UraLUID: Supporting data-driven (prosodic) research

Eszter Simon, Nikolett Mus, Ágnes Kalivoda, Eszter Ruttkay-Miklián

Research Institute for Linguistics, Hungarian Academy of Sciences {simon.eszter,mus.mikolett,kalivoda.agnes}@nytud.mta.hu, eruttkaym@gmail.com

> 2nd workshop on Uralic prosody September 28-29, 2017, Budapest

Outline

- Introduction
- The content of the database
- The structure of the database
- Future plans

Simon et al. (RIL HAS)

Background

- Languages under influence. Uralic syntax changing in an asymmetrical contact situation
- Research Institute for Linguistics, Hungarian Academy of Sciences
- January 2016 July 2017
- supported by the National Research Development and Innovation Office
- website: http://www.nytud.hu/depts/tlp/uralic/dbases.html

Languages under influence

The pilot-project had a twofold objective:

- Theoretical work
 - to describe and analyze potential syntactic changes due to Russian contact
 - examined languages: Udmurt, Northern Khanty (Synya dialect),
 Eastern Khanty (Surgut dialect group), and Tundra Nenets
- Practical work
 - to collect and structure primary data
 - to build a linguistically annotated database of written and spoken texts, as well as of "old" and "new" data

- Introduction
- The content of the database

- The structure of the database
- 4 Future plans

Source types

	Udmurt	Synya Khanty	Tundra Nenets	Surgut Khanty
"old"	folklore	folklore	folklore	folklore
"new"	blog post	interview	newspaper spontaneous speech folklore	interview spontaneous speech

Data types

	Udmurt	Synya Khanty	Tundra Nenets	Surgut Khanty
"old"	written	written	written	written spoken
"new"	written	spoken	written spoken	written spoken

Spoken data

	Speakers	Place of recording	Length	Format
Synya Kh.	1	Ovgort	01:08:22	.wma
Tundra Nenets	1	Moscow	01:57:40	.wav
Surgut Kh. (new)	2	Kogalym	00:28:58	.wav
Surgut Kh. (old)			00:36:17	.wav

The Synya Khanty spoken data

- the original data were collected within the frame of the project entitled In the Khanty way (http://hantisirn.nytud.hu/)
- source: interviews
- genre: expository, narrative texts

The Tundra Nenets and Surgut Khanty spoken data

- the spoken data were collected during our fieldworks
- source: spontaneous speech data (elicited by graphical materials), "read" speech
- genre: narrative, procedural, expository texts

Introduction

- 2 The content of the database
- The structure of the database
- 4 Future plans

Data available on our website

- Metadata
 - title of the text
 - text ID
 - pages
 - file name
 - genre
 - token number
 - (sub)dialect
 - age of the speaker
 - gender of the speaker

http://www.nytud.hu/depts/tlp/uralic/dbases.html

- 2 Text
 - original transcription
 - Cyrillic
 - FU transcription(s)
 - IPA
 - ⇒ the original text split into sentences

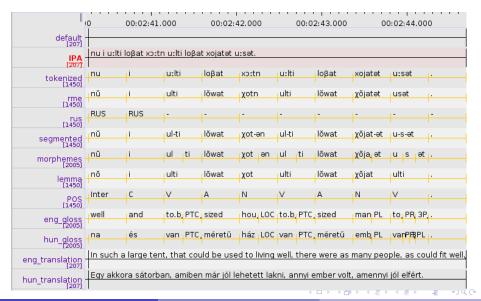
- Morphologically analyzed text: .tsv files consisting of 15 columns
 - 1-9: the token and its transcriptions
 - 10: (segmented) token
 - 11: lemma
 - 12: Hungarian gloss
 - 13: English gloss
 - 14: POS tag (and semantic label if there is any)
 - 15: RUS label (if the word has Russian origin)

- Translation
 - English
 - Russian
 - German
 - Hungarian
 - \Rightarrow every translation is sentence-level aligned with the original text

ELAN files

- .eaf file (containing data on sentence-, token- and morpheme-level)
- the .eaf files are created with a script written in Python3
- the script uses the Pympi module developed for creating and processing ELAN and Praat annotation files (http://github.com/dopefishh/pympi)
- the original sentences are aligned to the time slots of the audio file
- the other pieces of information are connected to the sentences via symbolic references

Data in ELAN



How to use the UraLUID database?

- use the .eaf files with the corresponding audio files:
 - download the latest version of ELAN:
 https://tla.mpi.nl/tools/tla-tools/elan/
 - download the Charis SIL font package: https://software.sil.org/charis/
 - download the .eaf and audio files: http://www.nytud.hu/depts/tlp/uralic/dbases.html
 - \bullet open ELAN \to Open... \to choose the needed .eaf file and the corresponding audio file
- use the .tsv files
 - download the .tsv files
 - use Unix commands or your dear old statistical tools

Introduction

- 2 The content of the database
- The structure of the database
- Future plans

Future plans

Tundra Nenets spoken data:

- Transcription
 - Cyrillic
 - FUT
 - IPA
- Morphological analysis ⇒ improving the Giellatekno's morphological analyzer
- POS-tags
- Translation
 - Russian
 - English

Acknowledgements

The authors wish to thank Elena Skribnik, Zsófia Schön, Agrafena Pesikova, Anisya Volkova, Katalin Gugán, Márta Csepregi, Olesya Khanina, Aleksey Kozlov, Khadry Okotetto, Orsolya Tánczos, Erika Asztalos, Yulia Speshilova, and Szilvia Németh for their contribution to the building of the database.

Thank you for your attention!

http://www.nytud.hu/depts/tlp/uralic/dbases.html



Simon et al. (RIL HAS) UraLUID 22 / 22