

Towards Corpus Annotation Standards The MATE Workbench¹

Laila Dybkjær, Niels Ole Bernsen
Natural Interactive Systems Laboratory
Science Park 10, 5230 Odense M, Denmark
E-post: laila@nis.sdu.dk, nob@nis.sdu.dk

Abstract

This paper describes the European MATE project and its work towards speech corpus annotation standards. Focus is on the MATE workbench which is thoroughly described and illustrated.

1. Introduction

The aim of the European Telematics project MATE (Multilevel Annotation Tools Engineering) has been to facilitate the use and reuse of spoken language resources, coding schemes and tools by addressing theoretical issues and implementing practical solutions. MATE was launched in 1998 in response to the increasing need - not least in the area of spoken language dialogue systems research and development - for tools and standardisation efforts in support of efficient markup of spoken dialogue corpora at multiple levels.

The main results of the project are the MATE markup framework which bridges between the theoretical and the practical activities of MATE and is proposed as a standard for the definition and representation of markup for spoken dialogue corpora (Dybkjær and Bernsen 2000b), and the MATE workbench which supports the use of the markup framework. Even if the project has formally come to an end, the project consortium has continued funding for improving the workbench and new versions continue to appear. The newest version of the software can always be downloaded from the MATE web site at

¹ This paper is also being published in the Proceedings of the COCODA meeting in Beijing, 21 October 2000.

<http://mate.nis.sdu.dk>. Both an executable version and the source code is available, the latter under the GNU open source LGPL license.

A discussion forum has been set up at the MATE web site for asking questions and sharing experience on the workbench, and for adding new tools to the MATE workbench to enhance its functionality.

In the following, Section 2 briefly describes the theoretical approach of MATE, Section 3 provides a detailed walkthrough of the MATE workbench, and Section 4 mentions related work.

2. Theoretical Background

To provide a solid basis for the coding standard to be proposed by MATE, more than sixty existing coding schemes belonging to five different coding levels (i.e. prosody, (morpho-) syntax, co-reference, dialogue acts, and communication problems) and their cross-level interaction were reviewed (Klein et al. 1998). The collected information served as background for establishing the MATE markup framework (Dybkjær et al. 1998, Dybkjær and Bernsen 2000b). In MATE, a coding level is some level of abstraction at which to conceptualise, tag, analyse, and retrieve information inherent in language corpora.

The MATE markup framework is a conceptual model which basically prescribes (i) how files are structured, for instance to enable multi-level annotation, (ii) how tag sets are represented in terms of elements and attributes, and (iii) how to provide essential information on markup, semantics, coding purpose etc.

The central concept of the MATE markup framework is the coding module. A coding module is an extended kind of coding scheme which prescribes what constitutes a coding, including the representation of markup and the relations to other codings. Coding modules incorporate (i), (ii) and (iii) above. Thus, the MATE coding module is a proposal for a standard description of coding schemes.

For each of the above-mentioned five annotation levels and the issues to do with cross-level annotation, one or several of the reviewed coding schemes were adopted as starting points for the definition, following the MATE markup framework, of best practice coding schemes for implementation in the MATE workbench (Mengel et al. 2000).

3. The MATE Workbench

The MATE workbench (Isard et al. 1998, Isard et al. 2000, Dybkjær and Bernsen 2000a) is a software tool set which supports the MATE markup framework, incorporates the MATE best practice coding schemes, and enables users to annotate corpora and extract information about annotated corpora via a user-friendly interface. The workbench is in continued development and, although the full functionality and desired user-friendliness has not yet been achieved, there is already considerable interest among colleagues from around the world in using the MATE workbench.

The workbench is implemented in Java. It has been tested on Unix (Solaris) and Windows (NT and 98) but should run on any platform for which Java 1.2 or newer is available. The workbench has a modular architecture which facilitates the addition of new modules and new tool functionality by its users. XML is used for internal file representation. Style sheets are used for specifying the visual presentation of data to users. Style sheets are written in the MATE Style sheet Language (MSL). The emerging standard in this area is XSLT. XSLT, however, was not fully defined when the workbench was being designed, and lacked various necessary functionalities at the time. It was therefore decided to implement MSL which uses the MATE query language but is otherwise similar to XSLT.

In the following, we describe and illustrate the MATE workbench functionalities.

3.1. Getting started

As mentioned, the MATE workbench can be downloaded from the MATE web site at <http://mate.nis.sdu.dk>. The download page describes how to start the workbench from a command tool. Starting the workbench will result in the following two windows being opened. The first window (Figure 1) allows the user to open new project windows (Figure 2), access available tools, such as the coding module editor (Figure 5), and access the online help function.

The second window (Figure 2) shows example projects which include best practice coding schemes and annotated examples for a number of different annotation levels.

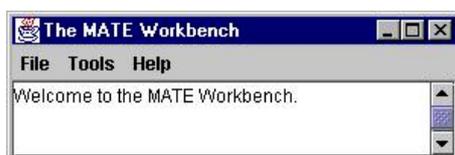


Figure 1. The control window.

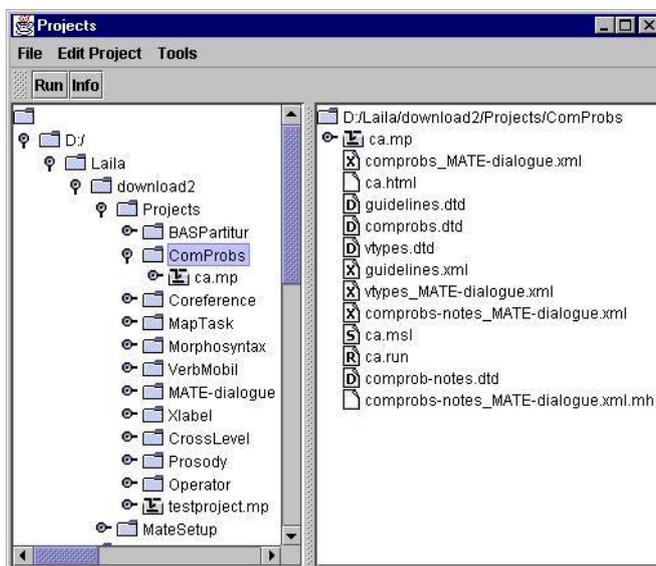


Figure 2. The projects window.

3.2. Annotation

The two main functionalities of the MATE workbench are support for annotation of spoken dialogue corpora and linguistic data, and information

extraction from annotated corpora. The workbench annotation support is described in the following.

3.2.1. Best practice coding schemes

Ready-to-use best practice coding schemes are included for most of the coding levels addressed by MATE. Prosody is not yet included. Style sheets allow the user to view and annotate files in an appropriate way using any of the implemented best practice schemes. Example dialogues are provided which illustrate how a tagged dialogue and the accompanying tag set will be shown when a particular best practice scheme is being applied. Internally in the workbench, an annotated dialogue is represented as a set of references to the transcription of the dialogue and possibly to other coding files. By default, the transcription refers to timeline information. The two windows in Figures 3 and 4 show annotation based on the MATE MapTask scheme and the MATE communication problems scheme, respectively. Figures 3 and 4 aptly illustrate the very different display requirements imposed by different coding schemes.

Clicking on the green ‘PLAY’ button in the Map Task window in Figure 3 will result in the audio file corresponding to the transcribed turn being played. The giver and the follower are the two speakers. Each speaker turn can be annotated with a speech act chosen from the list on the left by selecting the speaker and then clicking on the tag to be assigned. For example, the second giver utterance has been annotated with the ‘instruct’ dialogue act.

In the communication problems window in Figure 4, the orthographically transcribed dialogue is shown in the upper left-hand panel. To support the coder, guidelines for cooperative dialogue are shown in abbreviated form in the upper right-hand panel. Types of violation of particular guidelines are incrementally added by the coder in the lower right-hand panel. This panel is empty when annotation starts. The blue markup in the dialogue refers to the types of violation described in the lower right-hand panel and the violations themselves

refer to the guidelines. The lower left-hand panel shows annotator's notes. Again, this panel is empty when annotation starts. Notes can be added whenever the annotator needs to add an explanation of, e.g., why something went wrong in the dialogue so as to cause a communication problem.

There is no example coding module for transcription in the MATE workbench. Instead, a converter from Transcriber format (<http://www.etca.fr/CTA/gip/Projets/Transcriber/>) to MATE format enables transcriptions made using Transcriber to be exported to MATE format and annotated using the MATE workbench.

3.2.2. Adding a new coding module

Users may add new coding modules (coding schemes) themselves for existing or new coding levels via the coding module editor, cf. Figure 5. In order for a coding scheme and the dialogues annotated using it to be usable and understandable by people other than its creator, some key information must be provided. The MATE coding module which is the standard coding scheme description format proposed by MATE, serves to capture this information. A coding module includes the ten items shown in Figure 5. It prescribes what constitutes a coding, including markup representation and relations to other codings (module references).

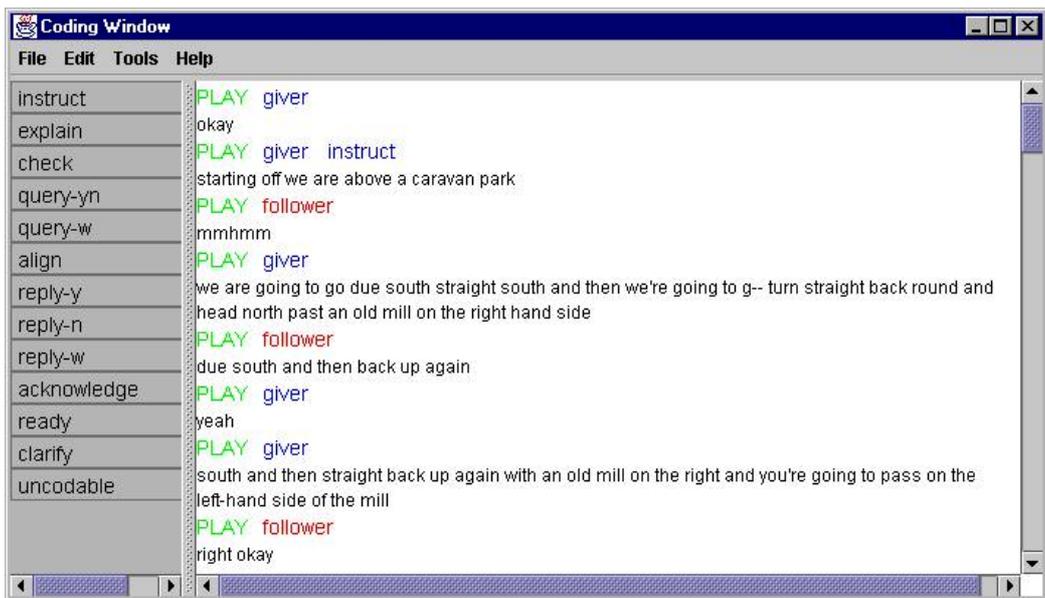


Figure 3. Annotation using the MATE Map Task scheme.

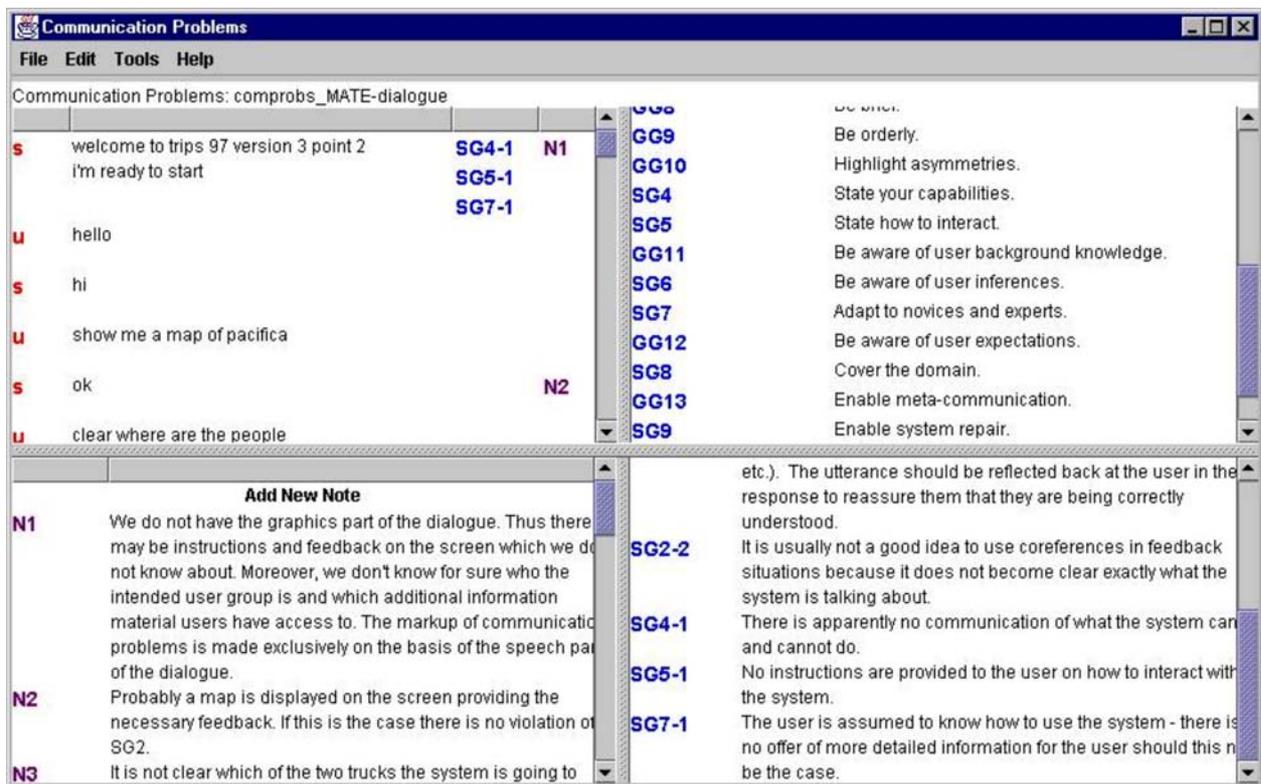


Figure 4. Annotation using the MATE communication problems scheme.

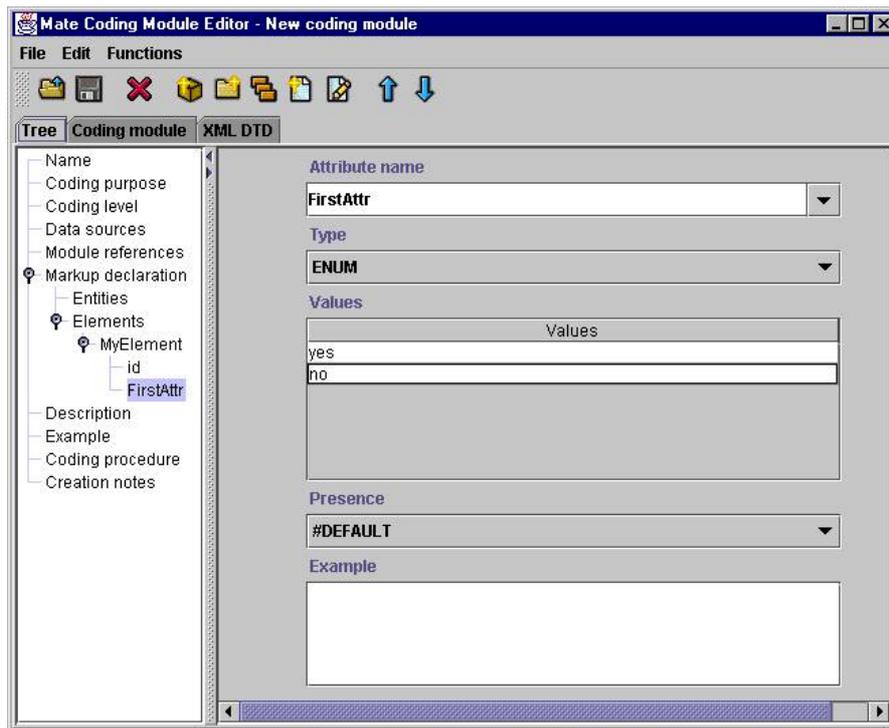


Figure 5. The coding module editor.

3.2.3. Adding a new project

A new MATE project folder can be created (Figure 6) via the 'File' menu in the project window in Figure 2.



Figure 6. Creating a new project.

3.2.4. Adding a new file

Figure 7 shows the browsing for an existing file to be added to a project folder. The window is invoked from the 'Edit Project' menu in the project window in Figure 2. New files in a particular format may be created

(Figure 8) via the same menu and files may also be removed via this menu. As a minimum, a new XML file which is being created must be given a name, and the DTD or coding module to be applied must be specified. However, it is recommended that further header information is provided as well, cf. Figure 9.

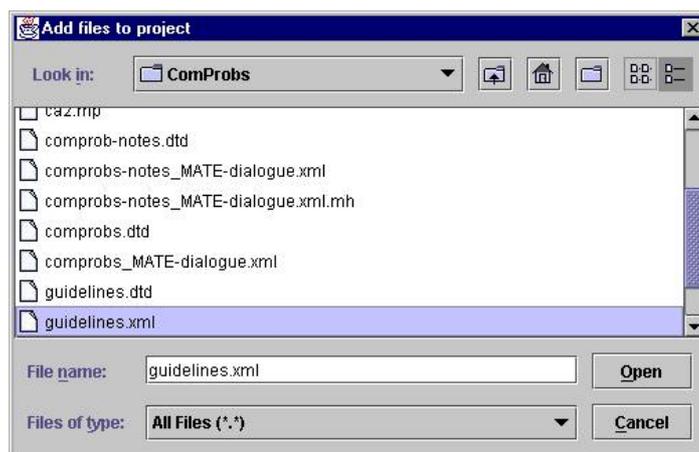


Figure 7. Browsing to add a file to a project folder.

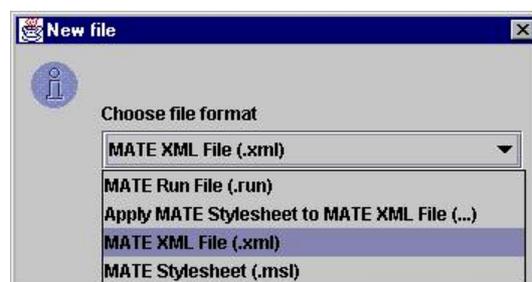


Figure 8. Creating a new file.

3.2.5. Listening to audio files

Audio files can be selected in the projects window (Figure 2) and played using the MATE audio tool. The window in Figure 10 will then appear displaying the sound file as a waveform. The audio tool might be used during transcription if a user has added an appropriate transcription coding

module. As it stands, the audio tool acts as a support tool during annotation and annotation review. For instance, when annotating communication problems there is sometimes a need for listening to the speech file in order to disambiguate an utterance and diagnose what went wrong.

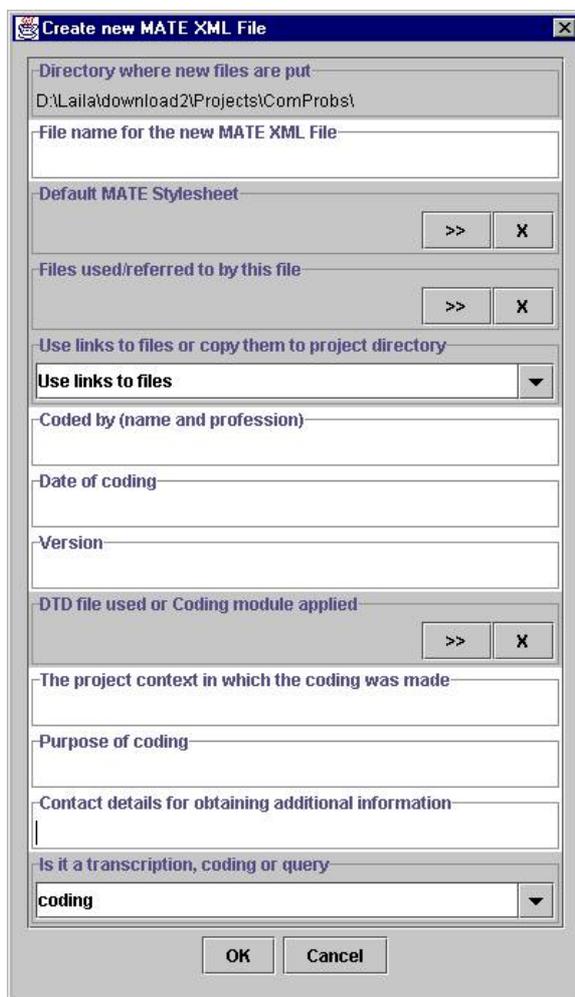


Figure 9. Creating a new XML file

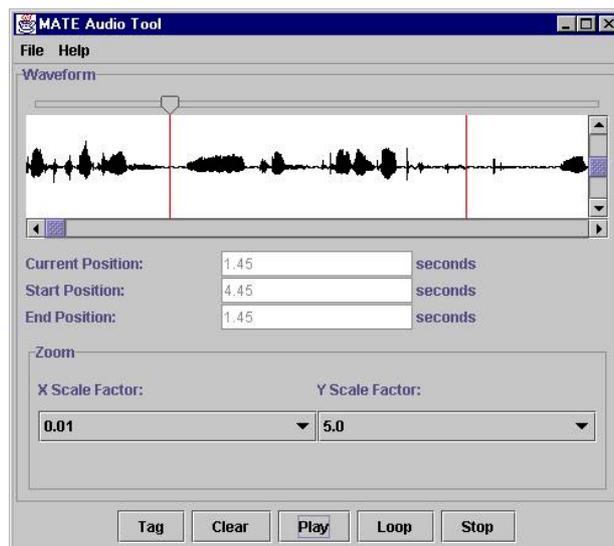


Figure 10. The MATE audio tool.

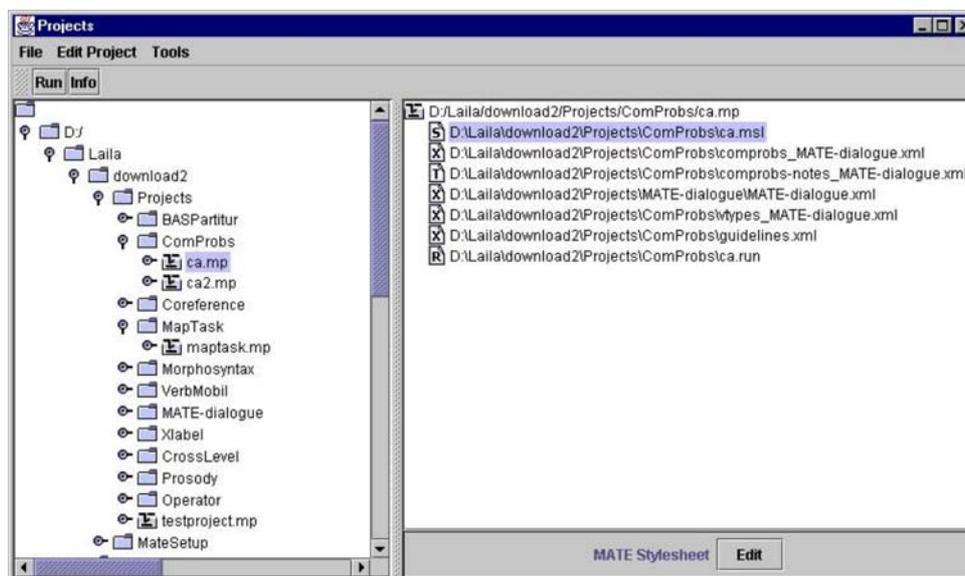


Figure 11. The edit option.

3.2.6. Editing a file

The MATE workbench uses three basic file types. Coding files are XML files. Style sheets are MSL files used for visualisation of codings and for annotation support. Runfiles specify which style sheet to apply to which XML file. Any XML, MSL or runfile in a MATE project folder can be

opened and edited (cf. Figures 12-14) by selecting the file and clicking on 'Edit' in the projects window, cf. Figure 11.

3.2.7. Import from and export to other file formats

Files may be imported from other formats to XML. For the moment, conversion from XLabels and BAS Partitur to XML is enabled. New converters can easily be added, including converters which export from XML to other formats, such as HTML. Figure 15 shows import from BAS Partitur.

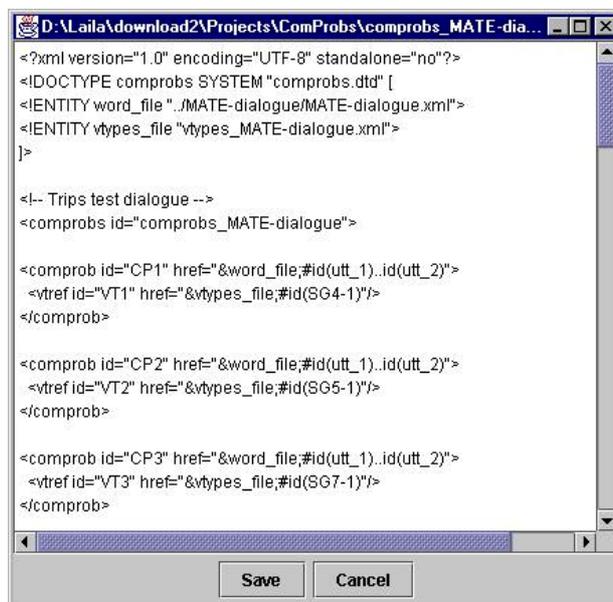


Figure 12. Editing an XML file.



Figure 13. Editing a style sheet.

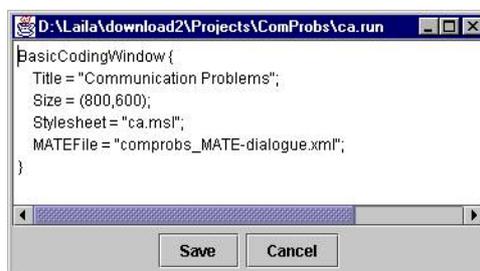


Figure 14. Editing a runfile.

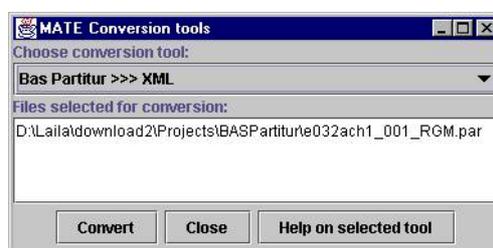


Figure 15. Import from BAS Partitur to XML.

3.3. Extracting information from annotated corpora

Once a corpus has been annotated, it must be possible to extract information from it for many different purposes. The MATE query tool is available for selecting the document(s) to be queried and for specifying the information to be subjected to a query, cf. Figure 16 (Isard et al. 2000). The query tool is activated from the tools menu in a coding window, cf. Figures

3 and 4. The information which can be extracted includes statistical information. Results are shown as illustrated in Figure 17. The interface for displaying results is not yet finalised. This is why the query results window for the moment only shows the raw XML data extracted from the queried XML file.

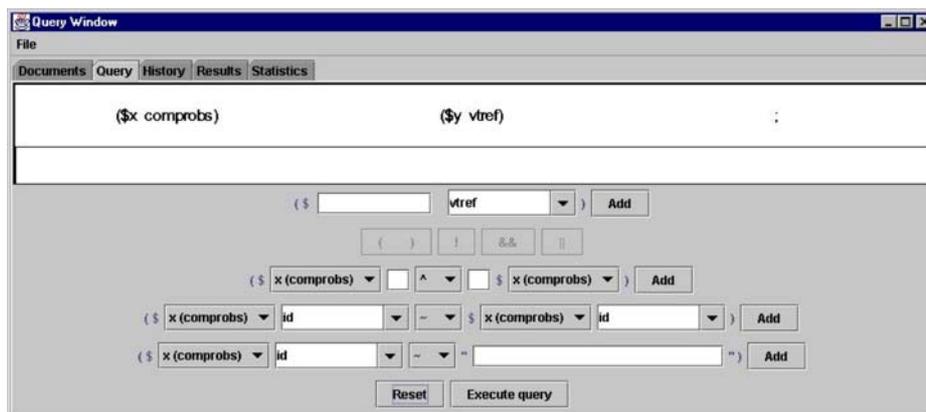


Figure 16. The query window.

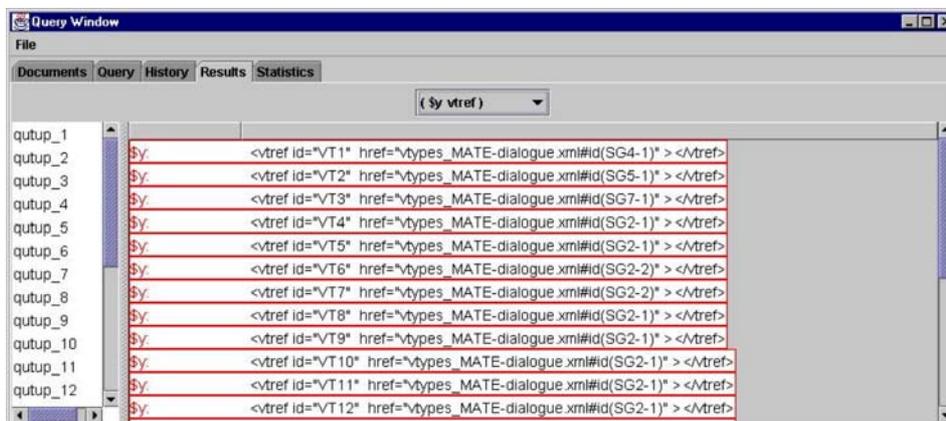


Figure 17. Results of a violation types query, cf. Figure 4.

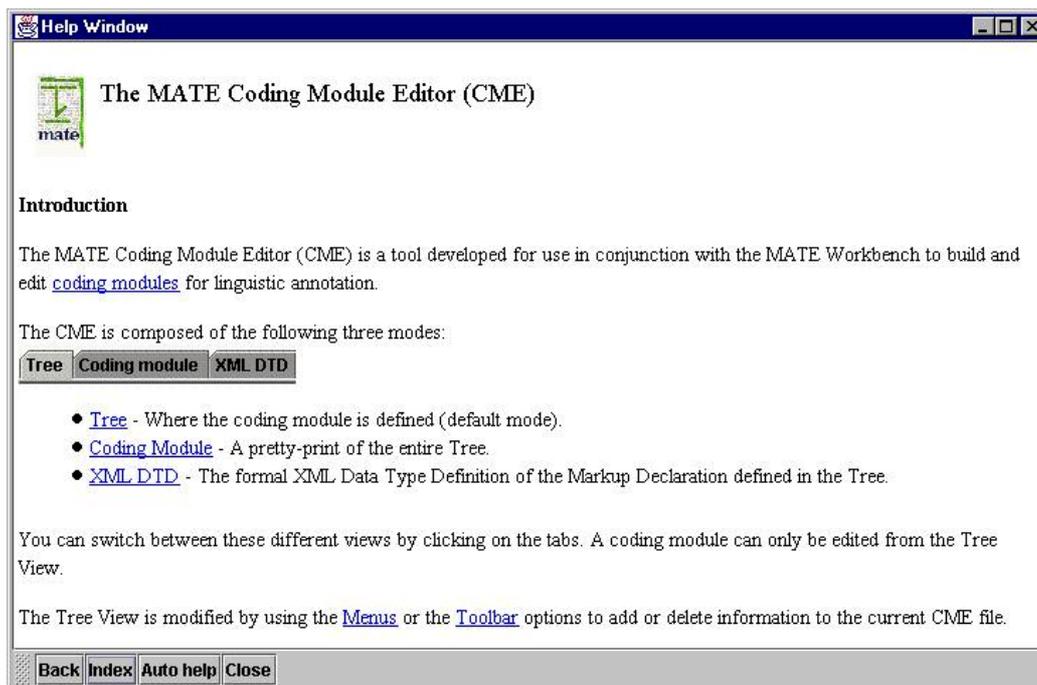


Figure 18. The help window.

3.4. Getting help

An online help facility may be consulted at any time during use of the MATE workbench. Figure 18 shows the topmost help page for the coding module editor. Explanation of particular coding schemes is available from the help menu in the coding window, cf. Figures 3 and 4.

3.5. Usability

Usability is a key concern in MATE. The focus on usability is reflected in the coding module concept and in the workbench coding module editor. The coding module prescribes a comprehensible standard description to be made of any coding scheme. The coding module editor makes it easy to specify coding modules – not least their markup declarations - also for non XML-literate users. The editor enables the user to specify the markup declaration for a new coding module almost without requiring any

knowledge of the underlying XML representation. The coding module editor automatically generates a DTD which is then used internally by the workbench. The coding module editor thus represents a major step forward compared to tools which require users to write DTDs.

There are still a couple of major usability issues to be solved, however. One issue is a user-friendly way of creating new coding visualisations. Writing style sheets for the workbench is cumbersome and requires programming skills because no editor is provided. The user must edit the raw style sheet code (or write new code), cf. Figure 13. It is high on our wishlist to enable users to easily define new visualisations. This may be done either by providing a style sheet editor comparable to the coding module editor as regards ease of use, or, alternatively, through a completely new interface concept replacing the need for style sheets and enabling users to easily define new visualisations.

A second major issue is the interface to the query tool and its results. As for the latter, Figure 17 makes it evident that usability improvements are needed. The query results could be presented far more transparently using an appropriate style sheet. This is also on the MATE to-do list. So is a more comprehensible interface for expressing queries than the present one (Figure 16).

4. State of the Art

Several frameworks for speech corpus annotation have been proposed but to our knowledge the MATE markup framework is still the more comprehensive framework around. An example of another framework is the annotation framework recently proposed by Bird and Liberman (1999)

which is based on annotation graphs. These are now being used in the ATLAS project (Bird et al. 2000) and in the Transcriber tool (Geoffrois et al. 2000). The annotation graphs serve as an intermediate representation layer between interface and internal data structures. Whilst Bird and Liberman do not consider coding modules or discuss the interface from a usability point of view, they present detailed considerations concerning timeline representation and timeline reference. The two frameworks may, indeed, turn out to complement each other nicely.

Acknowledgements

We gratefully acknowledge the support for the MATE project provided by the European Commission's Telematics/Language Engineering Programme. We would also like to thank all MATE partners. Without the very considerable joint efforts of the project consortium it would not have been possible to build the MATE workbench.

5. References

The MATE workbench is available in executable version and under the GNU open source LGPL license from the MATE web site at <http://mate.nis.sdu.dk>. MATE reports are also available from this web site.

Bird, S. and Liberman, M. 1999. *A Formal Framework for Linguistic Annotation*. Technical Report MS-CIS-99-01. Department of Computer and Information Science, University of Pennsylvania.

Bird, S., Day, D., Garofolo, J., Henderson, J., Laprun, C. and Liberman, M.: ATLAS 2000. A Flexible and Extensible Architecture for Linguistic

Annotation. In *Proceedings of the 2nd International Conference on Language Resources and Evaluation (LREC 2000)*. Athens:1699-1706.

Dybkjær, L. and Bernsen, N. O. 2000a. The MATE Workbench. *Proceedings of the LREC'2000 workshop on Data Architectures and Software Support for Large Corpora*, Athens: 33-37.

Dybkjær, L. and Bernsen, N. O. 2000b. The MATE Markup Framework. *Proceedings of the 1st SIGdial Workshop on Discourse and Dialogue*, Hong Kong.

Dybkjær, L., Bernsen, N. O., Dybkjær, H., McKelvie, D. and Mengel, A. 1998. *The MATE Markup Framework*. MATE Deliverable D1.2.

Geoffrois, E., Barras, C., Bird, S. and Wu, Z. 2000. Transcribing with Annotation Graphs. *Proceedings of the 2nd International Conference on Language Resources and Evaluation (LREC 2000)*, Athens:1517-1521.

Isard, A., McKelvie, D., Cappelli, B., Dybkjær, L., Evert, S., Fitschen, A., Heid, U., Kipp, M., Klein, M., Mengel, A., Møller, M. B. and Reithinger, N. 1998. *Specification of Workbench Architecture*. MATE Deliverable D3.1.

Isard, A., McKelvie, D., Mengel, A., Møller, M. B., Grosse, M. and Olsen, M. V. 2000. *Data Structures and APIs for the MATE Workbench*. MATE Deliverable D3.2.

Klein, M., Bernsen, N. O., Davies, S., Dybkjær, L., Garrido, J., Kasch, H., Mengel, A., Pirrelli, V., Poesio, M., Quazza, S. and Soria, C. 1998. *Supported Coding Schemes*. MATE Deliverable D1.1.

Mengel, A., Dybkjær, L., Garrido, J., Heid, U., Klein, M., Pirrelli, V.,
Poesio, M., Quazza, S., Schiffrin, A. and Soria, C. 2000. *MATE Dialogue
Annotation Guidelines*. MATE Deliverable D2.1.