

A Development Environment for an MTT-Based Sentence Generator

Bernd Bohnet, Andreas Langjahr and Leo Wanner
Computer Science Department
University of Stuttgart
Breitwiesenstr. 20-22
70565 Stuttgart, Germany
{bohnet|langjahr|wanner}@informatik.uni-stuttgart.de

1 Introduction

With the rising standard of the state of the art in text generation and the increase of the number of practical generation applications, it becomes more and more important to provide means for the maintenance of the generator, i.e. its extension, modification, and monitoring by grammarians who are not familiar with its internals. However, only a few sentence and text generators developed to date actually provide these means. One of these generators is KPML (Batesman, 1997). KPML comes with a Development Environment and there is no doubt about the contribution of this environment to the popularity of the systemic approach in generation.

In the generation project at Stuttgart, the realization of a high quality development environment (henceforth, DE) has been a central topic from the beginning. The DE provides support to the user with respect to writing, modifying, testing, and debugging of (i) grammar rules, (ii) lexical information, and (iii) linguistic structures at different levels of abstraction. Furthermore, it automatically generalizes the organization of the lexica and the grammar. In what follows, we briefly describe DE's main features. The theoretical linguistic background of the DE is the *Meaning-Text Theory* (Mel'čuk, 1988; Polguère, 1998). However, its introduction is beyond the scope of this note; the interested reader is asked to consult the above references as well as further literature on the use of MTT in text generation—for instance, (Iordanskaja *et al.*, 1992; Lavoie & Rambow, 1997; Coch, 1997).

2 Global View on the DE

In MTT, seven levels (or strata) of linguistic description are distinguished, of which five are relevant for generation: semantic (Sem), deep-syntactic (DSynt), surface-syntactic (SSynt), deep-morphological (DMorph) and surface-morphological (SMorph). In order to be able to generate starting from the data in a data base, we introduce an additional, the conceptual (Con) stratum. The input structure to DE is thus a conceptual structure (ConStr) derived from the data in the DB. The generation process consists of a series of structure mappings between adjacent strata until the SMorph stratum is reached. At the SMorph stratum, the structure is a string of linearized word forms.

The central module of the DE is a compiler that maps a structure specified at one of the five first of the above strata on a structure at the adjacent stratum. To support the user in the examination of the internal information gathered during the processing of a structure, a debugger and an inspector are available. The user can interact with the compiler either via a graphic interface or via a text command interface. For the maintenance of the grammar, of the lexica and of the linguistic structures, the DE possesses separate editors: a rule editor, a lexicon editor, and a structure editor.

2.1 The Rule Editor

The Rules. Most of the grammatical rules in an MTT-based generator are two-level rules. A two-level rule establishes a correspondence

between minimal structures of two adjacent strata. Given that in generation five of MTT's strata are used, four sets of two-level rules are available: (1) Sem \Rightarrow DSynt-rules, (2) DSynt \Rightarrow SSynt-rules, (3) SSynt \Rightarrow DMorph rules, and (4) DMorph \Rightarrow SMorph-rules.

Formally, a two-level rule is defined by the quintuple (\mathcal{L} , $\mathcal{C}xt$, $\mathcal{C}onds$, \mathcal{R} , $\mathcal{C}orr$). \mathcal{L} specifies the lefthand side of the rule—a minimal source substructure that is mapped by the rule onto its destination structure specified in \mathcal{R} , the righthand side of the rule. $\mathcal{C}xt$ specifies the wider context of the lefthand side in the input structure (note that by far not all rules contain context information). $\mathcal{C}onds$ specifies the conditions that must be satisfied for the rule to be applicable to an input substructure matched by \mathcal{L} . $\mathcal{C}orr$ specifies the correspondence between the individual nodes of the lefthand side and the righthand side structures.

Consider a typical Sem \Rightarrow DSynt-rule, which maps the semantic relation '1' that holds between a property and an entity that possesses this property onto the deep-syntactic relation ATTR. The names beginning with a '?' are variables. The condition 'Lex::(Sem::(?Xsem.sem).lex).cat = adj' requires that the lexicalization of the property is an adjective. '?Xsem \Leftrightarrow ?Xdsynt' and '?Ysem \Leftrightarrow ?Ydsynt' mean that the semantic node ?Xsem is expressed at the deep-syntactic stratum by ?Xdsynt, and ?Ysem by ?Ydsynt.

```
property(Sem_DSynt) {
  leftside:
    ?Xsem -1 $\rightarrow$  ?Ysem
  conditions:
    Sem::?Xsem.sem.type = property
    Lex::(Sem::(?Xsem.sem).lex).cat = adj
  rightside:
    ?Xds
    ?Yds
    ?Yds -ATTR $\rightarrow$ ?Xds
  correspondence:
    ?Xsem  $\Leftrightarrow$  ?Xds
    ?Ysem  $\Leftrightarrow$  ?Yds}
```

The rule editor (RE) has two main functions: (i) to support the maintenance (i.e. editing and examination) of grammatical rules, and (ii) to

optimize the organization of the grammar by automatic detection of common parts in several rules and their extraction into abstract 'class' rules. The theoretical background and the procedure of rule generalization is described in detail in (Wanner & Bohnet, submitted) and will hence not be discussed in this note.

While editing a rule, the developer has the standard commands at his/her disposal. Rules can be edited either in a text rule editor or via a graphic interface. Obviously incorrect rules can be detected during the syntax and the semantic rule checks. The syntax check examines the correctness of the notation of the statements in a rule (i.e. of variables, relations, conditions, etc.)—in the same way as a conventional compiler does. The semantic check examines the consistency of the conditions, relations, and attribute-feature pairs in a rule, the presence of an attribute's value in the set of values that are available to this attribute, etc. If, for instance in the above rule 'adj' is misspelled as 'adk' or erroneously a multiple correspondence between ?Yds and ?Xsem and ?Ysem is introduced, the rule editor draws the developer's attention to the respective error (see Figure 1).

Rule Testing. Rule testing is usually a very time consuming procedure. This is so partly because the generator needs to be started as a whole again and again, partly because the resulting structure and the trace must be carefully inspected in order to find out whether the rule in question fired and if it did not fire why it did not. The DE attempts to minimize this effort. With 'drag and drop' the developer can select one or several rules and apply them onto an input structure (which can be presented either graphically or in a textual format; see below). When a rule dropped onto the structure fires, the affected parts of the input structure are made visually prominent, and the resulting output (sub)structure appears in the corresponding window of the structure editor. If a rule did not fire, the inspector indicates which conditions of the rule in question were not satisfied. See also below the description of the features of the inspector.

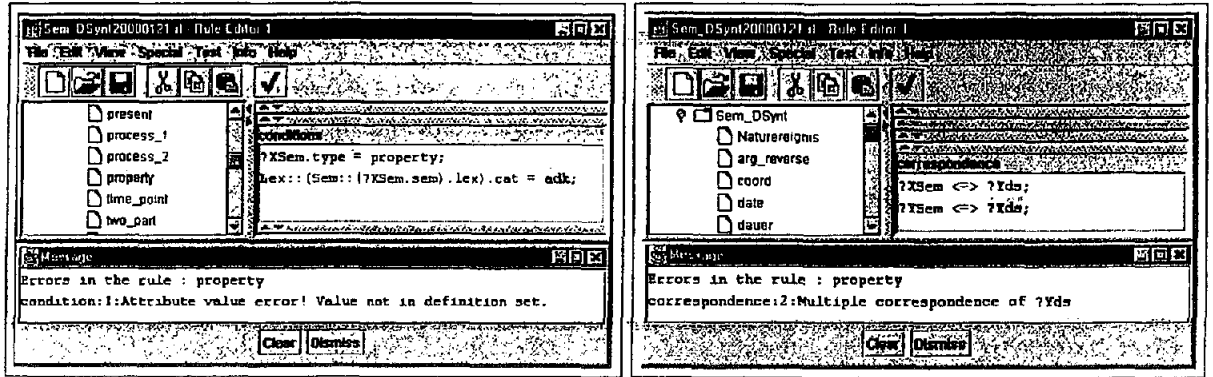


Figure 1: Error messages of the rule editor

2.2 The Structure Editor

The structure editor manages two types of windows: windows in which the input structures are presented and edited, and windows in which the resulting structures are presented. Both types of windows can be run in a text and in a graphic mode. The input structures can be edited in both modes, i.e., new nodes and new relations can be introduced, attribute-value pairs associated with the nodes can be changed, etc.

In the same way as rules, structures can be checked with respect to their syntax and semantics. Each structure can be exported into postscript files and thus conveniently be printed.

2.3 The Lexicon Editor

The main function of the lexicon editor is to support the maintenance of the lexica. Several types of lexica are distinguished: conceptual lexica, semantic lexica and lexico-syntactic lexica.

Besides the standard editor functions, the lexicon editor provides the following options: (i) sorting of the entries (either alphabetically or according to such criteria as 'category'); (ii) syntax check; (iii) finding information that is common to several entries and extracting it into abstract entries (the result is a hierarchical organization of the resource). During the demonstration, each of these options will be shown in action.

2.4 The Inspector

The inspector fulfils mainly three functions. First, it presents information collected during the application of the rules selected by the developer to an input structure. This informa-

tion is especially useful for generation experts who are familiar with the internal processing. It concerns (i) the correspondences established between nodes of the input structure and nodes of the resulting structure, (ii) the instantiation of the variables of those rules that are applied together to the input structure in question, and (iii) the trace of all operations performed by the compiler during the application of the rules.

Second, it indicates to which part of the input structure a specific rule is applicable and what its result at the destination side is. Third, it indicates which rules failed and why. The second and third kind of information is useful not only for generation experts, but also for grammarians with a pure linguistic background.

Figure 2 shows a snapshot of the inspector editor interface. Sets of rules that can simultaneously be applied together to an input structure without causing conflicts are grouped during processing into so-called *clusters*. At the left side of the picture, we see two such clusters (Cluster 13 and Cluster 22). The instances of the rules of Cluster 13 are shown to the right of the cluster pane. The cluster pane also contains sets of rules that failed (in the picture, the corresponding icon is not expanded). The left graph in Figure 2 is the input structure to which the rules are applied. For illustration, one of the rules, namely *date*, has been selected for application: the highlighted arcs and nodes of the input structure are the part to which *date* is applicable. The result of its application is the tree at the right. Beneath the graphical structures, we see the correspondence between input nodes

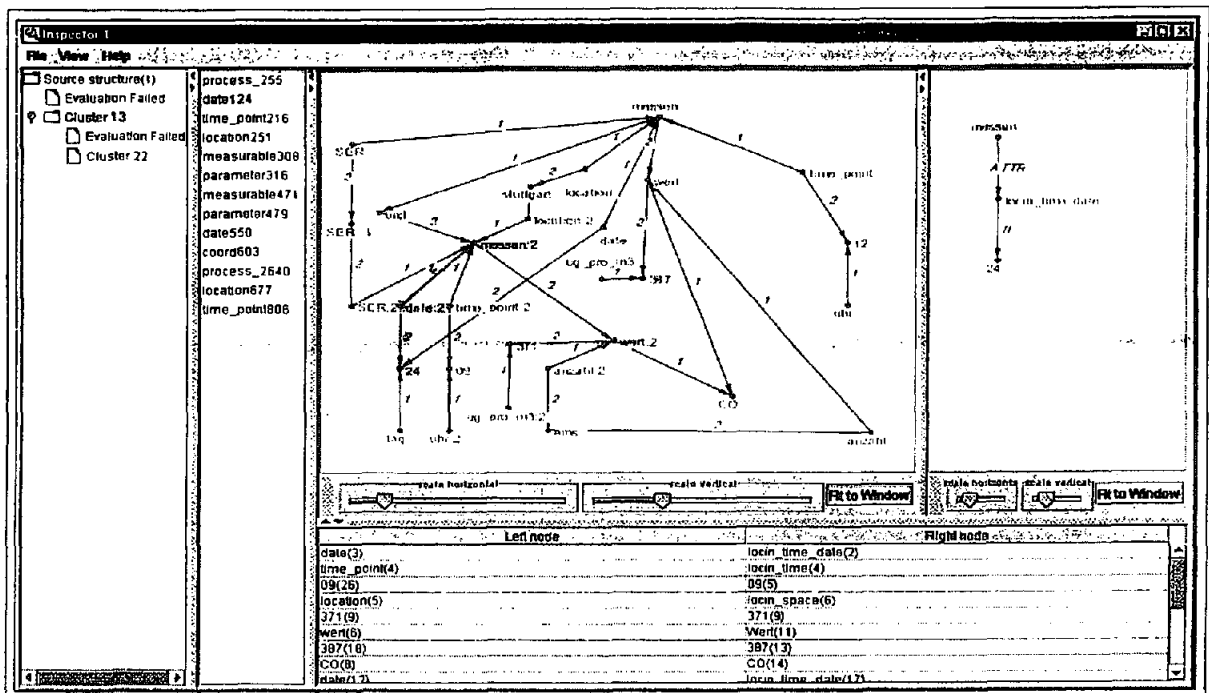


Figure 2: The inspector interface of the DE.

and result nodes. The numbers in parentheses are for system use.

2.5 The Debugger

In the rule editor, break points within individual rules can be set. When the compiler reaches a break point it stops and enters the debugger. In the debugger, the developer can execute the rules statement by statement. As in the inspector, the execution trace, the variable instantiation and node correspondences can be examined. During the demonstration, the function of the debugger will be shown in action.

3 Current Work

DE is written in Java 1.2 and has been tested on a SUN workstation and on a PC pentium with 300 MHz and 128 MB of RAM.

Currently, the described functions of the DE are consolidated and extended by new features. The most important of these features are the *import* and the *export* feature. The import feature allows for a transformation of grammatical rules and lexical information encoded in a different format into the format used by our generator. Tests are being carried out with the import of RealPro (Lavoie & Rambow, 1997) grammatical rules and lexical information (in particular

subcategorization and diathesis information) encoded in the DATR-formalism. The export feature allows for a transformation of the rules and lexical information encoded in our format into external formats.

Bibliography

- Bateman, J.A. 1997. Enabling technology for multilingual natural language generation: the KPML development environment. *Natural Language Engineering*. 3:2:15-55.
- Coch, J. 1997. Quand l'ordinateur prend la plume : la génération de textes. *Document Numérique*. 1.3.
- Iordanskaja, L.N., M. Kim, R. Kittredge, B. Lavoie & A. Polguère. 1992. Generation of Extended Bilingual Statistical Reports. *COLING-92*. 1019-1022. Nantes.
- Lavoie, B. & O. Rambow. 1997. A fast and portable realizer for text generation systems. *Proceedings of the Fifth Conference on Applied Natural Language Processing*. Washington, DC.
- Mel'čuk, I.A. 1988. *Dependency Syntax: Theory and Practice*. Albany: State University of New York Press.
- Polguère, A. 1998. La théorie sens-texte. *Dialangue*. 8-9:9-30.
- Wanner, L. & B. Bohnet. submitted. Inheritance in the MTF-grammar.