

A Prosodic Corpus of Non-Native Speech

Jan-Torsten Milde & Ulrike Gut

Department of Linguistics and Literary Studies
University of Bielefeld, Germany

milde@coli.uni-bielefeld.de, gut@spectrum.uni-bielefeld.de

Abstract

The paper describes the design and implementation of an XML-based corpus environment for prosodically annotated data. The TASX-environment (TASX: Time Aligned Signal data eXchange format) constitutes the technical basis for a corpus designed to explore the acquisition of prosody by second language learners. It supports all aspects of the corpus setup procedure: XML-based annotation of the speech data, all transformation of non XML-annotations, and the web-based analysis and dissemination of the data.

1. Introduction

In this paper we describe ongoing research in the design and implementation of an XML-based corpus environment for prosodically annotated data. The development of the corpus environment is part of the LeaP project, which explores the acquisition of prosody by second language learners of both German and English. In a period of two years a large set of recordings of second language learners' speech will be made and phonologically annotated. From this data an XML-annotated spoken language corpus will be set up. The model is based on a client/server approach. For performance reasons the XML-annotated data can be stored in a relational database. The XSL-T-based transformation of the data is a server sided process. The TASX-environment presented here supports the complete corpus setup procedure: XML-based annotation of raw speech data, the transformation of non XML-data and the analysis and dissemination of the corpus.

The paper is organized in five sections. First, a short overview of the LeaP project will be given, which explains the specific requirements for the TASX-environment. In the next section the underlying XML-based TASX format will be explained and the components of the TASX-environment will be described in more detail. In section 4, we will then explain how the LeaP corpus has been set up with the TASX-environment and how the data can be linguistically analysed. Finally, a short conclusion will be given.

2. The LeaP project

The LeaP (Learning Prosody) project¹ explores the acquisition of prosody by second language learners of both German and English. It focuses on three areas of prosody: stress assignment on both the word and the phrase level, sentence intonation and speech rhythm. So far, the acquisition of second language prosody has not attracted a large amount of research but it has nevertheless often been proclaimed to be nearly impossible [2]. However, this assumption of non-attainment so far has

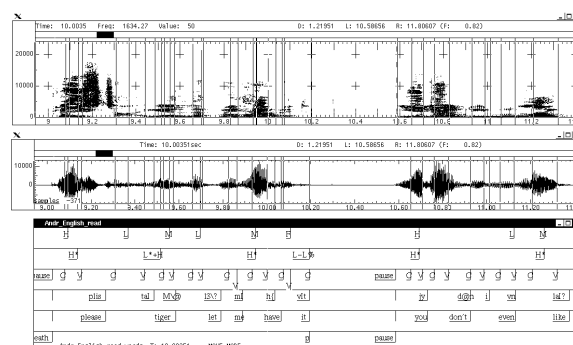


Figure 1: Prosodic annotation in the LeaP project is carried out on six tiers: a phrase, an orthographic tier, a syllable tier, a rhythm tier, a tone tier and a pitch tier.

only been supported by a few single case studies [1].

Research in the LeaP project will be based on a large corpus of second language learners' speech. The focus lies on two main research questions: first, a detailed description of the second learners' prosody within the latest theoretical frameworks and a comparison to native speakers' prosody will be carried out. It is assumed that second language prosody constitutes a good testing ground for theoretical concepts in prosody and might provide evidence for their further development. The second line of research aims to provide an assessment of the extralinguistic factors such as personal variables (e.g. native language, age at the beginning of language learning, motivation, musicality) and the type of teaching method that might enhance the outcome and speed of the acquisition process. The LeaP experimental design comprises three treatment groups, who will undergo intensive prosodic training in English and German of up to one and a half years duration.

For both research questions a multitude of data of various types will be collected: the corpus of spoken language will consist of at least 400 recordings of between 2 and 10 minutes length. It comprises three different speech styles: read speech, prepared speech (a retelling of a story) and free speech. In the first six months of the project, 107 recordings have already been made. In addition to this speech material, meta data have been collected for every speaker. This consists of personal data such as the learners' age, sex, native language and the onset of learning of the second language, as well as ratings of motivation and interest.

The prosodic annotation of the speech material is carried out using ESPS/waves+ with six different tiers (see figure 1). On the first tier, the phrase tier, phrases are annotated as well as anything occurring between them (e.g. pauses, laughter, noise).

¹<http://www.spectrum.uni-bielefeld.de/LeaP/>

On the second tier, the word tier, each word is transcribed orthographically. On the third tier, the syllable tier, each syllable is transcribed in SAMPA. On the fourth tier, the rhythm tier, the vocalic and consonantal parts of the speech are annotated. On the fifth tier, a transcription of intonation in a ToBI [10] style is carried out. The sixth tier, the pitch tier, contains an annotation of highs and lows in the pitch contour. This means that for each recording there are approximately 3000 time stamps. In the first six months of the project, 67 recordings have been annotated in this fashion.

3. The TASX format

A central aspect of our research is to explore up to which point current standard XML technology (XML, XSL-T, XSL-FO, XPATH, SVG, XQUERY) can be used to model linguistic databases, to transform, query and distribute the content of such databases and to perform adequate linguistic analysis. As a result, all linguistic data in our system is stored in an XML-based format called TASX: the Time Aligned Signal data eXchange format.

A TASX-annotated corpus consists of a set of sessions, each one holding an arbitrary number of descriptive tiers, called layers. Each layer consists of a set of separated events. Each event stores some textual information (e.g. a syllable) and is linked to the primary audio data by two time stamps. Relations between events on different tiers can be encoded by defining links using the ID/IDREFS mechanism of XML. Finally, arbitrary meta-data can be assigned to the complete corpus, each session, each layer and each event.

3.1. The TASX-annotator and the corpus engine

The complete TASX-environment consists of:

- tools for the annotation of empirical language data (video and audio material),
- an input mask for processing meta data
- programs for the transformation of various formats of linguistic standard software (Transcriber, Praat, ESPS/waves+, SyncWriter, Exmaralda etc.)
- a set of programs for linguistic analysis of the TASX-annotated data, and
- a corpus system for the distribution of language data via the internet, including interactive corpus query and multimodal data display in a standard web browser.

In the following sections these modules will be described in more detail (see also [6]).

3.2. The TASX-annotator

The TASX-annotator is a central component of the TASX-environment. The tool allows the annotation and transcription of video (multi-channel) and audio data (see figure 2).

The program is extremely user friendly and can be used without a high level of computer skills. It is possible to completely control the tool by either mouse *or* by keyboard shortcuts. Video and audio playback can be controlled by a foot switch. Different data views are programmed (time-aligned partiture, word-aligned partiture, sequential text view) to make annotation as effective as possible.

The time aligned view is organized as a two dimensional grid of infinite size. A layer is presented as a horizontal tier of events.

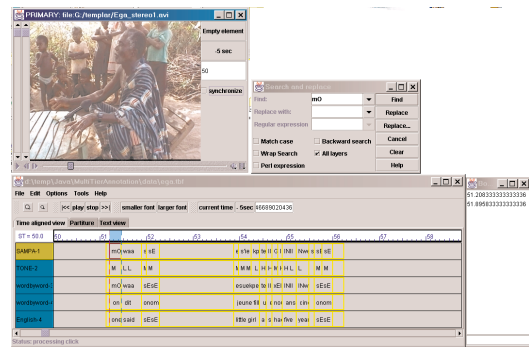


Figure 2: A screenshot of the TASX-annotator. In the bottom half, the main panel is visible, where the time aligned tier view has been selected. Next to the main panel the bookmark list is displayed and above it the find tool. In the upper left corner the video display can be seen, showing some data from the Ega corpus.

The order of the layers is arbitrary and can be changed instantly. The user is able to define time intervals by dragging the mouse. Each time interval represents an event. The event is displayed as a graphical box which can be selected and moved with the mouse.

In the text view the data can be manipulated in a standard text editor panel. The content of the editor represents the layer and each line represents an event. A list selection box allows to switch between different layers. It is possible to transfer text from standard text editors, e.g. Microsoft Word, by cut and paste operations. To additionally speed up the transcription process, a word completion function has been implemented for the text view. Entering the initial letter of a word and consecutively pressing CTRL+L will bring up all words starting with this letter. Once the text is transferred into the TASX-annotator, the events still have to be aligned with the primary audio and video data. Switching back to the time aligned view and moving the events with mouse makes this task quite simple.

In the partiture view the data cannot be edited. In practice this means that the data is transformed into an HTML table and then displayed to the user. A number of different HTML formatted views have been designed. The views can also be saved to external files and loaded into standard web browsers.

One potential strength of the TASX-annotator is its manner of handling the export/import of XML based information. A standard way of solving this problem would be the implementation of a set of format specific XML parsers which construct the internal representation (e.g. JDom) of the XML file. While powerful integrated development systems such as *Sun's Forte for Java* make the design of such XML handlers simpler, it still remains a complex task to implement such a parser. In the TASX-annotator we follow a different approach. The system integrates an XSL-T processor (saxon), making it easy to perform on the fly data transformations. The import of an XML-file is split into two steps: first an XSL-T stylesheet transforms the XML file into TASX, second another XSL-T stylesheet will transform the TASX file into a simple text oriented format. This format can be loaded efficiently.

Table 1: List of currently implemented transcoding tools. The table shows the programming languages used to implement the transcoders.

TASX	import	export
Annotation graphs	XSL-T	XSL-T
Exmaralda	XSL-T/Java	Java/XSL-T
HTML-table	–	XSL-T
HTML-partiture	–	XSL-T
RTF	–	XSL-T/Java
Anvil	XSL-T	–
Praat-label	Perl/XSL-T	XSL-T
ESPS-label	Perl	XSL-T
ESPS-freq	Perl/XSL-T/Java	XSL-T
SyncWriter	Perl	–

3.3. Transcoding tools

The development of tools for the TASX-environment is based on the concept that a re-implementation of functionalities already available in other speech processing software is not necessary. Established speech software such as Praat or ESPS/waves+ do not need to be duplicated. The TASX-environment therefore focuses only on the development of transcoding filters from and into various formats. These include: Praat/freq, Praat/label, ESPS/waves+, ESPS/F0-analysis, Transcriber, annotation graphs stored in XML, SyncWriter and basic text formats (see table 1). In addition, filters for data import and export of the Exmaralda system [9] are available. Most of these components are implemented in Java, transformations are defined in XSL-T and a smaller number of additional tools is written in Perl (mainly to transform non-XML data).

3.4. The corpus system

The main function of the corpus system constitutes the internet-based dissemination of the corpus data. With the currently implemented interface it is also possible to inspect and query the speech corpus, to listen to the audio material and to display the graphic representation of the waveforms in a standard web browser. We make use of the built-in features of the web browser here. Furthermore, the PAX-tools [3] for displaying the intonation contour, the intensity and the spectrogram of the selected regions in the audio file can be integrated.

When playing back the sound file, both the audio parts and the waveform images are generated automatically by a small Java servlet program. The servlet parses the XML-annotated corpus, extracts the time stamps of the relevant events and then cuts out the corresponding parts of the original sound file.

The corpus system is split into two larger subcomponents: the *information pool* and the *corpus engine* (see figure 3). The information pool stores the primary data (raw audio data) as well as the XML-annotated transcriptions of the audio files. The corpus engine consists of five subsystems:

1. Web-client: the interactive user interface is completely defined to run in a standard web browser. We are using HTML-query forms which activate services on the server side to generate XSL-T-filters processing the data. Waveforms are displayed using SVG. This will allow the user to select parts of the sound signal and to perform more complex phonetic analyses.

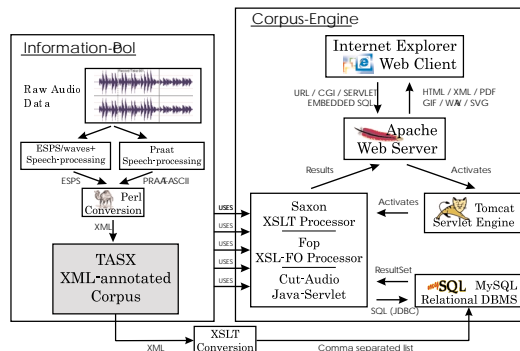


Figure 3: The system architecture of the corpus system. The corpus system is split into two subsystems: the information pool (left) storing the TASX-annotated data and the corpus engine (right) distributing the data over the internet.

2. Web-server: the web server distributes the corpus information in several standard formats (XML, HTML, PDF, SVG, WAV).
3. Servlet-engine: the servlet engine activates the suitable services on the server side (transformation of XML-annotated data, on-the-fly phonetic analysis, generation of graphics).
4. Servlets: a set of TASX/XML-aware servlets are used to transform the data in numerous ways: generating HTML to be displayed in the browser, generating PDF to be printed out, generating wavefiles and images of the waveforms. XSL-T and XSL-FO are used to perform the transformations. The servlets have access to the information pool and the relational database.
5. Relational database: in order to improve the system performance, the XML-annotated corpus data is stored in a relational database. The database basically replaces a standard file system. An XSL-T-program translates the XML-annotated corpus data into a suitable format for the DBMS.

The implementation of the corpus system is based on open source software. The TASX-Annotator is a pure Java application; all other tools are smaller XSL-T and perl scripts. As a result, the complete TASX-environment runs on Windows and Unix platforms. The software will be distributed under GPL and can be downloaded from our website².

4. Setting up the TASX-annotated corpus

Once the prosodic annotation as described in section 2 is completed the TASX-annotated corpus can be set up.

4.1. LePa data conversion

First, the ESPS/waves+ files are being converted into the XML-based TASX format. This is done automatically by a small perl program called `esps2tasx`. The converter is able to take a whole set of ESPS/waves+ files and transform them into one large TASX-annotated corpus. Next, the TASX-annotated data is stored in a relational database. This is done to improve the performance of the corpus system. Each session is stored as a binary large object.

²<http://coli.lili.uni-bielefeld.de/~milde/tasx/>

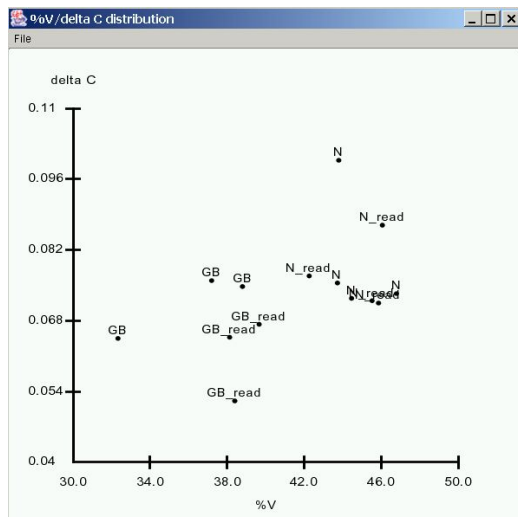


Figure 4: The graph shows an example for the distribution of speech rhythm of British English native speakers (GB) and Nigerian speakers of English (N). A clear separation between the two groups can be observed. The graph is generated in SVG format by directly performing calculations on the TASX-annotated corpus.

After the data has been transformed into the TASX-annotated form, it becomes possible to use the complete set of tools of the TASX-environment.

4.2. Analysis of the prosodic data

The phonetic analysis of the LeaP data is carried out in a semi-automatic style, supported by various TASX analysis tools. For the calculation of the speech rhythm according to Ramus et al. [8] for example, the information of the fourth tier as described in section 2 is taken. The length of all vocalic (V) and all consonantal (C) parts of the utterances in the recording are calculated and their standard deviation (ΔV and ΔC) is computed, as well as the percentage of vocalic intervals (%V) across the entire recording. These measurements have proved useful for the description of the difference in speech rhythm between languages and varieties of languages ([5]). The results for all speakers are illustrated in an automatically generated graph (see figure 4).

Similarly, the analysis of the speakers' pitch range is carried out semi-automatically. From the time stamps of the fifth tier the pitch height is taken from the corresponding ESPS/waves+get_f0 file and the following measurements for the pitch range and pitch span analysis according to Patterson [7] are calculated: mean initial highs, mean subsequent highs, mean lows, mean final lows.

A third area of prosodic analysis of the speech data is tonal alignment in stressed syllables. English and German differ in that respect ([4]) and it might be a useful feature for the description of non-native prosody. For the analysis, the time stamps on the tone tier, which provides information about the occurrence of stressed syllables, the time stamps on the pitch tier, which gives pitch maxima, and the time stamps on the rhythm tier, which indicates the vowel boundaries, are combined and the presence of pitch height in relation to the vowel boundary is calculated.

Other features important for the analysis of language learn-

ers' prosody such as speech rate, fluency, and intonation patterns are also supported by TASX analysis tools. All tools will be freely available to the scientific community as open source software.

5. Conclusions

Despite the early stage of the research the TASX-based approach has already proved to be highly efficient and reliable. The time consuming task of phonetically analysing speech data is partially substituted by automatic analysis. In the transformation process from non-XML to XML-annotated data some errors in the human annotations can be detected. Furthermore, due to the highly structured format of the TASX-converted data more complex research questions can be investigated in a systematic way.

The very good availability of XML aware software and tools enabled us to develop a powerful linguistic environment in a very short time. Even more important, the TASX-annotated data can be transformed into large number of different formats. The will hopefully lead to the creation of linguistic resources which can be used over a long period of time by different researchers with various goals.

6. References

- [1] Archibald, J., 1998. *Second language phonology*. Amsterdam: Benjamins.
- [2] Boyle, J., 1987. Perspectives on stress and intonation in language learning. *System*, 15(2):189–195.
- [3] Gibbon, D.; Trippel, T., 2001. PAX - an annotation based concordancing toolkit. In Peter Buneman, Steven Bird and Mark Liberman, editors, *IRCS Workshop on Linguistic Databases, University of Pennsylvania, Philadelphia, USA*.
- [4] Grabe, E., 1998. *Comparative intonational phonology: English and German*. MPI Series in Psycholinguistics 7, Wageningen, Ponsen en Looien.
- [5] Gut, U., 2001. The prosody of Nigerian English. In Dafydd Gibbon and Ulrike Gut, editors, *Proceedings of the TAPS 2001 workshop, Bielefeld*.
- [6] Milde, J-T.; Gut, U., 2001. The TASX-environment: an XML-based corpus database for time aligned language data. In Peter Buneman Steven Bird and Mark Liberman, (editors), *IRCS Workshop on Linguistic Databases, University of Pennsylvania, Philadelphia, USA*.
- [7] Patterson, D., 2000. *A Linguistic Approach to Pitch Range modelling*. Ph.D. thesis, University of Edinburgh.
- [8] Ramus, F.; Nespors, M.; Mehler, J., 1999. Correlates of linguistic rhythm in the speech signal. *Cognition*, 73:265–292.
- [9] Schmidt, T., 2001. Gesprächstranskription auf dem Computer - das System EXMARaLDA. *Gesprächsforschung*, <http://www.gespraechsforschung-ozs.de>, 2.
- [10] Silverman, K.; Beckman, M.; Pitrelli, J.; Ostendorf, F.; Wightman, C.; Pierrehumbert, J.; Hirschberg, J., 1992. Tobi: a standard for labeling english prosody. In *Second International Conference on Spoken Language Processing 2, Banff, Canada*, 867–870.