

MULTILINT

MULTILINT is a research and development project, sponsored by the German Ministry of Economy. The project partners are BMW AG (Bayerische Motorenwerke) and IAI (Institute for Applied Information Science at the University of the Saarland).

The project is developing a system in the frame of which an integrated form of supporting multilingual document production and management is realized and within which each linguistic component is used in the production chain as early as possible and as efficiently as possible.

1 Technical Documentation and Language Engineering

By means of integrating linguistic intelligence into the processing of technical documents, a considerable rate of improvement can be achieved: an automatic spell, syntax and style checking of the text and checking of terminological consistency will ensure a reliable information retrieval as well as the generation of corresponding notions in the foreign language and thus international compatibility.

Most of the MT systems available on the market, be they PC-based versions like GlobaLink, Trados-TM or the IBM-Translation Manager or 'large' MT systems like SYSTRAN, LOGOS, METAL, provide more or less sophisticated interfaces to the components supporting document production. However, the integration of such translation tools and other components into the user's environment as well as the linking of all tools as required by the user are not realized in a satisfying way by any of the existing systems, although there exist efforts offering better integrated solutions, such as within the ESPRIT project 'Translators Workbench' or within the EUROLANG optimizer, where at least a common user interface for all functions has been achieved. Rank Xerox is working on Compass and TANS, translation aids based on linguistic tools and architectures.

Although most of these systems dispose of a large amount of linguistic resources (built over years or decades), none of them uses linguistic intelligence optimally which is available to some extent either within or outside the system. Above all they don't account for the author's support, i.e. the person producing the text is neither supported wrt. an error-free document in his mother tongue nor is he supported in connecting expressions, terms etc. (e.g. synonyms, hypernyms, hyponyms, domain specific terms, ...).

The project described in this article is situated in the domain of multilingual document processing (including the domain of document translation) within the specialized field of automotive technical maintenance and support services. It is based on results of more than 10 years work in multilingual language processing in the European Union (EU) carried out within projects such as EUROTRA, MLAP and LRE. The theoretical findings of these projects have been used to build several stable

research prototypes for a broad range of linguistic tasks which are now ripe to be validated in an industrial context and to be further developed for their deployment in mobile information systems.

2 The Pilot User and His Needs

The pilot application site BMW has already made some efforts in the domain of document production and document management: they developed an authoring and documentation management system RS which allows the author only to use standardized processes when processing a document. There are about 80 to 100 persons of the central technical service department working daily with the system producing about 600 pages per year and person. Most of the authors are engineers and/or technical writers with good knowledge of text processing and layout; however, about 25% are beginners.

The system which is used at BMW (i.e. TIS = Technisches Informationssystem) allows the most important service related processes to be coordinated and aims at providing the BMW dealer organization with voluminous data bundled in a functional context by means of a CD-ROM. At the moment there are about 3500 systems installed all over the world. There exists a CD for German, English, French, Italian, Dutch, Swedish, Spanish, US English and Japanese which is up-dated four times a year.

The TIS system deals with the storage and the administration of technical information within document handling. A first application domain is determined on the basis of a corpus consisting of service information documents including repair messages and measures which must be available in German and English.

A special application within TIS is the Hotline Support System HUS. Problems occurring with cars which are reported to BMW by appointed dealers are dealt with by means of the HUS system. On the basis of type-specific data and problems Hotline Information (HLI) already stored contributing to the solution of the problem is searched for and reported to the dealer. In case there is no HLI (hotline information) linked to a specific problem, it is checked whether the problem is already described within the HUS system, if not, the problem will be added to HUS and a corresponding fullfledged HLI will be produced and entered into the database of the system, then being available for all users.

The user needs emerging from this scenario can be characterized as follows:

In order to achieve high quality documentation in the source language, which is a prerequisite for high quality texts in the target language, the user, i.e. the technical author, needs tools allowing the text to be controlled. This control should include grammatical, stylistic and terminological checkings, thus allowing the user to process the documents optimally, to structure them and to represent their content unambiguously. This leads to an optimization of the subsequent translation process in terms of time, costs and quantity.

By integrating the MULTILINT system into the specific application environment of TIS the goal of which is having access to the above described messages in various languages within reasonable time, it will also be integrated in the overall environment of the authoring system RS, thus in addition supporting other text processing activities, guiding and advising the user by means of a user-friendly front-end.

3 Components of MULTILINT

In the following chapters, the components of MULTILINT are described in more detail; currently, they all exist in a prototypical form. Experimental use by selected technical authors at the pilot

site is foreseen for the next months.

3.1 User Interface Module SYSCON

The user communicates with the tools via a common interface for which a data exchange format has been specified. The current prototypical version of the user interface has been realized in Tcl/Tk, thus allowing network access to the tools on UNIX workstations which are the platform used at the pilot site; a portation to JAVA/html (runnable also on PCs) is underway.

The user interface module guides the user by suggesting the functions to be activated next on the basis of the output of an accomplished task. Every tool can also be used on its own, either by a human user or as a particular function within a technical information system (integration in a specific environment of a user).

The linguistic resources are shared among all tools in order to ensure completeness, coherence and flexibility across linguistic applications.

Another important task of this component will be to examine and to realize the communication between the modules of the authoring system on the one hand, and the linguistic tools on the other.

To ensure optimal functionality, the linguistic tools are parameterized with respect to various functions. These parameters may either be pre-defined or may be chosen by the user:

- text format, e.g. SGML as used by the pilot user system TIS
- source and target language, e.g. German, English, etc.
- domain, e.g. sale, construction, logistics etc. This is organized in a hierarchical way.
- intended function for bilingual or monolingual treatment, e.g.
 - spell, syntax and style checking
 - indexation (storage of an expression, e.g. abbreviations and their meaning, main topics of the document etc.)
 - mono- or bilingual retrieval
 - terminology checking
 - automatic generation of glossary and production of additional material for the use of translation memories and terminological data bases:

As mentioned, the user can choose the source and target language out of a fixed set of languages. If he does not determine the language pair, the system will set the source language which has been determined by default in his user profile. The choice of the source language is checked by means of the morphological component. If the chosen language is not compatible with the text specified, then the language which has been identified by the morphological analysis will be chosen, and this decision will be indicated. The same holds for the choice of domain.

The morphological analysis and the determination of the technical domain of every piece of text are part of the automatic pre-processing, which also provides the user with a list of words not found in the database percentage. Commercial systems, too, have this function or a modified version of it. Additionally, the pre-processing function described here, gives information (interactively) about the number of words which do not belong to the determined or automatically identified domain. Information about the average length of sentences is provided as well, since this is of importance wrt. the translatability of the text.

3.2 Linguistic Resources

3.2.1 Morphological Analysis

The basic material for the morphological processing of texts already existed in various European languages: first, there is the system MPRO which decomposes the words and which provides their analysis and synthesis wrt. inflection, derivation and composition, and second, there exists a complete list of morphemes, for German, e.g. about 20.000, covering almost any text (MAAS - Mpro). These components have been evaluated with the selected domain specific texts and completed according to the needs of the specific application. For German an example of the automotive sector is given below:

```
>Seitenfenster
entry = {string='Seitenfenster',lex=fenster,clex= {lex=seite,clex=nil,
hyper= {f1=loc,f2=rand},head= {cat=n,deriv=nil,pref=nil,ehead= {gen=fem}}},
lemma=seitenfenster,hyper= {f1=loc,f2=way},head= {cat=n,deriv=nil,pref=nil,
ehead= {gen=neuter},ehead=({case=({nom};{dat};{acc}),num=sing};{case=({nom};
{gen};{acc}),num=plu})},graphics=cap}. [].

-- Surface structure: --
[compound,n,_1619,[hyper=loc,more=way,up=an],fenster,[n,[ns,n,a,seite,
seite,_1582,_1580,_1578,(0;n)],[no,_1567,_1565,[hyper=loc,more=rand],[[0,0,0],
[n,0]],fem,no]],_1531,seite]]

-- Semantic structure: --
[np,[n,[hyper=loc,more=way,up=an],_2984,fenster,_2922,_2924],[fuehren,[n,
[hyper=loc,more=way,up=an],c0,fenster,_2922,_2924],[n,[hyper=loc,_2912=
rand|_2910],wohin,seite,sg,def],no]]

=== Interpretation: ===
Ein Seiten|fenster ist ein Fenster, das zur Seite fuhrt .
```

The concatenation of multi word units (e.g. in Bezug auf) and of terminological expressions (e.g. air intake grill) is possible. In addition, we have taken into account abbreviations and acronyms which occur quite often in these specialized texts. Also complex compounds have been accounted for in a special way.

The results provided by this component are also used for the generation of detailed concordances (e.g. extracting pieces of text showing a word of a given syntactic category).

3.2.2 Syntactic Analysis

The syntactic processing, consisting of the description of grammatically correct structures as well as the transduction rules for interpreting these rules of analysis and synthesis, already available for several languages, has been evaluated and completed on the basis of the selected corpus. Special attention has been paid to the improvement of their efficiency. In the application domain foreseen (automotive sector, especially service information) e.g. formulations in telegraphic style which are characterized by verbless constructions, subjectless verbal groups and determinerless nominal groups, have been treated.

The syntactic analysis is organized in several steps; some homographies are reduced by simple 'neighborhood' tables and by the completion of a flat syntactic analysis which determines noun

groups and some syntactic chunks of the sentence. Special analysis rules which take frequent errors into account work with relaxed constraints for the indication of syntactic mistakes; a complete and unification based syntactic analysis is used during the attempt to obtain better translations. As this component makes frequent use of different readings and semantic constraints, such processes cannot be carried out online.

3.3 Linguistic Tools

As mentioned above, the NLP components described are being applied for various purposes:

- for checking spelling, syntax, style and terminology
- for checking the translatability of the text
- within the proper translation process, both in analysis and in the synthesis of the target language translation

- **Source text control**

This tool is a device which supports the author during the production of a document in the source language wrt. linguistic and terminological correctness. The source language text has to fulfill certain linguistic criteria - it would not make sense to support the production of a text version in a foreign language if it is done on the basis of a defective source language version. These criteria include orthographic, terminological and syntactic correctness, the avoidance of long sentences and of meaning ambiguities etc. They are important for a better understanding of the source language text. This type of monitoring function in text production is partially available with existing text systems which are endowed with spell checkers as well as programs checking syntax and style. The respective components in MULTILINT can however make use of an exhaustive linguistic analysis and of special features developed together with technical authors and developers of the user site systems. A few examples should illustrate the type of messages issued when inadequate input is found:

- 'Normal' spelling checking:

Die Kaefigmutter der Sitzschiene ist lohse.

Warning:

"lohse: unknown word".

- Syntax checking:

Die Rohbautoleranzen sind eingehalten nicht worden.

Warning:

In this structure the negation element should occur after the finite verb".

Die Aenderung wird auch an der anderen Spaltdichtungen durchgefuehrt.

Warning:

The noun and its determiners should agree in number".

Es ist der zeit in Bearbeitung.

Warning:

A noun begins with a capital letter".

Die Kaefigmutter der Sitzschiene sind lose.

Warning:

Subject and verb do not agree in number".

Dichtlippe an an der oberen Frontscheibenabdeckung krepelt ein.

Warning:

unknown structure".

– Technical writing:

Kuhlwasserthermostatgehalusedeckel

Warning:

Dieses Kompositum besteht aus zu vielen Gliedern (max. 4).

oder

Vermeiden Sie Wortmonster; verwenden Sie einen Bindestrich!

Auch in den darunterliegenden Geschwindigkeitsklassen ist man aus Sicherheitsgründen zurüchhaltender geworden.

In diesem Fall erhält man eine Unbedenklichkeitsbescheinigung des Fahrzeugherstellers.

Warning:

Unpersönliche Konstruktionen mit 'man' vermeiden.

oder

Vermeiden Sie Sätze mit 'man'!

Eine erneute Aktivierung kann nur in der Maske 'Stimullgeraite', wie oben beschrieben, erfolgen.

Warning:

Auf Funktionsverbgefüge verzichten.

oder

Vermeiden Sie Sätze mit Verben, die nicht mehr ihre eigentliche Bedeutung haben.

Den Innenraumschutz bietet ein Ultraschallmodul.

Warning:

Das Subjekt steht nicht im Vorfeld.

oder

Das Subjekt sollte vorne stehen!

Die Wulstkabel , verkupferte oder vermessingte und mit Gummi ummantelte Stahlkabel, sorgen fuir festen Sitz des Reifens auf der Felge.

Warning:

Zu viele attributierende Elemente.

oder

Vermeiden Sie Attributhäufungen!

• **Foreign language glossary**

This device produces on the basis of the input text a domain specific glossary in the target language which is similar to a morpho-lexical translation, e.g.

Abdeckleiste	außen	Seitenfenster	hinten	deformiert
⇓	⇓	⇓	⇓	⇓
cover strip	outside	side window	rear side	distorted

The user may also determine the lexical categories (noun, adjective, verb) for which he wants a glossary, and whether the glossary should be extended by derivational knowledge: the user may ask for instance for the translation of the corresponding verbal stem of an unknown noun.

German compounds ('Leder-lenk-rad') and multiword units in English ('brake ventilation adapter') are looked up in the respective terminology database. If they are not found in the database they are analyzed and translated compositionally.

If it is not recommended that the machine translation of an input sequence be performed after the pre-processing or if the user does not want an automatic translation at all, he may activate the FGLOSS function instead.

• **Translation memory control – TMCON**

As already mentioned above BMW is using the translation memory system TRANSIT, a tool by means of which it is checked whether parts of the text to be translated have already been translated before. Prerequisite for this function is that there are a lot of parallelized texts available in a database.

Due to the integration of linguistic tools, as envisaged in this project, it will be possible to search e.g. for lexical elements of morphologically and/or syntactically analyzed sentences in the database. In case there are too many matches wrt. the search text, the user has the further option of a complete linguistic analysis of the text, in order to specify priorities. Within this full analysis not only lexical items but also their semantic relations within the clause are checked wrt. their similarity.

• **Informative Machine translation – MT**

Finally, there is the possibility to generate on the basis of transfer rules and the syntactic processing in the target language possible translations of short sentences or parts of texts satisfying several requirements (e.g. short messages within the service information process). The following prerequisites should be fulfilled by the input text:

- it should not contain spelling or syntax errors
- it should not be too long, respectively easy to sub-divide
- there should be sufficient lexical material in the respective dictionaries

Texts fulfilling these requirements can be analyzed, transferred and generated in the target language, all this by means of all linguistic tools available. The result will be a grammatically correct sentence in the target language which should be a good basis for understanding and, if needed, for postediting.

If a complete translation is not recommendable (either on grounds of defective input or of text complexity), a quick translation can be done on the basis of morpho-lexical information having the same word order as the input sentence or a so-called phrase translation, e.g.

Deckel hinten zu tief, vorne zu hoch
[The caps] [back] [too deep] [,] [sein] [high] [vorne]

4 Conclusion: Testing and Evaluation

During the first year of the project, the linguistic resources and tools have been adapted to and tested with the service information material coming from the industrial pilot site. At the same time, long dialogues have been carried out between the developers and several groups of possible users of the MULTILINT system. They have expressed their wishes and preferences for their multiple purposes: technical writing of new pieces of service information in several languages, mono- and bilingual retrieval of these pieces, the control of terminological consistency for translation into new languages and the use within a future information network comprising several countries where business is growing fast. The current prototype is the result of this process and will undergo an exhaustive on-site testing during the year of 1997; in several evaluation loops, a gradual updating will take place in order to provide the best possible linguistic and ergonomic solutions for the different user groups. The project is accompanied by a working group whose members are all active in different German and other European car manufacturing companies; subject of this working group is the possibility of testing (quick) adaptability of MULTILINT to similar purposes in the same and other domains. The project is one of several attempts currently made in Europe to bring the results of long-term research in computational linguistics and language engineering closer to industrial application needs.

References

- [Adriaens & Macken (95)] Adriaens, G., 'Simplified English Grammar and Style Correction in an MT Framework: The LRE SECC Project', in: *Aslib Proceedings*, vol 47, no 3, März 1995, 73 - 82.
- [AECMA SE-Guide (95)] AECMA (eds.), 'AECMA Simplified English', AECMA Document: PSC-85-16598, AECMA, Bruxelles, 1995
- [Bolioli et al. (92)] Bolioli, A., Dini, L., Malnati, G., 'JDII: Parsing Italian with a Robust Constraint Grammar', in: *Proceedings of COLING 92*, Nantes, 1992, 1003 - 1007.
- [Chambers et al. (95)] Chambers, C., Malnati, G., Tesio, R., Soma, E., 'JDII - The Linguistic Syntax Checker', in: *Proceedings of Language Engineering '95*, Montpellier, 1995, 121 - 129.
- [Clémencin (96)] Clémencin, G., 'Integration of a CL-checker in an Operational SGML Authoring Environment: Methodological and Technical Issues', in: *Proceedings of CLAW 96*, KU Leuven, Leuven, 1996, 32 - 42.
- [Fottner-Top (96)] Fottner-Top, C., 'Workshop: Erstellung von verständlicher und benutzerfreundlicher technischer Dokumentation', München: Institut für technische Literatur, 1996.
- [Gingras (87)] Gingras, B., 'Simplified English in Maintenance Manuals', in: *Technical Communication*, First Quarter 1987, 1987, 24 - 28.
- [Hayes (1994)] Hayes, P. J. 'Automatic Machine Translation of Technical Information without Post-Editing', Pittsburg: Carnegie Group, 1994.
- [Heuler (93)] Heuler, M., 'Kontrollierte Sprache', in: *Tekom Nachrichten*, 4/93, 1993, 13 - 15.
- [Hoard et al. (92)] Hoard, J.E., Wojcik, R., Holzhauser, K., 'An Automated Grammar and Style Checker for Writers of Simplified English', in: O'Brian Holt, P., Williams, N., *Computers and Writing: State of the Art*, Kluwer Academic Publishers, Dordrecht/Boston/London, 1992, 278 - 296.
- [Kincaid et al. (91)] Kincaid, J.P., Kincaid, C.D., Kniffin, J.D., Thomas, M., 'Intelligent Authoring Aids for Technical Instructional Materials Written in Controlled English', in: *Journal of Artificial Intelligence in Education*, Vol. 2 (3), Spring 1996, 77 - 81.
- [Maas (96)] Maas, D. 'MPRO-Dokumentation', in: *Hausser, R. Morpholympics-Dokumentation*, Springer 1996.
- [Rank Xerox] Rank Xerox, 'Fact Sheet Locolex/Compass and XeLDA/TANS', Internal Document, Rank Xerox Company 1996.
- [SDD-Broschüre] NN, 'Siemens Dokumentations Deutsch, Syntaxregeln für Technische Dokumentation', Internes Dokument, Siemens AG, München.
- [Schachtl (96)] Schachtl, S., 'Requirements for Controlled German in Industrial Applications', in: *Proceedings of CLAW 96*, KU Leuven, Leuven, 1996, 143 - 149.
- [Wojcik & Holmback (96)] Wojcik, R.H., Holmback, H., 'Getting a Controlled Language Off the Ground at Boeing', in: *Proceedings of CLAW 96*, KU Leuven, Leuven, 1996, 22 - 31.