

# Encoding of Swedish dialect transcripts

The aim of this project is to develop methods and standards for data entry and encoding of dialect transcripts. Though independent, the project is associated with other Swedish dialect documentation projects.

## Background

The archives of Swedish dialect research institutions contain large collections of recordings and transcripts. Most of these transcripts are written with an unique phonetic alphabet, the Swedish dialect alphabet, specially designed in 1878 by Professor J. A. Lundell of Uppsala. Only a limited number of scholars are able to read these transcripts, which contain a large amount of information of both linguistic and ethnological interest. To make these dialect collections accessible to the general public, various kinds of software tools are needed for the encoding, conversion, search, and display of these texts.

## The project

The project is divided into three subprojects: Encoding tools, Character codes, and Text conversion.

### *1. Encoding tools*

This subproject involves the creation of detailed standards for dialect recording transcript formats and the development of software and other tools for entering dialect texts.

- Making an XML schema for dialect transcripts, containing tags for various types of metadata.
- Collecting information about existing software tools - fonts, keyboard modifiers, OCR software etc.
- Development of new software where existing software isn't sufficient.
- Creation of manuals and other types of training material for the persons carrying out the actual work of entering, importing, and cataloguing dialect transcripts.

### *2. Character codes*

This subproject involves the documentation of the phonetic values of all the characters of the Swedish dialect alphabet and including these characters in the ISO 10646/Unicode character code standard.

- Investigating actual character usage in different dialects and various types of dialect texts. This phase should include cataloguing the usage of similar dialect alphabets in Denmark and Norway.
- Making translation tables between the Scandinavian dialect alphabets and the International phonetic alphabet (the IPA alphabet).
- Making a proposal for including the characters of the dialect alphabets in ISO 10646/Unicode, the international character code standard.
- Making Unicode-encoded fonts for free distribution, containing all the Scandinavian dialect characters.

### 3. Text conversion

This subproject contains the development of sets of rules for conversion between various orthographic systems, e.g. conversion from the Swedish dialect alphabet into IPA or 'phonetic spelling', i.e. ordinary alphabetic characters indicating the pronunciation. There will probably be a large number of rule sets, due to phonetic differences between the dialects and the demand for different levels of accuracy.

*ə áá, ja míns min fár, də va min fá sòm hàde ə gádən dá. han sadə æt̪er ən fjòtənə máš, ə ə də  
 jík ju brá ... di fík ju. vərə ældrə sadə di sàdə alti t̪idiarə ən di jø\_ɲú fə\_ɲín. ... si di hàdə san  
 ɲèskap, ša di va tvùgnə te ə ə bòrja nàgərlúnda, sa di kunnə hinə mæ. nu há\_ɲi san ɲèskap ša nú  
 ə ... mæ traktórər šá sa də də gá\_ɲə mykə ɲàskarə ...*

*ə ɔ:'ə, jø mi'n:s min fɔ:'r, de vɔ min fɔ:' šæm ha'd:e ə go:'ɲən do:' . han soðə ɛ't̪:ɛ,r ən  
 fjɔ:tənə mɑ'š:, ə ə de jík ju brɔ:' ... di fík ju. vo:ra ɛ'l:drə soðə di so'd:ə alti ti:'diarə  
 ən di jø: ɲə' fce ɲi:'n. ... si di ha'd:ə son ɲe'škəp, šo di vɔ tvø'ɲ:nə te o o bæ'r:ja  
 no:'gərlø,n:da, so di kunnə hin:a mæ:'. nɛ hɔ:' ɲi son ɲe'škəp šo nɛ:' ə ... mæ traktɔ:'rər  
 šo:' so de:' de de go:' ɲə mykə ɲɑ:'s:karə ...*

*e ae, ja minns min far, də va min far sòm hadde e gården då. han såde ärter en fjortene mars, ə e  
 də jík ju bra ... di fík ju. våra äldre såde di sådde allti tidiare än di jör nu för tin. ... si di hadde  
 sån rerskap, sjä di va tvungne te å å börja någerlunnda, så di kune hinna mä. nu har di sån  
 rerskap sjä nu e ... mä traktorer sjä så då då gå de myke rasskare ...*

Figure 1. At the top, a transcript written with the Swedish dialect alphabet; in the middle, the same text automatically converted to IPA characters; at the bottom, the same text automatically converted to Swedish 'phonetic spelling'.