 'TypefaceChanges' mode is given in (1):

<sup>18</sup> For a detailed discussion of these operators refer to section 4.2.

<sup>19</sup> In this paper, we use the term 'typeface' meaning both fonts and typefaces.

<sup>20</sup> This operation mode is the default mode and corresponds to the configuration switch **[Settings] TypefaceStates = Off**.

```
(1) S ->  XFTbo,  A
          |  XFTit,  B
          |  [ XFTst ] C.
```

According to the syntax and semantics of the LEXPARSE grammar formalism to be specified in the following section, this rule is interpreted in the following way: Category **S** is expanded to category **A** if the control code for a typeface change to 'bold' (which is assigned to the XCode **XFTbo**) is available as the next token in the input; category **S** is expanded to category **B** if the control code for a typeface change to 'italics' (XCode **XFTit**) is available as the next token, and category **S** is expanded to category **C** if the control code for a typeface change to 'standard' (XCode **XFTst**) is available as the next token in the input. Since the terminal **XFTst** is optional, **S** may also be expanded to category **C** if no XCode token is available at all. This optionality is needed for cases in which the transition to typeface 'standard' has occurred before the rule is applied, because in these cases the respective XCode token is no longer available in the input.

- However, the optionality of category **C** may yield undesired results: If a transition to typeface 'bold' or 'italics' has occurred before the rule is applied, **S** will mistakenly be expanded to **C**, because the **XFTbo** or **XFTit** tokens are no longer available in the input and an expansion to **A** or **B** is, therefore, prevented.

*Operation mode 'TypefaceStates':* When the operation mode 'TypefaceStates'<sup>21</sup> is set, XCodes denoting typographic information (i.e. all XCodes of the class **XFT\*\***) cannot be used as terminals in the grammar. Instead they are interpreted as conditions controlling whether a rule may be expanded or not.

The parser maintains an internal 'typeface indicator' which stores the current typeface state. Each **XFT\*\*** token occurring in the input is consumed immediately and the value of the 'typeface indicator' is changed accordingly. Thus, the system can retrieve at any time the current value of the 'typeface indicator'.

In order to emphasize the special role of these XCodes, the grammar formalism demands that colons are placed right after all **XFT\*\*** codes. The rule following the colon can be expanded only if the condition (expressed by the **XFT\*\*** code and the colon) is true, i.e., if the 'typeface indicator' is set to the value corresponding to this **XFT\*\*** code. An example is given in (2):

```
(2) A -> XFTst: X.
```

Category **A** may be expanded to category **X** if and only if the 'typeface indicator' has the value 'standard', i.e., if the last XCode token **XFT\*\*** which occurred in the typesetting tape so far was **XFTst**. In contrast to the mode 'TypefaceChanges', it is of no importance whether the 'typeface indicator' has reached this state right before the expansion of category **A** to category **X** or at an earlier time. The crucial thing is that **XFTst** was the last **XFT\*\*** token occurring in the typesetting file, which means that the 'typefaces indicator' has the typeface 'standard' as its current value.

With the help of operation mode 'TypefaceStates', we can now reformulate the example grammar (1) in the following way:

```
(3) S -> XFTbo: A
          | XFTit: B
          | XFTst: C.
```

Category **S** may be expanded to category **A** if the parser has reached typeface state 'bold'; **S** may be expanded to **B** if the parser has reached typeface state 'italics' and **S** may be expanded to **C** if the parser has reached typeface state 'standard'. If the 'typeface indicator' is in any other state, category **S** cannot be expanded by the rules given in (3).

Typeface changes and the scope of typographic information cannot be directly expressed in context-free rewrite rules; this is why specific facilities overcoming these limitations are required. Managing typeface states with a 'typeface indicator' is more suitable than using terminal symbols for typeface changes, because it allows for a deterministic description of typographic information in the dictionary entry grammar. Moreover, this approach reflects the difference between the semiotic status of  $SI_{typ}$  and  $SI_{nryp}$ , as discussed in section 2.2.3.

<sup>21</sup> This operation mode must be explicitly set by the configuration switch [Settings] **TypeFace-States = On**.

### 3.2.4. The Accumulator

Preterminal categories of dictionary entry grammars may be expanded by rules containing terminals and/or operators. Terminals may match with tokens being generated by the parser's scanner; operators are symbols which represent specific patterns of terminals and may match with a single token or a sequence of tokens.

Tokens or sequences of tokens, which match the terminals or operators specified within a rule, are consumed by the parser and are collected in a container object called *accumulator*. Each category has its own accumulator which consecutively 'collects' all tokens that can be matched and consumed from the input while expanding this category. The result of the accumulation process for each category is finally displayed in the parse tree.<sup>22</sup> The management of the accumulator takes care that all tokens are correctly assigned to the accumulator of a category resp. to the accumulators of its daughter categories.

The following example illustrates how the accumulators of the category A and B are managed:

```
(1) A -> $$ ":" B $$ ".".
    B -> "(" $$ ")".
(2) a b c: (d e f) g h i.
(3) [a] [b] [c] [:] [(] [d] [e] [f] [)] [g] [h] [i] [.]
(4) +--> A: "a b c : g h i ."
    +--> B: "( d e f )"
(5) Accumulator of A: "a b c : g h i ."
(6) Accumulator of B: "( d e f )"
```

In (1), category A may be expanded to any sequence of literals (expressed by the \$\$-operator), a colon, category B, a second sequence of literals, and a dot. Category B as in (2) may be expanded into any sequence of literals enclosed by parentheses. The input (2) thus results in the token list (3). Expanding category A with this token list results in the parse tree (4).

By default, literal tokens are added to the accumulator of a category and displayed in the parse tree, whereas XCode tokens are not added to the accumulator and are thus not shown in the parse tree. This reflects the conceptual status of XCodes denoting *S*'s which – in contrast to items – are not elements of the hierarchical microstructure of a dictionary entry. In particular cases, this default setting may be inverted by marking an XCode or a string terminal with the circumflex character (^). Literal tokens, for example, which function as *S*'s<sup>nd</sup> and are thus not assigned to XCodes, may be marked with the circumflex so as to prevent them from being displayed in the parse tree. The effect of applying the circumflex character can be demonstrated by the following example:

```
(7) A1 -> XBRPO * XBRPC.
    A2 -> XBRPO^ * XBRPC^.
(8) (a b c)
(9) [XBRPO] [a] [b] [c] [XBRPC]
(10) +--> A1: "a b c"
    +--> A2: "( a b c )"
```

In (7), both categories A1 and A2 cover a sequence of literal tokens, which is enclosed by the XCode tokens XBRPO (*B*Racket*P*arenthesis*O*pen) and XBRPC (*B*Racket*P*arenthesis*C*lose) denoting a pair of parentheses. The input chain (8) results in the token list (9). As shown in

<sup>22</sup> In the configuration's option [Display] Format the variable \$\$ represents the contents of the category's accumulator and may be used in the format string accordingly.

(10), the resulting parse tree for category A1 only displays the literal tokens, whereas the parse tree for category A2 also contains the parentheses '(' and ')', because the respective XCodes are marked with the circumflex character in the expansion rule for A2. More examples for the application of the circumflex character are given in section 4.3.

## 4. The LEXPARSE Grammar Formalism

### 4.1. Formal Description of the LEXPARSE Grammar

In essence, the LEXPARSE grammar formalism uses simple context-free production rules supplemented by additional features for the handling of optionality, alternation and multiple repetition of rules as well as for the specification of closed sets of terminals. In addition, the formalism allows parser directives to be built into the grammar, e.g., to reset counters or to designate a dictionary entry as being incorrect.

#### 4.1.1. Formal Syntactic Description

In the following, we describe the LexParse grammar formalism using an EBNF (*Extended Backus Naur Form*)<sup>23</sup> notation:

```

Category      = CatSpec "->" [ Exceptions ":" ]
               [ SimpleRule "->" ] Rules ".".

Exceptions    = Exception { "," Exception }.

Exception     = "-" XCode
               | "+" XCode.

Rules         = Rule { "|" Rule }.

Rule          = SimpleRule [ "?" TermSpec ] [ "@" ].

SimpleRule    = RuleNode { "," RuleNode }.

RuleNode      = CatSpec
               | TermSpec
               | "[" SimpleRule "]"
               | "<" SimpleRule ">"
               | TextOp
               | Counter
               | "*" TermSpec
               | "*" TermSpec
               | "\" TermSpec
               | .

TermSpec      = Terminal
               | "{" TerminalSet "}".

TerminalSet   = Terminal { "," Terminal }.

TextOp        = "$"
               | "$$"
               | "#".

```

<sup>23</sup> This Extended BNF uses the following additional meta characters: The definition character "=" replaces the BNF definition characters ":", "=". Square brackets "[" indicate optionality of an item whereas curly brackets "{}" indicate the possible repetition of an item. A dot terminates each rule. Cf. JENSEN & WIRTH 1985.

```

Counter      = "%" [ "R" ] CounterType.
CounterType  = "I"
              | "i"
              | "A"
              | "a"
              | "1".

CatSpec      = Letter { Char }.

Char         = Letter
              | Digit
              | "_".

String       = SLetter { SLetter }
              | Digit { Digit }.
              | PunctuationChar
              | RegularExpression.

```

The term **Letter** denotes an element out of the set of letters **A** to **Z** and **a** to **z**.

The term **Digit** denotes an element out of the set of digits **0** to **9**.

The term **SLetter** denotes an element out of the set of letters **A** to **Z**, **a** to **z** and all characters as defined by the configuration option `[Scanner] AddLetters`.

The term **PunctuationChar** denotes any other character not being a **SLetter** or a **Digit**.

The term **RegularExpression** denotes a regular expression **rep** which must be of the form **^rep\$**. All special characters being part of a regular expression are described together with their semantics below.

Within the mode 'TypefaceChanges' an additional rule is defined:

```

Terminal     = XCode [ "^" ]
              | "" String "" [ "^" ].

```

The term **XCode** denotes an element out of the set of **XCodes**. The complete list of all **XCodes** and their semantics can be found in the appendix.

Within the mode 'TypefaceStates' an additional rule is defined:

```

Terminal     = XCode [ "^" ]
              | XCodeFT ":"
              | XCodeFT "!"
              | "" String "" [ "^" ].

```

The term **XCode** denotes an **XCode** of any class but class **XFT\*\***. The term **XCodeFT** denotes an **XCode** of class **XFT\*\***.

#### 4.1.2. Semantic Description

*Note:* All parser directives are  $\epsilon$  productions, that is, they do not consume any tokens from the input.

*XCode-Exceptions:* All **XCODE-Exceptions** are parser directives.

```
C -> -XCODE: R.
```

Within the scope of category **C**, the **XCode** **XCODE** is considered to be a literal and not an **XCode**.

```
C -> +XCODE: R.
```

Within the scope of category C, the XCode XCODE is considered to be an XCode and not a literal.

*Note:* XCode exceptions may be nested up to any level.

*Rule Abbreviation:*

$$C \rightarrow X \rightarrow R1 \mid R2.$$

This rule will be replaced by the rule

$$C \rightarrow X, R1 \mid X, R2.$$

which means that both rules are equivalent.

*Generals:*

$$C \rightarrow R1 \mid R2.$$

If the expansion of category C to R1 fails, C may be expanded to R2.

$$C \rightarrow [ R1 ] R2.$$

Category C may be expanded to R1 R2 or solely to R2. (R1 is optional and *may* be expanded at this position once.)

$$C \rightarrow < R1 > R2.$$

Category C may be expanded to multiple instances of R1 followed by R2 or solely to R2. (R1 is optional and may be expanded at this position several times.)

$$C \rightarrow \{ t1, t2, t3 \}.$$

Category C expands to *one* of the three terminals t1, t2 or t3.

$$C \rightarrow * t.$$

Category C expands to a sequence of tokens up to and including a specific token which matches the terminal t.

$$C \rightarrow * \backslash t.$$

Category C expands to a sequence of tokens up to a specific token which matches the terminal t. However, this matching token is not consumed by the parser and added to the accumulator of category C but remains in the input.<sup>24</sup>

$$C \rightarrow \backslash t, R.$$

If the next token available in the input does not match the specified terminal t, category C expands to R. The token matching the specified terminal t is not removed from the input.

$$C \rightarrow \$.$$

Category C expands to a single literal token.

$$C \rightarrow \$\$.$$

Category C expands to a sequence of (at least one) literal tokens.

$$C \rightarrow \#.$$

<sup>24</sup> The LEXPARSE rule  $C \rightarrow * \backslash t.$  with  $C \in N, t \in T$  corresponds to the context-sensitive rule  $Ct \rightarrow \alpha t$  with  $C \in N, t \in T$  where  $\alpha$  covers any sequence of tokens except the token matching t.

Category C expands to a single literal token which must be a number, i.e., a sequence of digits.

$C \rightarrow \text{"string"}$ .

Category C expands to a literal token which matches the quoted **string**.

*Note:* The **string** is considered to be a regular expression **rep** only, if it is specified in the form  $\wedge \text{rep} \$$ .

$C \rightarrow R ? t$ .

Category C expands to R if and only if C could be expanded to R t.<sup>25</sup>

*Counters:*

$C \rightarrow \%I$ .

Category C expands to a single literal if this literal matches an uppercase roman number (i.e., I, II, III, IV ...).

$C \rightarrow \%i$ .

Category C expands to a single literal if this literal matches a lowercase roman number (i.e., i, ii, iii, iv ...).

$C \rightarrow \%A$ .

Category C expands to a single literal if this literal matches an uppercase latin letter (i.e., A, B, C, D ...).

$C \rightarrow \%a$ .

Category C expands to a single literal if this literal matches a lowercase latin letter (i.e., a, b, c, d ...).

$C \rightarrow \%1$ .

Category C expands to a single literal if this literal matches an arabic number (i.e., 1, 2, 3, 4 ...).

*Note:* If a specific counter symbol could be derived successfully, the corresponding counter is incremented.

*Resetting the Counters:*

$C \rightarrow \%RI, R$ .

The counter for uppercase roman numbers is reset to I before category C is expanded to R.

$C \rightarrow \%Ri, R$ .

The counter for lowercase roman numbers is reset to i before category C is expanded to R.

$C \rightarrow \%RA, R$ .

The counter for uppercase latin letters is reset to A before category C is expanded to R.

$C \rightarrow \%Ra, R$ .

The counter for lowercase latin letters is reset to a before category C is expanded to R.

<sup>25</sup> The LEXPARSE rule  $C \rightarrow R ? t$ , with  $C \in N, R \in \Sigma^*, t \in T$  corresponds to the contextsensitive rule  $Ct \rightarrow Rt$  with  $C \in N, R \in \Sigma^*, R \neq \varepsilon, t \in T$ .



**C** -> %R1, R.

The counter for arabic numbers is reset to 1 before category C is expanded to R.

*Note:* All directives resetting a counter are parser directives.

*Typeface States:*

**C** -> XFT\*\* : R.

Category C is expanded to R if and only if the internal indicator for typeface states is in the state of XFT\*\*.

**C** -> XFT\*\* ! R.

Before category C is expanded to R, the internal indicator for typeface states is set to the specified value XFT\*\*.

*Note:* All directives handling typeface states are parser directives.

*Additional Parser Directives:*

**C** -> t^.

Category C expands to a token matching the specified terminal t. If this token is an XCode token and it is marked as indicated, it will be added to the accumulator of category C. By default, all XCode tokens are not added to the accumulator. If this token is a literal token and it is marked as indicated, it will not be copied to the accumulator. By default, all literal tokens are added to the accumulator.

**C** -> R @.

Category C expands to R. After C has been expanded, the current entry is marked as being non-wellformed and backtracking is terminated for the parsing of this entry.

#### 4.1.3. Regular Expressions in LexParse

The LEXPARSE program supports the following characters as being part of a regular expression.<sup>26</sup>

|                  |   |
|------------------|---|
| <b>^</b>         | This expression matches the beginning of the source string.   |
| <b>\$</b>        | This expression matches the end of the source string.   |
| <b>.</b>         | This expression matches any character.  |
| <b>[abcA-Z]</b>  | This expression matches a set consisting of the characters a, b, c and the sequence of characters A to Z.   |
| <b>[^abcA-Z]</b> | This expression matches any character except the set consisting of the characters a, b, c, and the sequence of characters A to Z.                             |
| <b>\c</b>        | This expression matches the character c. This may be used to override any special meaning given to a character.   |
| <b>*</b>         | If one of the expressions as given above is followed by this character, then this expression matches a sequence of zero or more instances of this expression. |

*Note:* A regular expression **rep** as being used within the specification of terminals for the

<sup>26</sup> The term 'source string' denotes a line as being read from the input data file as far as the regular expression is used in the context of a preprocessor or a postprocessor directive. The term 'source string' denotes a literal token as far as the regular expression is used within the specification of string terminals for the grammar.

grammar must always be specified in the form  $\wedge\text{rep}\$$  so that it can be distinguished from other strings.

## 4.2. Annotations to the Formalism

The LEXPARSE grammar formalism uses several features in writing dictionary entry grammars and in handling the special requirements of dictionary entry parsing. In this section, these features will be explained in detail.

### 4.2.1. Single and Multiple Optional Rules

The grammars in (1) and (2) generate the same language.

```
(1) A    -> [ X ].
(2) A    -> Xaux.
    Xaux -> X
        | .
```

However, the resulting parse trees differ in that the parse tree in (2) contains an additional node *Xaux*, whereas in (1) category *A* expands directly to category *X*. Optional rules, thus, simplify the grammar and lead to a representation of the dictionary entry structure, which is more adequate because the semantically empty category *Xaux* is not needed.

The grammars in (3) and (4) also generate the same language.

```
(3) B    -> < X >.
(4) B    -> Xaux.
    Xaux -> X, Xaux
        | .
```

However, the resulting parse trees differ in that the parse tree in (4) contains additional nodes *Xaux*, whereas in (3) category *B* expands directly to category *X*. Due to the right-recursion used in (4), the derivation of category *B* results in an asymmetric parse tree with a depth of  $n$ , whereas the expansion of category *B* in (3) results in a symmetric parse tree with a uniform depth of 2 and  $n$  nodes being immediate constituents of *B*. This greatly simplifies the grammar and leads to a representation of the dictionary entry structure, in which no auxiliary category *Xaux* is required. The parse trees (5) and (6), both obtained by deriving  $X\ X\ X$  from category *B*, may illustrate the difference: The parse tree in (5) corresponds to the grammar given in (4); the parse tree in (6) corresponds to the grammar given in (3).

|   |  |
|---|--|
| <pre>(5) +--&gt; B       +--&gt; Xaux           +--&gt; X               +--&gt; Xaux                   +--&gt; X                       +--&gt; Xaux                           +--&gt; X</pre> | <pre>(6) +--&gt; B       +--&gt; X           +--&gt; X               +--&gt; X</pre> |
|---|--|

Another example is the category *Number* which can be defined by a pure context-free grammar, as in (7), or – much more easily – by using the LEXPARSE formalism, as in (8).<sup>27</sup>

```
(7) Number -> Sign, Digit, Digits.
    Sign    -> "-"
           | .
```

<sup>27</sup> The category *Digit* will be defined in the next section.

```

    Digits -> Digit, Digits
           | .
(8) Number -> [ "-" ] Digit < Digit >.

```

#### 4.2.2. Using a Set of Terminals

The grammars in (1) and (2) generate the same language:

```

(1) A    -> X { t1, t2, t3 }.
(2) A    -> X, Tset.
    Tset -> t1
           | t2
           | t3.

```

Set descriptions as used in (1) lead to an additional simplification of the grammars and to a representation of the dictionary entry structure, which is more adequate because additional categories like *Tset* are avoided. In (1), category *A* may be directly expanded to one of those terminals which are a member of the respective set.

Set descriptions offer the possibility to specify finite sets with a given domain, as used for the definition of the notion *digit* in (3):

```

(3) Digit -> { "0", "1", "2", "3", "4", "5", "6", "7", "8", "9" }.

```

Using a regular expression for the terminal string, one could also replace rule (3) by rule (4):

```

(4) Digit -> "[0-9]".

```

In the context of dictionary entry parsing, a distinction between numbers and words is quite helpful. Therefore, the LEXPARSE grammar formalism offers a special operator to consume numerical tokens, which makes the derivation rule for category *Digits* again more simple. The *#*-operator in (5) corresponds to the regular expression string in (6).

```

(5) Digits -> #.
(6) Digits -> "[0-9]*".

```

Therefore, the rule in example (8) of the previous section can be replaced by the following rule:

```

(7) Number -> [ "-" ] #.

```

#### 4.2.3. Using the 'Skip' Operators

The skip-operators enable the derivation of sequences of tokens up to a specific 'delimiter' token. 'Delimiter' tokens generally function as *Sl*s marking the boundary of a lexicographic item and can be both, XCode or literal tokens. Two types of skip-operators are provided:

- The inclusive skip-operator (*\**) reads up to and including the 'delimiter' token.
- The exclusive skip-operator (*\*\*) reads up to the 'delimiter' token without removing it from the input. In contrast to the inclusive skip-operator, the 'delimiter' token remains in the input and may be consumed by another operator.

The inclusive skip-operator is useful for error handling routines – as in (1) – or for categories covering any sequence of tokens enclosed by two tokens representing *Sl*<sup>def</sup>, as in (2).

```

(1) A    -> X1, X2, X3 ":"
           | * ":".
(2) B    -> "(" * ")" .

```

In (1), category A may be expanded to a sequence of categories X1 X2 X3, and a trailing colon. If this expansion fails, an alternate rule can be applied: A expands to any sequence of tokens (expressed by the inclusive skip-operator \*) up to the colon which terminates the consumption. This enables the expansion of A although the constituents X1, X2 and X3 could not be derived from A. In this case, the alternate derivation rule for A serves as an error handler rule.

In (2), category B expands to a left parenthesis and the inclusive skip-operator which consumes all tokens up to the next literal token matching the right parenthesis.

The inclusive skip-operator is quite useful for the step-by-step development of dictionary entry grammars: the grammar writer can rough in the dictionary entry structure and deal with the structural details of complex item classes at a later stage.

The exclusive skip-operator leads to the expansion of a category into an arbitrary sequence of tokens up to a specific 'delimiter' token, which is not included by the expansion but indicates the next token of the input stream.

```
(4) S      -> A, X1
          |  A, X2
          |  A, X3.
(5) A      -> *\ { "(", "<", "[" }.
(6) X1     -> "(" * ")".
      X2     -> "<" * ">".
      X3     -> "[" * "]".
```

In (4), category S may be expanded to category A followed by either category X1, X2, or X3. X1, X2, and X3 can be distinguished by the delimiter tokens specified in (6). If the exclusive skip-operator is used, category A covers any sequence of tokens up to one of the delimiter tokens (defined in the set under (5)) which terminate the consumption of this sequence. However, the delimiter token is not consumed by the parser but remains in the input. In the consecutive expansion of either X1, X2, or X3 this specific token will be consumed and either X1, X2, or X3 may be further expanded.

*Note:* If an XCode token `X_EOF` (*EndOfFile*) occurs during a skip operation, the current entry is marked as being erroneous. The skip operation terminates, the current rule fails, and the parser issues an error message (error #404: **Critical XCode detected while skipping**).

*Note:* If an XCode token `XFLEN` (*FormatLemmaEnd*) occurs during a skip operation, the skip operation terminates and the current rule fails. In this case, the configuration switch `[Parser] SkipCriticalIsError` determines whether the current entry is marked as being erroneous and whether an error message will be issued or not.

#### 4.2.4. Using the 'Except' Operator

The except-operator (\) – also called the accept-all-except-operator – offers the possibility to specify a condition on which rules may be applied as an exception. By naming a terminal or a set of terminals, a rule may be expanded only if a token is not available as the next token of the input.

```
(1) A      -> ":" X1
          |  X2.
(2) B      -> \ ":" X2
          |  X1.
```

The grammars in (1) and (2) generate the same language in that they both expand a cate-

gory A to a category X1 if the next token available in the input is the literal ':'. If the next token available in the input is *not* the literal ':' (but any other literal token), category A is expanded to category X2.

But there is a semantic difference between the grammars in (1) and in (2): If the next token available in the input were the literal ':', grammar (1) would derive this token directly from category A, and the literal ':' would be added to the accumulator of category A. The grammar in (2), by contrast, would derive the literal ':' from category X1, because the first rule of category B in (2) cannot be expanded with this token being available in the input.<sup>28</sup>

#### 4.2.5. Using the 'Test' Operator

The test-operator (?) can be used to test whether the next token available in the input is of a specific token. If the test is positive, i.e., the next token available is of the specified type, it is not consumed by the parser but remains in the input until it can be consumed by another operator.

(1) A      -> XBRPO \* XBRPC ? "!".

In (1), category A expands to any sequence of tokens enclosed by the XCodes XBRPO (*BRacketParenthesisOpen*) and XBRPC (*BRacketParenthesisClose*) only, if the next token following this sequence is the literal '!'. The parser does not consume this token; it remains in the input and can be consumed by another operator.

Note that the test-operator can occur only at the end of a rule. Only the Error-directive may be placed right after the test-operator.

#### 4.2.6. Using the 'Error' Directive

The error-directive (@) offers the possibility to construct error handler routines which mark dictionary entries as being incorrect (non-wellformed) even if they were completely derived.

An error handler routine is a category with a single derivation rule expanding to the skip-operator and the error-directive. Error handler routines enable the parser to generate maximum parse trees.

(1) A      -> X1, X2, X3, XPCOL  
          | \* XPCOL, @.

In (1), category A may be expanded to a sequence of categories X1, X2, X3, and a trailing colon which is assigned to the XCode XPCOL (*PunctuationCOLon*). If this expansion fails, an alternate rule may be applied: A expands to the skip-operator (\*) and the colon (XPCOL). If this expansion succeeds, the error-directive is applied to the current entry, i.e., the current entry is marked as being incorrect.

This construction offers the possibility to continue the parsing of the entry, even when it is regarded as being non-wellformed because the categories X1, X2, and X3 could not be derived from A.

Marking an entry as being non-wellformed is also necessary if the derivation of a specific rule indicates an invalid input.

(2) B      -> XPCOL, X  
          | XPSEM, X, @.

In (2), category B may be expanded to category X if the required XCode token XPCOL is

<sup>28</sup> Assuming – of course – that X1 can expand a ':' token at all.

available as the next token of the input. However, if the next token being available is the invalid XCode token **XPSEM**, then **B** also expands to **X** but the current entry must be marked as being incorrect.

#### 4.2.7. Using Counters

Dictionary entries describing homonymous and/or polysemous lemma signs generally use different types of counters to identify different levels of homonym and sense distinctions. Each of these levels correspond to a different degree of semantic proximity.

If homonym/sense counters are treated as 'normal' categories within a context-free grammar, they cannot be incremented after being derived and, therefore, cannot be checked for correctness.

For this reason, LEXPARSE offers several built-in counters each of which can be incremented and reset automatically or manually by using a directive in the grammar. In example (1), we assume a three level sense distinction, in which roman numbers are used at the first level, arabic numbers at the second, and lowercase letters at the third level.

- ```
(1) I.
    1.
      a)
      b)
    2.
      a)
      b)
  II.
    1.
    2.
```

During the parsing process, LEXPARSE must set the counters to their initial value each time a deeper sense level is reached. The procedure is illustrated in (2):

- ```
(2) I.  {Reset the counter: arabian numbers}
    1.  {Reset the counter: lowercase letters}
        a)
        b)
    2.  {Reset the counter: lowercase letters}
        a)
        b)
  II. {Reset the counter: arabian numbers}
    1. {Reset the counter: lowercase letters}
    2. {Reset the counter: lowercase letters}
```

The three-level system of sense counters in (1) can be described using the LEXPARSE grammar formalism by the rules given in (3).<sup>29</sup>

- ```
(3) NUM -> %I, ".", %R1 [ NUM ]
      | %1, ".", %Ra [ NUM ]
      | %a, ")" .
```

#### 4.2.8. Using the '\$\$' Operator

The \$\$-operator consecutively consumes literals and adds them to the category's accumulator.

<sup>29</sup> For reasons of clarity, this example is kept simple. It is incomplete since category **NUM** would also accept invalid input patterns such as 'I. II. 1. 2. 3. a) b)'.

The \$\$-operator does not consume any XCode tokens, i.e., if an XCode token occurs in the input it stops the consumption of tokens.

Input (1) thus results in the token list (2). The consumption of this token list by the grammar in (3) results in (4):<sup>30</sup>

```
(1) (This-and more- is text!)
(2) [XBRPO] [This] [-] [and] [more] [-] [is] [text] [!] [XBRPC]
(3) A    -> XBRPO $$ XBRPC.
(4) +--> A: "This - and more - is text !"
```

As illustrated in (5) and (6), the \$\$-operator may be followed by a string constant. In this case, all literals up to and including this string constant are consumed by the parser.

```
(5) B1    -> $$ ";".
(6) B2    -> $$ ".".
(7) example1; example2; example3.
(8) [example1] [;] [example2] [;] [example3] [.]
(9) +--> B1: "example1 ;"
(10)+--> B2: "example1 ; example2 ; example3 ."
```

The input (7) results in the token list (8). Given this input, the expansion of B1 results in (9), whereas the expansion of B2 results in (10).

Note that there is no upper limit on the number of literals that may be consumed by the \$\$-operator, but at least one literal is required for its successful application.

#### 4.2.9. Using the '\$' Operator

The \$-operator consumes single literals and may therefore be used to derive a specific number of literals from a category.

```
(1) A    -> $, $$.
```

In (1), category A can be expanded only if at least two literals are available in the input. The first literal is consumed by the \$-operator, the next and all consecutive literal tokens are consumed by the \$\$-operator.

#### 4.2.10. Using XCode Exceptions

XCode exceptions are required when a  $SI^{nd}$ , assigned to an XCode, does not function as  $SI$  in a specific and limited context. Within this context, XCode exceptions can be used to disable XCodes with the effect that the scanner, instead of generating XCode tokens, generates one (or more) of those literal tokens which are denoted by the respective XCode. XCode exceptions may be nested without any restrictions, i.e., an XCode being disabled by a specific category may be enabled again by a nested daughter category.

```
(1) S    -> X1, A, X2 [ B ].
    A    -> -XPCOL: [ B ] $$ .
    B    -> +XPCOL: "(" * XPCOL * ")" .
```

Assuming that the colon is defined as the XCode XPCOL, the derivation rule for category A in (1), disables XPCOL, i.e., in the scope of category A, the colon will not be generated as the XCode token XPCOL but as the literal token ':'.

<sup>30</sup> As already mentioned in section 3.2.4., XCode tokens are by default not added to the accumulator and, thus, do not appear in the resulting parse tree in (4).

Within the scope of category **B** (which is an immediate daughter of category **A**), the XCode **XPCOL** is treated in the usual way: when the token **XPCOL** occurs in the input, the consumption of tokens for the inclusive skip-operator is terminated because the colon functions as a  $SI_{ntyp}$  within this context.

XCode exceptions can only be specified for XCodes of the type **XB\*\*\*** (brackets), **XP\*\*\*** (punctuation characters) or **XC\*\*\*** (other characters), i.e., defined XCodes the denotation of which is known to the scanner. Disabling or enabling undefined XCodes does not change the parser's behaviour in any way.

### 4.3. Applying the LEXPARSE Grammar Formalism

In the following, we will show how the LEXPARSE grammar formalism can be applied. The examples are taken from the LEXPARSE grammar developed for the DUDEN-STILWÖRTERBUCH (cf. ENGELKE 1994).

With respect to the DUDEN-STILWÖRTERBUCH, the symbol 'semicolon' is considered to be a  $SI_{ntyp}^{def}$  and is therefore assigned to the corresponding XCode **XPSEM** (*PunctuationSEMicolon*). The symbol 'colon' is an  $SI_{ntyp}^{ind}$  which, in the majority of its occurrences, functions as *SI* and is, in consequence, also assigned to its corresponding XCode **XPCOL** (*PunctuationCOLon*).

#### 4.3.1. Applying Optional and Multiple Optional Rules

The application of (multiple) optional rules can be illustrated by the derivation rule for category **BeiGA** (*Beispielgruppenangabe*):

**BeiGA** -> [ **UGrA** ] **BeiA** < **XPSEM** [ **UGrA** ] **BeiA** > [ "."<sup>def</sup> ] .

**BeiGA** expands to an optional instance of category **UGrA** (*Unspezifizierte Grammatische Angabe*) followed by one instance of category **BeiA** (*Beispielangabe*). This sequence can be repeated any number of times; the sequences are separated from each other by the  $SI_{ntyp}^{def}$  'semicolon' (**XPSEM**) the series is terminated by an optional dot. The series of sequences of **UGrA BeiA** is expressed using the symbols for multiple optionality < >.

With respect to the DUDEN-STILWÖRTERBUCH, the symbol 'dot' is an  $SI_{ntyp}^{ind}$  which, in most cases, is part of a lexicographic item and not a  $SI_{ntyp}$ . In consequence, the dot is not assigned to an XCode, but treated as a string terminal. In the scope of the category **BeiGA**, however, where the dot actually functions as  $SI_{ntyp}$ , the token representing the dot must be suppressed in the resulting parse tree by means of the '^' operator.

In the dictionary entry for the lemma **treiben**, the category **BeiGA** is expanded to category **BeiA**, category **UGrA**, and another **BeiA**:

```
+--> BeiGA
+--> BeiA "ich ließ mich von den Verhältnissen treiben"
+--> UGrA "jmdm. etwas treiben; mit Raumangabe"
+--> BeiA "der Sturm trieb mir den Schnee ins Gesicht"
```

A series of six **BeiA** derived from **BeiGA** can be found in the entry for the lemma **Treffen**:

```
+--> BeiGA
+--> BeiA "regelmäßige, seltene Treffen"
+--> BeiA "ein Treffen der Abiturienten"
+--> BeiA "ein Treffen der Außenminister"
+--> BeiA "ein Treffen verabreden, veranstalten"
+--> BeiA "an einem Treffen teilnehmen"
```



+--> BeiA "zu einem Treffen kommen"

#### 4.3.2. Applying a Set of Terminals

Sets of terminals are useful for item classes having a closed set of possible values. In the DUDEN-STILWÖRTERBUCH, the gender of German nouns is specified by means of the definite articles *der*, *die*, and *das*. Therefore, the rule for category GA (*Genusangabe*) which denotes a class of items giving the gender of nouns contains the following set description:

GA -> XFTst: { "der", "die", "das" }.

Category GA (*Genusangabe*) expands to one of the string constants specified in the set description. Note that GA can be expanded only if the current value of the typeface indicator is in the typeface state 'standard'. This condition is expressed by the typeface condition XFTst:

#### 4.3.3. Applying the 'Skip' Operators

Category UGrA (*Unspezifizierte Grammatische Angabe*) occurs within the scope of category BeiGA (*Beispielgruppenangabe*) and denotes a class of items which are enclosed by angle brackets. Angle brackets are  $Sl_{ntyp}^{def}$  with regard to the DUDEN-STILWÖRTERBUCH and are therefore assigned to the XCodes XBRAO (*BRacketAngleOpen*) and XBRAC (*BRacketAngleClose*) respectively. Brackets generally appear in pairs enclosing items of specific item classes – in our example the class of items giving information on grammatical properties. The derivation rule for category UGrA may thus be formulated as follows:

UGrA -> XBRAO \* XBRAC.

The *inclusive skip*-operator consumes all tokens up to and including the XCode token XBRAC; as a consequence, UGrA expands to any sequence of tokens beginning with an XCode token XBRAO and ending with an XCode token XBRAC.

An application of the exclusive skip-operator will be shown in the following section and in section 4.3.5.

#### 4.3.4. Applying the 'Test' Operator

In the dictionary grammar for the DUDEN-STILWÖRTERBUCH, the test-operator was used in the derivation rule for category BeiA (*Beispielangabe*), which may be expanded to a sequence of tokens and the optional category SAA (*Sprechaktangabe*) enclosed by two slashes. Examples for items belonging to the class BeiA containing 'embedded' items belonging to the class SAA are given in (1) and (2):

- (1) 'du kriegst die Tür nicht zu! /ugs.; Ausruf des Erstaunens/,'
- (2) 'meiner Treu! /veraltet; Beteuerungsformel/,'
- (3) 'an die Tür/an der Tür klopfen;'
- (4) 'er tat so, als wäre nichts gewesen, als ob/als wenn/wie wenn er nichts wüßte, als wüßte er nichts;'

Examples (3) and (4) show, that slashes do not always enclose SAA items; they are also used to separate two lexical or morphological variants within an example sentence. But SAA items can automatically be distinguished from lexical/morphological variants: category SAA is always enclosed by two slashes and can occur only at the end of an example sentence, which means that a semicolon must follow immediately. This semicolon may not be covered by category SAA itself, but rather by its mother category BeiGA; a constellation which cannot be expressed by means of context-free rules. This is where the test-operator comes into play.

- (5) (a) **BeiA**  $\rightarrow$  \* \ "/" **SAA**  
 (b) | \* \ **XPSEM**.  
 (c) **SAA**  $\rightarrow$  "/"<sup>^</sup> \$\$ **XPSEM**<sup>^</sup> \$\$ "/"<sup>^</sup> ? **XPSEM**.

In (5a), **BeiA** may be expanded to any sequence of tokens up to the literal '/' possibly indicating category **SAA** (acting as  $S_{nyp}^{ind}$ , the symbol 'slash' is not assigned to an XCode). However, due to the exclusive skip-operator, the literal '/' is not derived from category **BeiA**, but from category **SAA**.

In (5c), category **SAA** expands to the literal '/' followed by a sequence of literals, a semicolon (**XPSEM**), more literals, and a terminating literal '/'. Both instances of the symbol 'slash' should be marked with the '^' character so as to prevent the corresponding literal tokens from being displayed in the resulting parse tree. The test-operator makes sure, that **SAA** is expanded if and only if the next token which is available in the input after the expansion of **SAA** is an XCode token **XPSEM**. However, this XCode token is not added to the accumulator for category **SAA** but remains in the input.

If the application of the rules in (5a) and (5c) fails, **BeiA** expands to a sequence of tokens up to (but not including) the XCode token **XPSEM**, as in (5b).

Given the grammar in (5), the examples (1) to (4) result in the structures (6) to (9):

- (6)  $\rightarrow$  **BeiA** : "du kriegst die Tür nicht zu!"  
 $\rightarrow$  **SAA** : "ugs.; Ausruf des Erstaunens"  
 (7)  $\rightarrow$  **BeiA** : "meiner Treu!"  
 $\rightarrow$  **SAA** : "veraltet; Beteuerungsformel"  
 (8)  $\rightarrow$  **BeiA** : "an die Tür/an der Tür klopfen"  
 (9)  $\rightarrow$  **BeiA** : "er tat so, als wäre nichts gewesen, als ob/als  
 wenn/wie wenn er nichts wüßte, als wüßte er nichts"

#### 4.3.5. Applying the 'Error' Directive

In the typesetting tape for the DUDEN-STILWÖRTERBUCH, the beginning and the end of dictionary entries is indicated by specific control codes. These control codes can be used to specify error recovery rules which help the parser in locating the beginning of the next entry when an error has occurred.

- (1) **WA**  $\rightarrow$  **XFLBE**, **FK**, **SK** [ **XPAST**, **PKP** ] **XFLEN**  
 | **WA\_Err**.  
 (2) **WA\_Err**  $\rightarrow$  \* **XFLEN** @.

Category **WA** (*Wörterbuchartikel*), representing dictionary entries, is expanded to an XCode indicating the beginning of an entry (**XFLBE**), the category **FK** (*Formkommentar*), the category **SK** (*Semantischer Kommentar*), the optional category **PKP** (*Postkommentar zur Phraseologie*) being preceded by an asterisk ('\*'; **XPAST**), and an XCode indicating the end of the dictionary entry (**XFLEN**).

If the derivation of the first rule in (1) fails, an alternate rule is applied: In this case, **WA** expands to a category **WA\_Err** which itself expands to the inclusive skip-operator, the XCode indicating the end of the dictionary entry (**XFLEN**) and the error-directive (2). The use of the inclusive skip-operator together with the XCode **XFLEN** leads to the consumption of all tokens up to and including the XCode indicating the end of the entry, i.e., the consumption of the complete current entry up to the beginning of the next entry. The inclusive skip-operator does not assign any structure to this sequence but collects the flat stream of tokens in the accumulator of category **WA\_Err**.

The error-directive, as being part of the derivation rule for category **WA\_Err**, has two

consequences:

- The current entry is marked as being non-wellformed and an appropriate error message is generated by the parser.
- Backtracking is rejected, i.e., all categories derived so far are displayed as a maximum parse tree.

Another error recovery rule can be built for category PKP (*Postkommentar zur Phraseologie*) of the DUDEN-STILWÖRTERBUCH by using the *exclusive skip*-operator:

```
(3) PKP  -> X1, X2, X3
      |   PKP_Err.
(4) PKP_Err -> *\ XFLEN @.
```

If none of the daughter categories X1, X2, or X3 of category PKP can be derived, PKP expands to an error handling category PKP\_Err which itself expands to any sequence of tokens up to the XCode indicating the end of the dictionary entry XFLEN, as in (4). However, this XCode token XFLEN is *not* consumed by the parser and derived from category PKP\_Err. Instead, it remains in the input and may be consumed by another operator, e.g., the one to be found in the derivation rule for category WA in (1). The XCode terminal XFLEN being part of this derivation rule matches with the XCode token XFLEN and may then be consumed by the parser.

#### 4.3.6. Applying Counters

The DUDEN-STILWÖRTERBUCH distinguishes between three different levels of semantic proximity. This distinction is indicated by means of three different types of counters: roman numbers, arabian numbers, and lowercase letters. The LEXPARSE grammar for the DUDEN-STILWÖRTERBUCH reflects these levels as follows, using three distinct built-in counters:

```
PAR  -> XFTbo: %I, ".", %R1.
PAA  -> XFTbo: %1, ".", %Ra.
PAB  -> XFTbo: %a, %BRPC.
```

The categories PAR, PAA, and PAB denote the different homonym/sense levels. Each category can be expanded if and only if the typeface state condition XFTbo: is met, i.e., the typeface indicator is in the state 'bold'.<sup>31</sup> All categories cover a separate counter symbol and provide the following counting scheme:

- ```
I. {Reset the counter: arabian numbers}
  1. {Reset the counter: lowercase letters}
     a)
     b)
  2. {Reset the counter: lowercase letters}
     a)
     b)
II. {Reset the counter: arabian numbers}
  1. {Reset the counter: lowercase letters}
  2. {Reset the counter: lowercase letters}
```

The counter for roman numbers, like all other counters, is reset automatically by the parser at

<sup>31</sup> Since the counter symbol must be in typeface 'bold', tokens covering the counter cannot be confused with tokens covering abbreviated lemmata which can, in principle, also appear at this position. This is particularly useful for lemmata with initial letters / or V which are abbreviated as "I." and "V." and could be easily confused with roman counter symbols.

the beginning of a new dictionary entry.<sup>32</sup>

#### 4.3.7. Applying the '\$\$' Operator

As pointed out already, example groups (denoted by the category **BeiGA**) may consist of several example sentences (denoted by the category **BeiA**) which are separated by the  $SI_{ntyp}^{def}$ , 'semicolon'. Simplified versions of the derivation rules for **BeiGA** and **BeiA** are given in (1):

- ```
(1) BeiGA -> BeiA < XPSEM BeiA > [ "." ^ ] .
    BeiA  -> $$.
(2) This is example 1;
    This is example 2;
    This is example 3.
(3) [This] [is] [example] [1] [XPSEM]
    [This] [is] [example] [2] [XPSEM]
    [This] [is] [example] [3] [. ]
(4) +--> BeiGA:
      +--> BeiA: "This is example 1"
      +--> BeiA: "This is example 2"
      +--> BeiA: "This is example 3"
```

Category **BeiA** expands to the \$\$-operator consuming any sequence of literals up to the next XCode token; input (2) thus results in the token list (3). Using the grammar in (1), one can successfully expand **BeiGA** and obtain the structure in (4).

All literal tokens being covered by category **BeiA** are collected in the accumulator of **BeiA**. Note that neither the semicolon nor the dot are added to the accumulators of the categories **BeiGA** or **BeiA**. This is due to the fact that, on the one hand, XCode tokens like **XPSEM** are, by default, not added to accumulators and that, on the other hand, the dot is explicitly marked with the '^' character not to be added.

#### 4.3.8. Applying the '\$' Operator

The application of the \$-operator can be illustrated by the expansion rule for the category **LZGA** (*Lemmazeichengestaltangabe*) which denotes items giving the form of the lemma sign.

```
LZGA -> XFTxb: $ [ "," ^ ] .
```

Category **LZGA** expands to a single literal token followed by an optional comma. The comma is an  $SI_{ntyp}^{ind}$ , with regard to the DUDEN-STILWÖRTERBUCH and is not assigned to an XCode. In order not to be added to the accumulator, the comma must be marked with the '^' character.

Note that **LZGA** can be expanded only if the typeface indicator is in the state 'extra bold' at the time of the expansion. This is expressed by the typeface condition **XFTxb**:

#### 4.3.9. Applying XCode Exceptions

The colon is a  $SI_{ntyp}^{ind}$ , with respect to the DUDEN-STILWÖRTERBUCH which, in most cases, functions as  $SI_{ntyp}$  separating items or groups of items from each other. For this reason, the colon is assigned to its corresponding XCode **XP COL**. In particular cases only, the colon occurs as part of an example group, like in the example sentence 'die beiden Mannschaften trennten sich 0:0', being part of the dictionary entry for **trennen**.

<sup>32</sup> The configuration switch [**Parser**] **AutoResetCounter** determines this behaviour. This switch is enabled by default.

To cope with these cases, the XCode XPCOL may be disabled by an XCode exception in the scope of category *BeiA*. Herein, the colon is generated as a literal token ':' rather than as an XCode XPCOL.

In the following example, semicolons and parentheses are considered to be  $SI_{ntyp}^{def}$ , and are therefore assigned to their corresponding XCodes. The colon, as well, is considered to be a  $SI_{ntyp}^{ind}$ , and is assigned to the XCode XPCOL. Category *BeiGA* can then be derived using the simplified grammar in (1):

- ```
(1) BeiGA -> BeiA < XPSEM, BeiA > [ "." ^ ].
    BeiA  -> -XPCOL:
           [ PragA ] $$ .
    PragA -> +XPCOL:
           XBRPO $$ XBRPC.
(2) (sl.) This is example 1;
    The teams left 0:0;
    This is example 3.
(3) [XPBRPO] [sl] [.] [XBRPC]
    [This] [is] [example] [1] [XPSEM]
    [The] [teams] [left] [0] [:] [0] [XPSEM]
    [This] [is] [example] [3] [.]
(4) +--> BeiGA:
    +--> BeiA: "This is example 1"
    +--> PragA: "sl."
    +--> BeiA: "The teams left 0:0"
    +--> BeiA: "This is example 3"
```

Category *BeiGA* expands to a sequence of at least one *BeiA* (*Beispielangabe*). Category *BeiA* expands to an optional category *PragA* (*Pragmatische Angabe*) and the \$\$-operator. Within the scope of *BeiA*, the XCode XPCOL is disabled by the XCode exception -XPCOL: in the second rule of (1). Within the scope of category *PragA* however, the colon functions as  $SI_{ntyp}$ , therefore it must be enabled by the XCode exception +XPCOL: in the third rule of (1) (as well as in other rules for daughter categories of *BeiA*).

The input (2) results in the token list (3). Using the grammar in (1), one can successfully expand *BeiGA* and obtain the structure given in (4).

## 5. The Implementation

LEXPARSE Version 1.10 is implemented in the C++ language and was developed using the BORLAND C++ Compiler for OS/2 Version 1.0. This compiler supports the full AT&T C++ standard version 2.1 (incl. streams and templates). The source code is compatible with the BORLAND C++ Compiler for DOS Version 3.1 and with the GNU C++ Compiler Version 2.4.1.

The development took place on an IBM compatible PC with Intel '486 CPU under OS/2 2.0. From the beginning, the program was designed to be easily ported to other platforms. The program currently runs under the MS-DOS, OS/2, and UNIX (SUN OS 4.1) operating systems.

The LEXPARSE program consists of several modules. Each module defines at least one object. Each object allows several operations to be applied, the so-called methods. The set of methods which may be applied from outside of a module are called interface methods. An interface

method consists of a name, an arity or argument list, a return value and – sometimes – of a precondition.

The object-oriented design has several advantages: Modules are easier to maintain and may be improved or replaced subsequently as long as the interface methods do not change. Furthermore, this design supports concepts such as data abstraction/encapsulation and structural/behavioral inheritance.

In the following, we give an overview of all program modules, the objects they define and their most important interface methods.

### 5.1. Module LEXPARSE

This module contains the main routine which initializes data structures and installs signal and error handlers. It reads and interpretes the configuration file by consulting the object *infile*. The grammar as described in the configuration file is passed to the object *grammar*. The main routine then opens the file containing the typesetting tape and executes a loop in which the parsing engine reads the input file and identifies consecutive entries according to the user-supplied grammar (cf. module PARSE). A special mechanism detects endless loops and aborts the program upon detecting one. All parse trees being generated by the parsing engine and being stored in the object *output* are sent to the output file if the corresponding entry is well-formed. If an entry is non-wellformed or if it is marked as being so, it is sent to the error output file.

The main routine counts all well-formed and non-wellformed entries and runs a timer to measure the time elapsed for parsing. At the end of the parsing process, the program generates a status report to inform the user about the process.

### 5.2. Module GRAMMAR

This module receives the grammar from the configuration file, converts each rule and stores all rules in the object *grammar*. The object *grammar* consists of an object 'starting symbol' and a set of objects 'non-terminals'. Each object 'non-terminal' consists of a list of objects 'alternate rules'. Each object 'alternate rule' consists of a list of objects 'nodes'. Each object 'node' corresponds to a specific element of the language as being defined by the grammar.

The following table gives an overview of all nodes which are distinguished by the object *grammar* and which must be treated by the parsing engine according to their semantics.

| Table of all Grammar Nodes    |   |
|-------------------------------|---|
| Node                          | Function                                |
| <code>typ_error</code>        | Error directive (@)                     |
| <code>typ_xcode</code>        | XCode                                   |
| <code>typ_category</code>     | Category                                |
| <code>typ_lt_single</code>    | A single literal (\$-operator)          |
| <code>typ_lt_multiple</code>  | A sequence of literals (\$\$-operator)  |
| <code>typ_lt_numerical</code> | Numerical literal (\#-operator)         |
| <code>typ_lt_expect</code>    | A specific literal (string constant)    |
| <code>typ_lt_rep</code>       | A specific literal (regular expression) |
| <code>typ_skip_incl</code>    | Inclusive skip operator (*)             |
| <code>typ_skip_excl</code>    | Exclusive skip operator (*\)            |
| <code>typ_except</code>       | Except operator (\)                     |
| <code>typ_tf_state</code>     | Typeface state                          |
| <code>typ_tf_state_set</code> | Set a typeface state                    |

| Table of all Grammar Nodes     |                                 |
|--------------------------------|---------------------------------|
| Node                           | Function                        |
| <code>typ_optional</code>      | Optional rule                   |
| <code>typ_opt_end</code>       | End of an optional rule         |
| <code>typ_mul_optional</code>  | Multiple optional rule          |
| <code>typ_mul_opt_end</code>   | End of a multiple optional rule |
| <code>typ_set_terms</code>     | Set of terminals                |
| <code>typ_set_terms_end</code> | End of a set of terminals       |
| <code>typ_cnt_roman</code>     | Counter: i, ii, iii, iv, ...    |
| <code>typ_cnt_ROMAN</code>     | Counter: I, II, III, IV, ...    |
| <code>typ_cnt_latin</code>     | Counter: a, b, c, d, ...        |
| <code>typ_cnt_LATIN</code>     | Counter: A, B, C, D, ...        |
| <code>typ_cnt_arab</code>      | Counter: 1, 2, 3, 4, ...        |
| <code>typ_res_roman</code>     | Reset counter (i, ii, iii, ...) |
| <code>typ_res_ROMAN</code>     | Reset counter (I, II, III, ...) |
| <code>typ_res_latin</code>     | Reset counter (a, b, c, ...)    |
| <code>typ_res_LATIN</code>     | Reset counter (A, B, C, ...)    |
| <code>typ_res_arab</code>      | Reset counter (1, 2, 3, ...)    |
| <code>typ_res_all</code>       | Reset all counters              |

The object *grammar* does not rearrange grammar rules, i.e., *grammar* does not transform an optional rule as specified under (1) into a pure context-free rule as specified under (2) by generating an additional category AUX.

```
(1) A    -> [ X ].
(2) A    -> AUX.
    AUX -> X
        | .
```

This forces the parsing engine to handle all specified grammar nodes, especially the nodes `typ_optional`, `typ_mul_optional` or `typ_set_terms`. This procedure is taken up with a view to greater clarity: The parser does not automatically generate new categories which might confuse a user on reading the logfiles during the development (or debugging) of a grammar. The grammar may be read in the logfiles as it was specified by the user.

The object *grammar* offers several interface methods such as `add_rule/2`, `add_label/2` and `add_start/1` with which rules, labels or a starting symbol may be added to the grammar. The object *grammar* may be passed to other modules as a constant object, i.e., a 'read-only' object which cannot be modified.

### 5.3. Module PPROCESS

This module implements the generic object *pprocess* for replacing string patterns. The preprocessor and the postprocessor are derived from this object. The object supports four different actions to be applied on a string:

- Delete a specific pattern within a string.
- Delete the whole string in which a specific pattern could be found.
- Replace a specific pattern by another pattern within a string.
- Delete a specific range of columns from a string.

The pattern specified for an action may be either a string constant or a regular expression string.

*pprocess* offers the interface method `add_action/2` with which complete actions may

be specified.

## 5.4. Module PREP

This module contains the preprocessor (object *prep*) for the parsing engine. The preprocessor reads the input file on a line-by-line basis and is able to apply seven different kinds of actions (*directives*) to the input:

*Delete*: Delete a specific pattern in the input line.

*DeleteRep*: Delete a specific regular expression in the input line.

*DelLine*: Delete the whole input line in which the specified pattern could be found.

*DelLineRep*: Delete the whole input line in which the specified regular expression could be found.

*Change*: Replace a specific pattern by another pattern within the input line.

*ChangeRep*: Replace a specific regular expression by a replace pattern within the input line.

*DeleteCol*: Delete a specific range of columns from the line.

*Note*: The order of all directives as specified in the configuration file is preserved during execution.

The **Change** directive offers the possibility of inserting XCode tokens into the input stream. For doing so, the desired XCode must be enclosed by two '\#' characters.

```
Change = "&12&34" -> "#XFLBE# #XFTxb#"
```

*prep* offers the interface method **add\_action/2** with which complete directives may be specified.

The interface method **get\_next/1** returns the next line as being read from the input file and being worked on by the preprocessor.

## 5.5. Module SCAN

This module contains the scanner (or tokenizer) of the parser. The scanner object *scan* receives lines of characters from the preprocessor and splits them into separate tokens. The scanner distinguishes two classes of tokens: XCode tokens and literal tokens.<sup>33</sup>

By configuring the scanner, each XCode may be defined by specifying a pattern which is replaced by an XCode token if it occurs in the input. The object *scan* supports a large table of XCodes which all can be connected with a pattern.<sup>34</sup>

*scan* offers the interface method **add\_xcode/2** to define an XCode together with its pattern.

The interface method **lookahead/1** tells the caller the next token, that is available in the input. The token is not removed from the input so that subsequent calls to **lookahead/1** all result in the same token.

**get/1** reads the next token from the input. The token is removed from the input.

**unget/1** puts the specified token back onto the input, so that the next call to **get/1** results in this token. There are no limits concerning the number of calls to **unget/1**.

The object *scan* also supports several interface methods to expand abbreviations: **expand\_define/2** allows the definition of a string which is returned as a literal token whenever a sequence of the initial letter of this string followed by a dot occurs in the input. This expansion must be enabled using **expand\_enable/0** before any calls to **get/1** and may be disabled using method **expand\_disable/0** after such a call.

<sup>33</sup> For a complete discussion of these different tokens refer to section 3.2.

<sup>34</sup> A list of all supported XCodes can be found in section A in the appendix.



## 5.6. Module PARSE

This module contains the object *parse* representing the parsing engine. The parsing engine consults the grammar representation stored in the object *grammar* and reads tokens from the scanner. All categories which can be expanded are sent to the object *output* in the module OUTPUT.

The parsing engine uses the algorithm of recursive descent parsing. This is a top-down, depth-first algorithm which heavily uses backtracking, i.e., reading and consulting the tokens repeatedly.<sup>35</sup>

*parse* offers the interface methods *accept/1* and *accept\_eof/0*. *accept/1*, is a polymorphic function and may be executed with the following different arguments: *grammar*, *category*, *rule* and *grammar node*. Both methods have a boolean return value indicating whether the specified argument could be accepted or not. This enables the following main loop for parsing:

```
WHILE ( parser.accept_eof() = false ) DO
    result := parser.accept( Grammar )
    ...
    display parse tree
    ...
END
```

## 5.7. Module OUTPUT

This module contains the object *output* which is capable of storing and displaying the generated parse trees. The object *output* supports many switches and options to control the style of display not to be mentioned here in more detail.

*output* offers the interface methods *cat\_announce/3* (call whenever a category is predicted by the parsing engine) and *cat\_complete/1* (call whenever a category could be completed by the parsing engine) to construct the parse tree. Each time a category has to be removed from the parse tree (e.g., because of backtracking) the method *cat\_undo/1* may be called to remove subtrees from the whole parse tree. The complete parse tree may be displayed by sending the object *output* to an output stream (e.g., logfile or console).

The object *output* is partly derived from the object *pprocess* which represents the post-processor of the output.

## 5.8. Module ERROR

This module contains the error handler of the LEXPARSE program. The object *error* distinguishes between six different error classes.

*Panic*: Serious error – the program terminates immediately. Output files and logfiles are unusable.

*Critical*: Serious error – the program terminates immediately. Output files and logfiles may be incomplete.

*Internal*: Internal program error – the program terminates immediately. Output files and logfiles may be evaluated.

*Fatal*: Program error – the program terminates immediately. Output files and logfiles may be evaluated.

*Normal*: Program error – the program terminates later on because of this error condition. Output files and logfiles may be evaluated.

*Warning*: Warning message only – the program does not terminate because of this message. Output files and logfiles may be evaluated completely.

<sup>35</sup> For a detailed discussion of this refer to AHO, SETHI & ULLMAN 1988 p. 181 ff.

*error* offers the interface method `add_error/6` with which an error message and an error condition may be raised. `exec/0` displays all messages accumulated so far and executes the corresponding action, e.g., the program termination if a normal error was generated.

Each error message consists of an error class (see above), the name of the module in which the error occurred, an error code, an error description and an error location.

The error code contains three digits whereby the first digit identifies the module in which the error occurred:

| Error Code Classes |                           |
|--------------------|---------------------------|
| Error codes        | Error occurred in module: |
| 0xx                | Could not identify module |
| 1xx                | Error in module LEXPARSE  |
| 2xx                | Error in module PREP      |
| 3xx                | Error in module SCAN      |
| 4xx                | Error in module PARSE     |
| 5xx                | Error in module GRAMMAR   |
| 6xx                | Error in module OUTPUT    |
| 7xx                | Error in module INI       |
| 8xx                | Error in module PPROCESS  |
| 9xx                | Error in module GREP      |

## 5.9. Module GREP

This module contains the object *grep*, a generic object to search for regular expressions in strings.

## 5.10. Module INI

This module contains the object *inifile* handling the .INI configuration files. The object offers several interface methods to access the sections and options of the configuration file.

## 5.11. Module TOOLS

This module contains additional generic functions, e.g., string functions and (simple) user interface functions.

# 6. Conclusion and Outlook

The LEXPARSE program is a powerful tool for dictionary entry parsing and the developing of dictionary entry grammars. The grammar formalism offers various features to meet the specific needs of dictionary entry parsing, e.g., the handling of structure indicators, font codes and counters. LEXPARSE promotes a quick development of dictionary entry grammars providing a rich set of facilities for identifying inadequate grammars as well as detecting non-wellformed dictionary entries. The parse trees being generated by the LEXPARSE program may be displayed in different, user-definable styles, offering both easy-to-read and comprehensive graphical representations and tagged structures for a computer evaluation.

Depending on the complexity and the size of the grammar and the entries, the parser runs at a speed of up to 10 entries per second on a PC with '486/33 CPU under the OS/2 operating system.

The program was tested on several dictionaries: a grammar was developed for the DUDEN-

STILWÖRTERBUCH at the Department of Linguistics at the University of Tübingen. This grammar covers over 98% of all entries.<sup>36</sup> A grammar for the DUDEN-BEDEUTUNGSWÖRTERBUCH is currently being worked out at the same institute. Another LEXPARSE grammar is being developed within the project "DEUTSCHES RECHTSWÖRTERBUCH" carried out at the Academy of Science in Heidelberg. In the German Department of the University of Heidelberg work has begun in elaborating a grammar for the "FRÜHNEUHOCHDEUTSCHES WÖRTERBUCH".<sup>37</sup>

We have several plans for enhancing the LEXPARSE system in the future:

- We are aiming at an integration with the *ELWIS-TOOLBOX* – a collection of different programs for intellectually analyzing dictionary entry structures.
- An elegant and efficient handling of nesting and niching dictionaries will be integrated into the system.
- An interactive debugger facilitating an easier development of dictionary grammars will soon be added to the program.

We hope that the experiences gathered in future applications of the system on varied types of dictionaries will assist us in further developing LEXPARSE into a flexible, efficient, multi-task tool in dictionary entry parsing.

## 7. References

- ALFRED V. AHO, RAVI SETHI & JEFFREY D. ULLMAN (1988): *Compilers, Principles, Techniques and Tools*. Reading, Massachusetts: Addison Wesley.
- BRIGITTE BLÄSER & MATTHIAS WERMKE (1990): Projekt 'Elektronische Wörterbücher / Lexika' Abschlußbericht der Definitionsphase. IWBS Report 145. Wissenschaftliches Zentrum / Institut für Wissensbasierte Systeme der IBM Deutschland GmbH.
- CHRISTOPH BLÄSI & HEINZ-DETLEV KOCH (1991): Maschinelle Strukturerschließung von Wörterbuchartikeln mit Standardmethoden. In: *Lexicographica*. Vol. 7, (1991) 1992, pp. 182-201. Tübingen: Niemeyer.
- BRAN BOGURAEV (1991): Building a Lexicon: The Contribution of Computers. In: *International Journal of Lexicography* Vol. 4, (1991), pp. 227-260
- BRAN BOGURAEV & TED BRISCOE (Eds.) (1989): *Computational Lexicography for Natural Language Processing*. London, New York: Longman.
- MARTIN BRYAN (1988): *SGML – An Author's Guide to the Standard Generalized Markup Language*. Waltham, Massachusetts: Addison Wesley.
- GÜNTHER DROSDOWSKI (Ed.) (1988): *DUDEN STILWÖRTERBUCH der deutschen Sprache*. Mannheim: Duden.
- SABINE ENGELKE (1994, forthcoming): *Grammatikentwicklung für Wörterbücher mit dem LEXPARSE Formalismus*. SFS Report, Universität Tübingen.
- RALF HAUSER (1993): *LEXPARSE User's Manual*. SFS Report 10-93, Seminar für Sprachwissenschaft der Universität Tübingen.
- RALF HAUSER & ANGELIKA STORRER (1993): *LEXPARSE – Ein Parser zur maschinellen Analyse von Wörterbuchstrukturen*. SFS Report 9-93, Seminar für Sprachwissenschaft der Universität Tübingen.
- FRANZ JOSEF HAUSMANN, OSKAR REICHMANN, HERBERT ERNST WIEGAND & LADISLAV ZGUSTA (Eds.) (1989): *Wörterbücher. Ein internationales Handbuch zur Lexikographie*. Berlin: de Gruyter.
- FRANZ JOSEF HAUSMANN & HERBERT ERNST WIEGAND (1989): *Component Parts and Structures of General Monolingual Dictionaries: A Survey*. In: *Wörterbücher. Ein internationales Handbuch zur Lexikographie*. Edited by FRANZ J. HAUSMANN, OSKAR REICHMANN, HERBERT E. WIEGAND & LADISLAV ZGUSTA. Berlin: de Gruyter, Vol. 1, pp. 328-360.
- K. JENSEN & N. WIRTH (1985): *Pascal User Manual and Report*. Berlin, New York: Springer.
- MARY S. NEFF & BRANIMIR K. BOGURAEV (1990): *From Machine-Readable Dictionaries to Lexical Data Bases*. Research Report RC #16080 (71353) 8/31/90, IBM THOMAS J. WATSON Research Center, Yorktown

<sup>36</sup> About 10% of the entries were non-wellformed, i.e., structure indicators were absent or counters were invalid. Cf. ENGELKE 1994.

<sup>37</sup> As far as can be seen right now, both grammars will cover distinctly more than 90% of all entries.

Heights, New York..

- MARY S. NEFF, R. J. BYRD & O. A. RIZK (1988): Creating and querying lexical data bases. In: Proceedings of the Second ACL Conference on Applied Natural Language Processing Austin, Texas; pp. 84-92.
- ANGELIKA STORRER (1994, forthcoming): Methodisches Vorgehen beim Entwickeln von Artikelstrukturgrammatiken. SFS Report, Universität Tübingen.
- ANGELIKA STORRER, HELMUT FELDWEIG & ERHARD HINRICHS (1993): Korpusunterstützte Entwicklung lexikalischer Wissensbasen. To appear in: Sprache und Datenverarbeitung 17, 1993, pp. 59-72.
- HERBERT ERNST WIEGAND (1989a): Aspekte der Makrostruktur im allgemeinen einsprachigen Wörterbuch: alphabetische Anordnungsformen und ihre Probleme. In: Wörterbücher. Ein internationales Handbuch zur Lexikographie. Edited by FRANZ J. HAUSMANN, OSKAR REICHMANN, HERBERT E. WIEGAND & LADISLAV ZGUSTA. Berlin: de Gruyter, Vol. 1, pp. 371-409.
- HERBERT ERNST WIEGAND (1989b): Der Begriff der Mikrostruktur: Geschichte, Probleme, Perspektiven. In: Wörterbücher. Ein internationales Handbuch zur Lexikographie. Edited by FRANZ J. HAUSMANN, OSKAR REICHMANN, HERBERT E. WIEGAND & LADISLAV ZGUSTA. Berlin: de Gruyter, Vol. 1, pp. 409-462.
- HERBERT ERNST WIEGAND (1989c): Formen von Mikrostrukturen im allgemeinen einsprachigen Wörterbuch. In: Wörterbücher. Ein internationales Handbuch zur Lexikographie. Edited by FRANZ J. HAUSMANN, OSKAR REICHMANN, HERBERT E. WIEGAND & LADISLAV ZGUSTA. Berlin: de Gruyter, Vol. 1, pp. 462-501.
- HERBERT ERNST WIEGAND (1989d): Die deutsche Lexikographie der Gegenwart. In: Wörterbücher. Ein internationales Handbuch zur Lexikographie. Edited by FRANZ J. HAUSMANN, OSKAR REICHMANN, HERBERT E. WIEGAND & LADISLAV ZGUSTA. Berlin: de Gruyter, Vol. 2, pp. 2100-2246.
- HERBERT ERNST WIEGAND (1991): Printed Dictionaries and their Parts as Text. An Overview of More Recent Research as an Introduction. In: Lexicographica. Vol. 6, (1990) 1991, pp. 1-124. Tübingen: Niemeyer.

Ralf Hauser, Seminar für Sprachwissenschaft der Universität Tübingen, e-mail: affie@sfs.nphil.uni-tuebingen.de

Angelika Storrer, Institut für deutsche Sprache Mannheim, e-mail: storrer@ids-mannheim.de

## Appendixes

### A List of All Available XCodes

The following table contains a list of all XCodes available in the LEXPARSE system together with a description of their semantics.<sup>38</sup>

*Note:* All pre- and userdefined XCodes are case sensitive.

| List of All Available XCodes |           |                            |
|------------------------------|-----------|----------------------------|
| XCode                        | Character | Description                |
| XBRPO                        | (         | BracketParenthesisOpen     |
| XBRPC                        | )         | BracketParenthesisClose    |
| XBRCB                        | {         | BracketCurlyOpen           |
| XBRCE                        | }         | BracketCurlyClose          |
| XBRAO                        | <         | BracketAngleOpen           |
| XBRAC                        | >         | BracketAngleClose          |
| XBRSO                        | [         | BracketSquareOpen          |
| XBRSC                        | ]         | BracketSquareClose         |
| XPDOT                        | .         | PunctuationDot             |
| XPCOM                        | ,         | PunctuationComma           |
| XPSEM                        | ;         | PunctuationSemicolon       |
| XPCOL                        | :         | PunctuationColon           |
| XPEXC                        | !         | PunctuationExclamationMark |

<sup>38</sup> The list of all available XCodes may be extended in future releases of LEXPARSE.

| List of All Available XCodes |           |                           |
|------------------------------|-----------|---------------------------|
| XCode                        | Character | Description               |
| XPQUE                        | ?         | PunctuationQuestionMark   |
| XPHYP                        | -         | PunctuationHyphen         |
| XPRDT                        | •         | PunctuationRoundedDot     |
| XPCDT                        | ·         | PunctuationCenteredDot    |
| XPAST                        | *         | PunctuationAsterisk       |
| XCSLA                        | /         | CharacterSlash            |
| XCBSL                        | \         | CharacterBackslash        |
| XCAMP                        | &         | CharacterAmpersand        |
| XCPAR                        | §         | CharacterParagraph        |
| XCSEP                        |           | CharacterSeparator        |
| XCARR                        | ^         | CharacterArrow            |
| XCQUO                        | “         | CharacterQuoteOpen        |
| XCQUC                        | ”         | CharacterQuoteClose       |
| XCSCO                        | _         | CharacterUnderscore       |
| XCCTL                        | ~         | CharacterTilde            |
| XCACC                        | ˘         | CharacterAccent           |
| XCATS                        | @         | CharacterATSign           |
| XFLBE                        |           | FormatLemmaBegin          |
| XFLEN                        |           | FormatLemmaEnd            |
| XFPBE                        |           | FormatParagraphBegin      |
| XFPEN                        |           | FormatParagraphEnd        |
| XFTAB                        |           | FormatTabulator           |
| XFTst                        |           | FormatTypefaceStandard    |
| XFTit                        |           | FormatTypefaceItalic      |
| XFTsl                        |           | FormatTypefaceSlanted     |
| XFTti                        |           | FormatTypefaceTiny        |
| XFTca                        |           | FormatTypefaceCapital     |
| XFTtt                        |           | FormatTypefaceTypewriter  |
| XFTsp                        |           | FormatTypefaceSuperscript |
| XFTsb                        |           | FormatTypefaceSubscript   |
| XFTbo                        |           | FormatTypefaceBold        |
| XFTxb                        |           | FormatTypefaceExtraBold   |
| XFTwi                        |           | FormatTypefaceWide        |
| XFTf0                        |           | FormatTypefaceFont0       |
| XFTf1                        |           | FormatTypefaceFont1       |
| XFTf2                        |           | FormatTypefaceFont2       |
| XFTf3                        |           | FormatTypefaceFont3       |
| XFTf4                        |           | FormatTypefaceFont4       |
| XFTf5                        |           | FormatTypefaceFont5       |
| XFTf6                        |           | FormatTypefaceFont6       |
| XFTf7                        |           | FormatTypefaceFont7       |
| XFTf8                        |           | FormatTypefaceFont8       |
| XFTf9                        |           | FormatTypefaceFont9       |
| X_EOF                        |           | EndOfFile                 |

Note: The semantics of XCodes of classes XF\*\*\* (Format) and XU\*\*\* (User-definable) cannot be expressed by a character.

## B A Sample '.INI' Configuration File for LEXPARSE

```

;
;   General INI file for LEXPARSE
;
;   sections PREPROCESS, SCANNER, PARSER
;   and GRAMMAR are exported to include file
include    = STILWORT.19      ; grammar

[Settings]
; Debug          = A
TypefaceStates = On
Preprocess      = On
SkipNLchar      = No
; StepEntry      = Yes
; StepSingle     = Yes
FirstEntry      = 3
LastEntry       = 9

[Logfiles]
preprocess      = c:\tmp\prep.log
scan            = c:\tmp\scan.log
parse           = c:\tmp\parse.log
output          = parsed\output.log
error           = parsed\error.log
ErrorsOnly      = Yes

[Display] ; for a display in TREE style
Parselevel      = Yes
Mode            = T
Unlabeled       = Yes
EmptyCat        = Yes
Tree            = " | "
Trees           = " +-->"
TreeLD          = " +-->"
Title           = Yes
Header          = "---- ParseTree ----\n"
Footer          = "-----\n"
Format          = " %c[ - %l][ \"%$\"]\n"
RMargin        = 79      ; for 80 column display
Progress        = Off
ToParseLog      = Yes

```

## C A Sample LEXPARSE Grammar

The following sample grammar describes a very small fragment of the DUDEN STILWÖRTERBUCH.

```

[PreProcess]
ChangeRep = "}}$" -> " " ; end-of-line marker
; begin of a lemma and boldface
Change    = "òT2ûò5û" -> "#XFLBE# #XFTxb#"
DeleteRep = "}-}$"      ; hyphen
Delete    = "òvr10û"    ; italic correction

```

```

[Scanner]
AddLetters = "äöüÄÖÜ"
XPSEM    = ";"
XPCOL    = ":"
XFLEN    = "es"
; control codes for typefaces ('Settings:TypefaceStates')
XFTst    = "ò1û"
XFTit    = "ò2û"
XFTbo    = "ò3û"
XFTxb    = "ò5û"

[Parser]
Start    = WA
IndentLog = On
RecoverCounterError = Yes
ExpandAbbrev = LZGA
SkipCriticalIsError = Yes

[Grammar]
;--- Wörterbuchartikel
WA      -> XFLBE, FK, SK, XFLEN
        | WA_Err.                ; error handler
WA_Err  -> * XFLEN @.            ; error handler
;--- Formkommentar
FK      -> LZGA [ GrA ] XPCOL
        | FK_Err.                ; error handler
FK_Err  -> * XPCOL @.            ; error handler
LZGA    -> XFTxb: $ [ ", "^ ].    ; Lemmazeichengestaltangabe
GrA     -> GA.                    ; Grammatische Angabe
GA      -> XFTst: { "der", "die", "das" }.
;--- Semantischer Kommentar
SK      -> < PAA, PragSema, BeiGA >.
;--- Typen von Polysemieangaben
PAA     -> XFTbo: $1, ".", $Ra.    ; Polysemieangabe (arabisch)
;--- Pragmatische Semantische Angabe
PragSema-> BPA.
;--- Bedeutungsangaben
BPA     -> XFTit: * XPCOL.
;--- Beispielgruppenangabe
BeiGA   -> BeiA < XPSEM, BeiA > [ "." ].
;--- Beispielangabe
BeiA    -> -XPCOL:                ; no XCodes XPCOL herein!
        XFTst: $$.
```

```

[Labels]
WA      = "Wörterbuchartikel"
WA_Err  = "Fehler in Wörterbuchartikel"
FK      = "Formkommentar"
FK_Err  = "Fehler in Formkommentar"
LZGA    = "Lemmazeichengestaltangabe"
GA      = "Genusangabe"
GrA     = "Grammatische Angabe"
SK      = "Semantischer Kommentar"
PAA     = "Polysemieangabe (arabisch)"
PragSema = "Pragmatisch-Semantische Angabe"
```

```

BPA      = "Bedeutungsparaphrasenangabe"
BeiGA    = "Beispielgruppenangabe"
BeiA     = "Beispielangabe"

```

## D Sample Sessions

In the following we present several sample sessions of running the LEXPARSE program to parse dictionary entries being part of a typesetting tape for the Duden-Stilwörterbuch.

### D.1. Treffen, DUDEN-STILWÖRTERBUCH, Page 703

The following dictionary entry was parsed by LEXPARSE successfully. The parse tree for the entry **Treffen** is complete. Abbreviated lemma signs were expanded automatically by the program.

```

[D:\FOR\SNS\LEXPARSE]lexparse /entry -1 lexparse.ini data\9-entries.dat
LEXPARSE V1.10 [Nov 05 1993] <DEBUG-SRC> A Dictionary Entry Parser
(C) 1992,1993 Seminar für Sprachwissenschaft der Universität Tübingen (C) 1992,1993
(C) 1992,1993 Ralf Hauser
Created by 'Affie' (k.f. c|o)
Reading configuration: lexparse.INI
LEXPARSE: (W-303) Warning occurred in module 'Scanner':
Unknown XCode results in a new user-defined XCode:
"XURED"
Inputfile.....: data\9-entries.dat
Parsing.....: Enabled
First Entry.....: 1
Last Entry.....: 1
Time of day.....: Fri Nov 05 12:22:22 1993

```

Entry #1 accepted.

```

--- ParseTree -----
WA - Wörterbuchartikel
+--> FK - Formkommentar
| +--> LZGA - Lemmazeichengestaltangabe "Treffen"
| +--> GrA - Grammatische Angabe
|   +--> GA - Genusangabe "das"
+--> SK - Semantischer Kommentar
| +--> SSK1 - Semantischer Subkommentar 1. Stufe
| | +--> PAA - Polysemieangabe (arabisch) "1."
| | +--> PragSemA - Pragmatisch-Semantische Angabe
| | | +--> BA - Bedeutungsangabe
| | |   +--> BPA - Bedeutungsparaphrasenangabe "Zusammenkunft, Begegnung"
| | | +--> BeiGA - Beispielgruppenangabe
| | |   +--> BeiA - Beispielangabe "regelmäßige, seltene Treffen"
| | |   +--> BeiA - Beispielangabe "ein Treffen der Abiturienten"
| | |   +--> BeiA - Beispielangabe "ein Treffen der Außenminister"
| | |   +--> BeiA - Beispielangabe "ein Treffen verabreden, veranstalten"
| | |   +--> BeiA - Beispielangabe "an einem Treffen teilnehmen"
| | |   +--> BeiA - Beispielangabe "zu einem Treffen kommen"
| +--> SSK1 - Semantischer Subkommentar 1. Stufe
| | +--> PAA - Polysemieangabe (arabisch) "2."
| | +--> PragSemA - Pragmatisch-Semantische Angabe
| | | +--> PragA - Pragmatische Angabe "militär. veraltet"
| | | +--> BA - Bedeutungsangabe
| | |   +--> BPA - Bedeutungsparaphrasenangabe "Gefecht"
| | +--> BeiGA - Beispielgruppenangabe
| |   +--> BeiA - Beispielangabe "frische Truppen ins Treffen führen"
+--> SSK1 - Semantischer Subkommentar 1. Stufe
| +--> PAA - Polysemieangabe (arabisch) "3."
| +--> PragSemA - Pragmatisch-Semantische Angabe
| | +--> PragA - Pragmatische Angabe "Sport"
| | +--> BA - Bedeutungsangabe

```



```

| | +--> BPA - Bedeutungsparaphrasenangabe "Wettkampf"
| +--> BeiGA - Beispielgruppenangabe
| | +--> BeiA - Beispielangabe "ein faires, spannendes Treffen"
| | +--> BeiA - Beispielangabe "das Treffen endete unentschieden"
| | +--> BeiA - Beispielangabe "sie konnte das Treffen für sich entscheiden"
+--> PKP - Postkommentar zur Phraseologie
| +--> SKP - Subkommentar zur Phraseologie
| | +--> Praga - Pragmatische Angabe "geh."
| | +--> PhrasA - Phrasenangabe "etwas ins Treffen führen"
| | +--> KPB - Kommentar zur Phrasenbedeutung
| | +--> PBA - Phrasenbedeutungsangabe "etwas als Argument vorbringen"

```

-----  
Entry #1: Accepted!

```

Inputfile.....: data\9-entries.dat
First entry.....: 1
Last entry.....: 1
Entries parsed...: 1
Errors.....: 0
Time of day.....: Fri Nov 05 12:22:23 1993
Time of start....: Fri Nov 05 12:22:22 1993
Elapsed time.....: 00h 00m 01s

```

```

LEXPARSE V1.10 [Nov 05 1993] <DEBUG-SRC>
Program terminated successfully
LEXPARSE: (W-105) Warning occurred in module 'LEXPARSE':
Parsing terminated due to 'LastEntry' specification

```

## D.2. trennen, DUDEN-STILWÖRTERBUCH, Page 704

The following dictionary entry was parsed by LEXPARSE successfully. The parse tree for the entry **trennen** is not complete. Abbreviated lemma signs were expanded automatically by the program.

The entry is not well-formed since the item *Beispielangabe* (BeiA) *ein Gebirgszug trennt das Land in zwei Regionen.* does not contain a trailing semicolon although there are more categories BeiA.

LEXPARSE generated a maximum parse tree and all tokens which could be consumed by the parser following the error were collected in the accumulator of category *Fehler in Wörterbuchartikel* (WA\_Err).

```

[D:\FOR\SNS\LEXPARSE]lexparse /entry 7-7 lexparse.ini data\9-entries.dat
LEXPARSE V1.10 [Nov 05 1993] <DEBUG-SRC> A Dictionary Entry Parser
(C) 1992,1993 Seminar für Sprachwissenschaft der Universität Tübingen
(C) 1992,1993 Ralf Hauser
Created by 'Affie' (k.f. c|o)
Reading configuration: lexparse.INI
LEXPARSE: (W-303) Warning occurred in module 'Scanner':
Unknown XCode results in a new user-defined XCode:
"XURED"
Inputfile.....: data\9-entries.dat
Parsing.....: Enabled
First Entry.....: 7
Last Entry.....: 7
Time of day.....: Fri Nov 05 12:25:06 1993

```

6 Entries skipped!

N O T E: Entry #7 could not be accepted by grammar!

```

--- ParseTree -----
WA - Wörterbuchartikel
+--> FK - Formkommentar
| +--> LZGA - Lemmazeichengestaltangabe "trennen"
+--> SK - Semantischer Kommentar
| +--> SSK1 - Semantischer Subkommentar 1. Stufe
| | +--> PAA - Polysemieangabe (arabisch) "1."
| | +--> PragSema - Pragmatisch-Semantische Angabe
| | +--> SSK2 - Semantischer Subkommentar 2. Stufe
| | | +--> PAB - Polysemieangabe (Buchstabe) "a"
| | | +--> PragSema - Pragmatisch-Semantische Angabe

```

```

| | | | | ++-> GrA - Grammatische Angabe
| | | | | | ++-> Sma - Satzmusterangabe "jmdn., sich, etwas von jmdn., von etwas/ aus\
| | | | | | etwas trennen"
| | | | | | ++-> BA - Bedeutungsangabe
| | | | | | ++-> BPA - Bedeutungsparaphrasenangabe "lösen, entfernen, abtrennen"
| | | | | ++-> BeIGA - Beispielgruppenangabe
| | | | | | ++-> BeIA - Beispielangabe "eine Borte vom Kleid, das Futter aus dem Mantel\
| | | | | | trennen"
| | | | | | ++-> BeIA - Beispielangabe "den Kopf vom Rumpf trennen"
| | | | | | | ++-> GZB - Glossat zur Bedeutung "abschlagen"
| | | | | | ++-> BeIA - Beispielangabe "ein Tier von der Herde trennen"
| | | | | | ++-> BeIA - Beispielangabe "das Kind von seiner Mutter, von seiner Familie\
| | | | | | trennen"
| | | | | | ++-> BeIA - Beispielangabe "das Erz vom Gestein trennen"
| | | | | | ++-> BeIA - Beispielangabe "das Eigelb vom Eiweiß trennen"
| | | | | ++-> SSK2 - Semantischer Subkommentar 2. Stufe
| | | | | | ++-> PAB - Polysemieangabe (Buchstabe) "b"
| | | | | | ++-> PragSemA - Pragmatisch-Semantische Angabe
| | | | | | | ++-> GrA - Grammatische Angabe
| | | | | | | | ++-> Sma - Satzmusterangabe "jmdn., etwas trennen"
| | | | | | | ++-> BA - Bedeutungsangabe
| | | | | | | | ++-> BPA - Bedeutungsparaphrasenangabe "auseinanderbringen"
| | | | | | ++-> BeIGA - Beispielgruppenangabe
| | | | | | | ++-> BeIA - Beispielangabe "Eigelb und Eiweiß trennen"
| | | | | | | ++-> BeIA - Beispielangabe "Sauerstoff und Wasserstoff trennen"
| | | | | | | ++-> BeIA - Beispielangabe "die Bestandteile einer Mischung sorgfältig trennen"
| | | | | | | ++-> BeIA - Beispielangabe "die Nähte trennen"
| | | | | | | ++-> BeIA - Beispielangabe "die Streitenden mußten getrennt werden"
| | | | | ++-> SSK2 - Semantischer Subkommentar 2. Stufe
| | | | | | ++-> PAB - Polysemieangabe (Buchstabe) "c"
| | | | | | ++-> PragSemA - Pragmatisch-Semantische Angabe
| | | | | | | ++-> GrA - Grammatische Angabe
| | | | | | | | ++-> Sma - Satzmusterangabe "etwas trennen"
| | | | | | | ++-> BA - Bedeutungsangabe
| | | | | | | | ++-> BPA - Bedeutungsparaphrasenangabe "in seine Bestandteile zerlegen"
| | | | | | ++-> BeIGA - Beispielgruppenangabe
| | | | | | | ++-> BeIA - Beispielangabe "ein Kleid trennen"
| | | | | | | ++-> BeIA - Beispielangabe "ein Stoffgemisch chemisch, durch Kondensation\
| | | | | | trennen"
| | | | | | | ++-> BeIA - Beispielangabe "ein Wort nach Silben trennen"
| | | | | ++-> SSK2 - Semantischer Subkommentar 2. Stufe
| | | | | | ++-> PAB - Polysemieangabe (Buchstabe) "d"
| | | | | | ++-> PragSemA - Pragmatisch-Semantische Angabe
| | | | | | | ++-> GrA - Grammatische Angabe
| | | | | | | | ++-> Sma - Satzmusterangabe "jmdn. trennen"
| | | | | | | ++-> BA - Bedeutungsangabe
| | | | | | | | ++-> BPA - Bedeutungsparaphrasenangabe "auseinanderreißen"
| | | | | | ++-> BeIGA - Beispielgruppenangabe
| | | | | | | ++-> BeIA - Beispielangabe "die beiden Geschwister sollten nicht getrennt werden"
| | | | | | | ++-> BeIA - Beispielangabe "der Krieg hat die Familie getrennt"
| | | | | | | ++-> BeIA - Beispielangabe "nichts konnte die Liebenden trennen"
| | | | | ++-> SSK1 - Semantischer Subkommentar 1. Stufe
| | | | | | ++-> PAA - Polysemieangabe (arabisch) "2."
| | | | | | ++-> PragSemA - Pragmatisch-Semantische Angabe
| | | | | | | ++-> BA - Bedeutungsangabe
| | | | | | | | ++-> BPA - Bedeutungsparaphrasenangabe "klar unterscheiden, auseinanderhalten"
| | | | | ++-> SSK2 - Semantischer Subkommentar 2. Stufe
| | | | | | ++-> PAB - Polysemieangabe (Buchstabe) "a"
| | | | | | ++-> PragSemA - Pragmatisch-Semantische Angabe
| | | | | | | ++-> GrA - Grammatische Angabe
| | | | | | | | ++-> Sma - Satzmusterangabe "jmdn., etwas trennen"
| | | | | | | ++-> BeIGA - Beispielgruppenangabe
| | | | | | | | ++-> BeIA - Beispielangabe "die Begriffe klar, sauber trennen"
| | | | | | | | ++-> BeIA - Beispielangabe "wir müssen Person und Sache strikt trennen"
| | | | | ++-> SSK2 - Semantischer Subkommentar 2. Stufe
| | | | | | ++-> PAB - Polysemieangabe (Buchstabe) "b"

```

```

| | +--> PragSema - Pragmatisch-Semantische Angabe
| | | +--> Gra - Grammatische Angabe
| | | +--> Sma - Satzmusterangabe "jmdn., etwas von etwas trennen"
| | +--> BeiGA - Beispielgruppenangabe
| | | +--> BeiA - Beispielangabe "wir müssen die Person streng von der Sache trennen"
| | | +--> BeiA - Beispielangabe "mein Beruf kann von meiner Freizeit nicht streng\
| | | | getrennt werden"
| +--> SSK1 - Semantischer Subkommentar 1. Stufe
| | +--> PAA - Polysemieangabe (arabisch) "3."
| | | +--> PragSema - Pragmatisch-Semantische Angabe
| | +--> SSK2 - Semantischer Subkommentar 2. Stufe
| | | +--> PAB - Polysemieangabe (Buchstabe) "a"
| | | +--> PragSema - Pragmatisch-Semantische Angabe
| | | | +--> Gra - Grammatische Angabe
| | | | | +--> Sma - Satzmusterangabe "sich trennen"
| | | | +--> BA - Bedeutungsangabe
| | | | +--> BPA - Bedeutungsparaphrasenangabe "auseinandergehen"
| | | +--> BeiGA - Beispielgruppenangabe
| | | | +--> BeiA - Beispielangabe "wir trennten uns am Bahnhof"
| | | | +--> BeiA - Beispielangabe "unsere Wege trennen sich hier"
| | | | | +--> GZB - Glossat zur Bedeutung "jeder nimmt einen anderen Weg"
| | | | +--> BeiA - Beispielangabe "nach drei Stunden Diskussion trennte man sich"
| | | +--> KSV - Kommentar zur speziellen Verwendung
| | | | +--> SKSV - Subkommentar zur speziellen Verwendung
| | | | +--> VSA - Verwendungsspezifizierende Angabe
| | | | | +--> FGA - Fachgebietsangabe "Sport."
| | | | +--> BeiGA - Beispielgruppenangabe
| | | | | +--> BeiA - Beispielangabe "die beiden Mannschaften trennten sich 0: 0"
| +--> SSK2 - Semantischer Subkommentar 2. Stufe
| | +--> PAB - Polysemieangabe (Buchstabe) "b"
| | | +--> PragSema - Pragmatisch-Semantische Angabe
| | | | +--> Gra - Grammatische Angabe
| | | | | +--> Sma - Satzmusterangabe "sich von jmdm. trennen"
| | | | +--> BA - Bedeutungsangabe
| | | | +--> BPA - Bedeutungsparaphrasenangabe "weggehen"
| | | +--> BeiGA - Beispielgruppenangabe
| | | | +--> BeiA - Beispielangabe "vor der Haustür trennte er sich von mir"
| +--> SSK1 - Semantischer Subkommentar 1. Stufe
| | +--> PAA - Polysemieangabe (arabisch) "4."
| | | +--> PragSema - Pragmatisch-Semantische Angabe
| | +--> SSK2 - Semantischer Subkommentar 2. Stufe
| | | +--> PAB - Polysemieangabe (Buchstabe) "a"
| | | +--> PragSema - Pragmatisch-Semantische Angabe
| | | | +--> Gra - Grammatische Angabe
| | | | | +--> Sma - Satzmusterangabe "sich trennen"
| | | | +--> BA - Bedeutungsangabe
| | | | +--> BPA - Bedeutungsparaphrasenangabe "eine Partnerschaft, Gemeinschaft\
| | | | | auflösen"
| | | +--> BeiGA - Beispielgruppenangabe
| | | | +--> BeiA - Beispielangabe "wir haben uns endlich, nach zwei Jahren,\
| | | | | freundschaftlich, im guten getrennt"
| | | | +--> BeiA - Beispielangabe "die beiden Teilhaber haben sich getrennt"
| | | +--> KSV - Kommentar zur speziellen Verwendung
| | | | +--> SKSV - Subkommentar zur speziellen Verwendung
| | | | +--> VSA - Verwendungsspezifizierende Angabe
| | | | | +--> WKA - Wortartenkonversionsangabe "adj."
| | | | | +--> WKA - Wortartenkonversionsangabe "Part."
| | | | +--> BeiGA - Beispielgruppenangabe
| | | | | +--> BeiA - Beispielangabe "die Eheleute leben getrennt"
| +--> SSK2 - Semantischer Subkommentar 2. Stufe
| | +--> PAB - Polysemieangabe (Buchstabe) "b"
| | | +--> PragSema - Pragmatisch-Semantische Angabe
| | | | +--> Gra - Grammatische Angabe
| | | | | +--> Sma - Satzmusterangabe "sich von jmdm. trennen"
| | | | +--> BA - Bedeutungsangabe
| | | | +--> BPA - Bedeutungsparaphrasenangabe "sich loslösen"

```

```

| | +--> BeiGA - Beispielgruppenangabe
| | | +--> BeiA - Beispielangabe "sieht sich von ihrem Mann getrennt"
| | | +--> BeiA - Beispielangabe "von meinem Gesangspartner habe ich mich getrennt"
| | +--> KSV - Kommentar zur speziellen Verwendung
| | +--> SKSV - Subkommentar zur speziellen Verwendung
| | | +--> VSA - Verwendungsspezifizierende Angabe
| | | | +--> MA - Metaphernangabe "übertr."
| | | +--> PragA - Pragmatische Angabe "verhüll."
| | | +--> BeiGA - Beispielgruppenangabe
| | | +--> BeiA - Beispielangabe "die Firma hat sich von diesem Mitarbeiter\
| | | getrennt"
| +--> SSK1 - Semantischer Subkommentar 1. Stufe
| | +--> PAA - Polysemieangabe (arabisch) "5."
| | +--> PragSemA - Pragmatisch-Semantische Angabe
| | | +--> GrA - Grammatische Angabe
| | | | +--> Sma - Satzmusterangabe "sich von etwas trennen"
| | | +--> BA - Bedeutungsangabe
| | | | +--> BPA - Bedeutungsparaphrasenangabe "etwas hergeben"
| | | +--> BeiGA - Beispielgruppenangabe
| | | +--> BeiA - Beispielangabe "sich von Erinnerungsstücken nur ungern trennen, nicht\
| | | trennen können"
| | | +--> BeiA - Beispielangabe "sich von jeglichem Besitz trennen"
| | +--> KSV - Kommentar zur speziellen Verwendung
| | +--> SKSV - Subkommentar zur speziellen Verwendung
| | | +--> VSA - Verwendungsspezifizierende Angabe
| | | | +--> MA - Metaphernangabe "übertr."
| | | +--> BeiGA - Beispielgruppenangabe
| | | +--> BeiA - Beispielangabe "sich von einem Gedanken, einem Wunsch, einer\
| | | Vorstellung trennen müssen"
| | | +--> BeiA - Beispielangabe "sich von einem Anblick nicht trennen können"
| +--> SSK1 - Semantischer Subkommentar 1. Stufe
| | +--> PAA - Polysemieangabe (arabisch) "6."
| | +--> PragSemA - Pragmatisch-Semantische Angabe
| | +--> SSK2 - Semantischer Subkommentar 2. Stufe
| | | +--> PAB - Polysemieangabe (Buchstabe) "a"
| | | +--> PragSemA - Pragmatisch-Semantische Angabe
| | | | +--> GrA - Grammatische Angabe
| | | | | +--> Sma - Satzmusterangabe "etwas trennt etwas"
| | | | +--> BA - Bedeutungsangabe
| | | | | +--> BPA - Bedeutungsparaphrasenangabe "etwas bildet eine Grenze, ein\
| | | | | Hindernis zwischen etwas"
| | | +--> BeiGA - Beispielgruppenangabe
| | | | +--> BeiA - Beispielangabe "ein Stacheldraht trennt die Grundstücke"
| | | | +--> BeiA - Beispielangabe "ein Zaun trennt die Gärten"
| | +--> KSV - Kommentar zur speziellen Verwendung
| | +--> SKSV - Subkommentar zur speziellen Verwendung
| | | +--> VSA - Verwendungsspezifizierende Angabe
| | | | +--> MA - Metaphernangabe "übertr."
| | | +--> BeiGA - Beispielgruppenangabe
| | | | +--> BeiA - Beispielangabe "uns trennen Welten"
| | | | | +--> GZB - Glossat zur Bedeutung "wir sind äußerst verschieden"
| | | | +--> BeiA - Beispielangabe "die verschiedene Herkunft trennte sie"
| +--> SSK2 - Semantischer Subkommentar 2. Stufe
| | +--> PAB - Polysemieangabe (Buchstabe) "b"
| | +--> PragSemA - Pragmatisch-Semantische Angabe
| | | +--> GrA - Grammatische Angabe
| | | | +--> Sma - Satzmusterangabe "etwas trennt jmdn., etwas von jmdn., von etwas"
| | | +--> BA - Bedeutungsangabe
| | | | +--> BPA - Bedeutungsparaphrasenangabe "etwas grenzt jmdn., etwas gegen\
| | | | jmdn., etwas ab"
| | +--> BeiGA - Beispielgruppenangabe
| | | +--> BeiA - Beispielangabe "der Kanal trennt England vom Kontinent"
| | | +--> BeiA - Beispielangabe "nur ein Graben trennt die Zoobesucher von den\
| | | Elefanten"
| | | +--> BeiA - Beispielangabe "eine Glaswand trennt ihn von seinem Verteidiger"
| | | +--> UGrA - Unspezifizierte grammatische Angabe "auch ohne Präp.- Obj."

```

```

|      +--> BeiA - Beispielangabe "ein Gebirgszug trennt das Land in zwei Regionen"
+--> WA_Err - Fehler in Wörterbuchartikel "übertr.: nur noch wenige Tage trennen uns von den\
Wahlen. XFTbo 7. XFTst( Rundf.)< etwas trennt etwas; mit Artangabe> XFTit etwas besitzt\
eine bestimmte Trennschärfe: XFTst das Radio trennt die Sender gut, scharf, nicht\
richtig, genügend. XFTbo 8. XFTst< jmdn., etwas trennen> XFTit( eine telefonische\
Verbindung unterbrechen: XFTst die Verbindung wurde getrennt; man hat uns getrennt."

```

-----

Entry #7: Error caused by directive ('@' Operator)

Inputfile.....: data\9-entries.dat

First entry.....: 7

Last entry.....: 7

Entries parsed....: 1

Errors.....: 1

Time of day.....: Fri Nov 05 12:25:12 1993

Time of start....: Fri Nov 05 12:25:06 1993

Elapsed time.....: 00h 00m 06s

LEXPARSE V1.10 [Nov 05 1993] <DEBUG-SRC>

Program terminated successfully

LEXPARSE: (W-105) Warning occurred in module 'LEXPARSE':

Parsing terminated due to 'LastEntry' specification

### D.3. überhaupt, DUDEN-STILWÖRTERBUCH, Page 716

The following dictionary entry was parsed by LEXPARSE successfully. The parse tree for the entry **überhaupt** is complete. Abbreviated lemma signs were expanded automatically by the program.

However, the entry is not well-formed since the enumeration in the item *Polysemieangabe* (PAA) is invalid. The parser did find an invalid symbol '5' instead of the expected symbol '4' in the scope of the forth instance of item *Semantischer Subkommentar* (SSK).

Nevertheless, the parser generated a complete parse tree and issued an error message (warning #403) indicating the invalid enumeration.<sup>39</sup>

```

[D:\FOR\SNS\LEXPARSE]lexparse lexparse.ini data\counter_error.dat
LEXPARSE V1.10 [Nov 05 1993] <DEBUG-SRC> A Dictionary Entry Parser
(C) 1992,1993 Seminar für Sprachwissenschaft der Universität Tübingen
(C) 1992,1993 Ralf Hauser
Created by 'Affie' (k.f. c/o)
Reading configuration: lexparse.INI
LEXPARSE: (W-303) Warning occurred in module 'Scanner':
Unknown XCode results in a new user-defined XCode:
"XURED"
Inputfile.....: data\counter_error.dat
Parsing.....: Enabled
First Entry.....: 1
Last Entry.....: no spec.
Time of day.....: Fri Nov 05 12:28:42 1993

```

N O T E: Entry #1 could not be accepted by grammar!

--- ParseTree -----

WA - Wörterbuchartikel

+--> FK - Formkommentar

| +--> LZGA - Lemmazeichengestaltangabe "überhaupt"

+--> SK - Semantischer Kommentar

| +--> SSK0 - Semantischer Subkommentar 0. Stufe

| | +--> PAR - Polysemieangabe (roemisch) "I."

| | +--> PragSemA - Pragmatisch-Semantische Angabe

| | | +--> GRA - Grammatische Angabe

| | | +--> WAA - Wortartangabe "Adverb"

| | +--> SSK1 - Semantischer Subkommentar 1. Stufe

| | | +--> PAA - Polysemieangabe (arabisch) "1."

| | | +--> PragSemA - Pragmatisch-Semantische Angabe

| | | | +--> BA - Bedeutungsangabe

<sup>39</sup> The switch [Parser] RecoverCounterError = Yes does achieve such a behaviour.

```

| | | +--> BPA - Bedeutungsparaphrasenangabe "insgesamt, aufs Ganze gesehen"
| | +--> BeiGA - Beispielgruppenangabe
| | +--> BeiA - Beispielangabe "ich habe ihn gestern nicht angetroffen, er ist\
| | | | überhaupt selten zu Hause"
| | +--> BeiA - Beispielangabe "mir gefällt es in Madrid, überhaupt in Spanien"
| +--> SSK1 - Semantischer Subkommentar 1. Stufe
| | +--> PAA - Polysemieangabe (arabisch) "2."
| | +--> PragSemA - Pragmatisch-Semantische Angabe
| | | +--> GrA - Grammatische Angabe
| | | | +--> UGrA - Unspezifizierte grammatische Angabe "verstärkend bei Verneinungen"
| | | +--> BA - Bedeutungsangabe
| | | +--> BPA - Bedeutungsparaphrasenangabe "ganz und gar"
| | +--> BeiGA - Beispielgruppenangabe
| | +--> BeiA - Beispielangabe "das ist überhaupt nicht möglich, nicht wahr"
| | +--> BeiA - Beispielangabe "davon kann überhaupt keine Rede sein"
| | +--> BeiA - Beispielangabe "er hat heute überhaupt noch nichts gegessen"
| | +--> BeiA - Beispielangabe "das geht ihn überhaupt nichts an"
+--> SSK1 - Semantischer Subkommentar 1. Stufe
| +--> PAA - Polysemieangabe (arabisch) "3."
| +--> PragSemA - Pragmatisch-Semantische Angabe
| | +--> BA - Bedeutungsangabe
| | +--> BPA - Bedeutungsparaphrasenangabe "abgesehen davon, überdies"
| +--> BeiGA - Beispielgruppenangabe
| +--> BeiA - Beispielangabe "du kannst einmal nachfragen, und überhaupt solltest\
| | | du dich/ und überhaupt, du solltest dich mehr darum kümmern"
+--> SSK1 - Semantischer Subkommentar 1. Stufe
| +--> PAA - Polysemieangabe (arabisch) "5."
| +--> PragSemA - Pragmatisch-Semantische Angabe
| | +--> BA - Bedeutungsangabe
| | +--> BPA - Bedeutungsparaphrasenangabe "gerade, besonders"
+--> BeiGA - Beispielgruppenangabe
| +--> BeiA - Beispielangabe "wir gehen gerne im Wald spazieren, überhaupt im\
| | | Herbst"
| +--> BeiA - Beispielangabe "man wird, überhaupt im Alter, nachlässiger"
+--> SSK0 - Semantischer Subkommentar 0. Stufe
+--> PAR - Polysemieangabe (roemisch) "II."
+--> PragSemA - Pragmatisch-Semantische Angabe
| +--> GrA - Grammatische Angabe
| | +--> UGrA - Unspezifizierte grammatische Angabe "Partikel"
| +--> BA - Bedeutungsangabe
| +--> BPA - Bedeutungsparaphrasenangabe "eigentlich"
+--> BeiGA - Beispielgruppenangabe
+--> BeiA - Beispielangabe "was willst du überhaupt hier?"
+--> BeiA - Beispielangabe "wie ist das überhaupt passiert?"
+--> BeiA - Beispielangabe "du könntest überhaupt etwas freundlicher sein"
+--> GZB - Glossat zur Bedeutung "ruhig"

```

```

-----
Entry #1: Error caused by an invalid counter value
LEXPARSE: (W-403) Warning occurred in module 'Parser':
Probably incorrect value for counter detected:
"5";
'4'

```

```

Inputfile.....: data\counter_error.dat
First entry.....: 1
Last entry.....: 1
Entries parsed...: 1
Errors.....: 1
Time of day.....: Fri Nov 05 12:28:43 1993
Time of start....: Fri Nov 05 12:28:42 1993
Elapsed time.....: 00h 00m 01s

```

```

LEXPARSE V1.10 [Nov 05 1993] <DEBUG-SRC>
Program terminated successfully

```