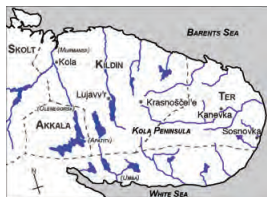# NEWS FROM SAAMI DOCUMENTARY LINGUISTICS

## Michael Rießler

Ä′vv Sää′mi mu′zei Njauddâm

## Budapest, 17th May 2012
## Piliscsaba, 18th May 2012

## INTRO

- Language documentation has a long tradition. Comprehensive language documentations of Saami languages were produced already in the beginning of the last century.
- Documentary linguistics evolved from traditional methodology in language documentation, but has become a linguistic sub-discipline of its own. The primary aim of the field is providing more and better data on the world's linguistic diversity for future research on and for endangered languages.
- The presentation focuses on the interfaces of documentary linguistics and language technology specifically from the perspective of applied research for endangered Saami languages.
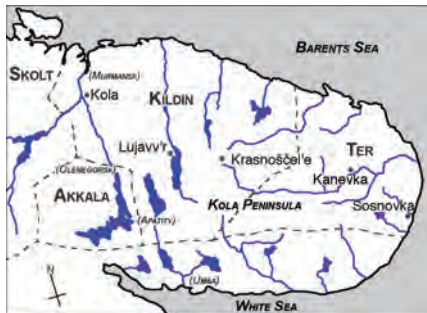
## BAKGROUND

- Skolt Saami museum in Neiden (currently)
- Universität Freiburg (July 2012)
- Foci: Kola Saami languages
  - Documentation
  - Description
  - Revitalization
- Other research:
  - Noun phrase syntax (history and typology) in northern European languages



*Bodø, March 2010*
from left: *M. Rießler, G. Lukin, D. Nathan*
(training in Saami documentary linguistics)

# KOLA SAAMI



- East-Saamic < Uralic, NW-Russia, Finland, (Norway)
  - Kildin Saami (< 800 speakers)
  - Skolt Saami (< 500 speakers)
  - Ter Saami (< 20 speakers)
  - Akkala Saami (< 5 speakers)

# KOLA SAAMI DOCUMENTATION PROJECT (KSDP)



- DoBeS project (2005–2011)
  Documentation
- DFG project (2012–2015)
  Description
- Collaboration with Giellatekno (U Tromsø)
  Language technology

## CONTENT

- Documentary linguistics
  - Definition
  - Short history of the field
- Kola Saami Documentation Project (KSDP)
  - Archive structure
  - Workflow practices
- Documentary linguistics as an applied discipline
  - Documentary linguistics vs. Language documentation
- Practical questions (language documentation and theoretical linguistics)
  - Metadata, equipment, data formats and archiving
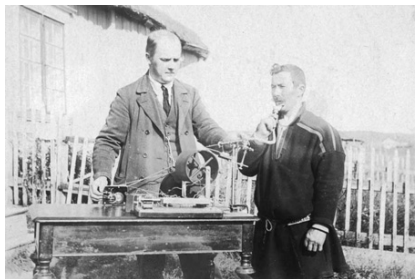
# DOCUMENTARY LINGUISTICS



*Edison Phonograph (1877)*

## DOCUMENTARY LINGUISTICS

- New and evolving field in linguistics
- Primarily concerned with language documentation (i.e. "a comprehensive, multi-facetted and multi-purpose record of linguistic practices characteristic of the investigated speech community"), which must be
  - comprehensive
  - multi-faceted
  - multi-purpose
- Two important purposes of language documentations are
  - being a data pool for theoretical research (in typology, anthropology, etc.)
  - being a data a pool for practical research (on revitalization, pedagogy, language technology, etc.)
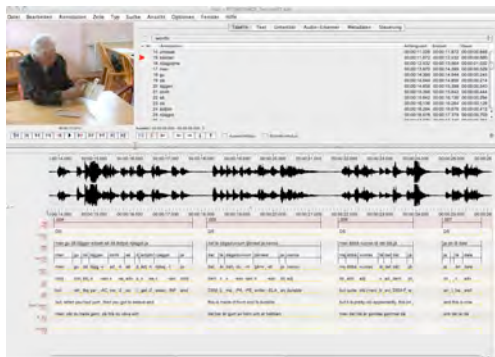- Documentary linguistics as primarily an applied discipline

# DOCUMENTARY LINGUISTICS (E.G. PITE SAAMI)



*Eliel Lagercrantz (1894–1973) documenting Pite Saami*

- Classical documentary trilogy ("Boasian Trilogy"), e.g.
  - *Sprachlehre des Südlappischen* 1923.
  - *Wörterbuch des Südlappischen* 1926.
  - *West- und südlappische Texte* 1957 / *Texte aus den see-, nord-, west-, und südlappischen Dialekten* 1963.

# DOCUMENTARY LINGUISTICS (E.G. PITE SAAMI)



*Pite Saami Documentation Project (PSDP), cf. www2.hu-berlin.de/psdp/*

- "Language documentation is a lasting, multipurpose record of a language" (Himmelmann 2006:1).

# SAAMI DOCUMENTARY LINGUISTICS

- Endangered languages – here Saami – must be documented before it is too late. However,
  - Only gathering language data does not result in useful documentation.
  - Storing language data is not the same as creating an archive.
  - Digitization alone is not sufficient preservation.

## SAAMI DOCUMENTARY LINGUISTICS

- With this in mind, it becomes obvious that the biggest problems with traditional methodology in language documentation are:
    - the lack of information about the existence and relevance of recordings (missing metadata);
    - their unaccessability (missing annotations and/or physical unavailability);
    - as well as the question of long-term storage and preservations (use of open, non-proprietary formats and ensuring longevity of data carriers).



nflrc.hawaii.edu/ldc/

# KOLA SAAMI DOCUMENTATION PROJECT (KSDP)

- Our main goals were (are):
  - documenting (record, annotate, archive) as much text recordings from different genres as possible, especially (formerly undocumented) conversations, spontaneous speech and procedurals
  - creating a corpus and other materials useful for future theoretical and practical research
  - working closely together with the Kola Saami communities



*Elicitation session*
*Lujavv'r, August 2009*
from left: *A. Antonova, M. Rießler*

## ARCHIVE STRUCTURE

- **Kola Saami Documentation Project** (introduction and general info)
- **Commentary** (access rights, bibliography, annotation conventions, project biography)
- **Misc. data** (fieldnotes, consultants' handwritten notes and wordlists, local newspaper articles, etc.)
- **Pictures**
- **Recordings**
    - —External
    - —KSDP
- **Studies** (papers and presentations on phonology, language sociology, etc.)
- **Teaching materials** (school grammar, text collection, posters, etc.)

# ARCHIVE STRUCTURE



- General archive structure
- www.mpi.nl/dobes/

# ARCHIVE STRUCTURE



- Example of a session with different files
- www.mpi.nl/dobes/

## WORKFLOW PRACTICES

- Annotations include minimally catalogue metadata, preliminary orthographic transcriptions and a translation into Russian.

- First transcriptions and translations are done by native speaker assistants (teachers, other interested language workers).



*Annotation, Lujavv'r, February 2006*
from top: *A. Mozolevskaya, E. Scheller, N. Zolotuchina*

## WORKFLOW PRACTICES

- Recording
    1. Recording on video and/or audio
    2. Written notes: catalogue metadata, preliminary annotations, elicitation stimuli, etc.

- Processing of data
    1. Digitizing video (if applicable)
    2. Preparing data for further processing (collecting session data, session cutting, labeling, storing on local server)
    3. Original transcription (by native speaker assistants) in a preliminary orthography and Russian translation
    4. Edited transcription (of selected recordings) in normative orthography, English translation, morphology, phonology, etc.

- Archiving

# DOCUMENTARY LINGUISTICS AND APPLIED RESEARCH

Is *documentary linguistics*

- a new evolving subfield of applied linguistics?
  *documentary linguistics ≥ language documentation*

or is it

- just a (newly improved) empirical method for linguistic data collection?
  *documentary linguistics = language documentation*

## ANNOTATION DEPTH IN KSDP

### STILL THE SAME QUESTION…

How deeply should we annotate our Kola Saami corpus?

… like Lagercrantz

pu̯aᴰz̧a̧ ji̯all vāŗeṣt

## ANNOTATION DEPTH IN KSDP

### STILL THE SAME QUESTION…

How deeply should we annotate our Kola Saami corpus?

… like contemporary typologists

| puaz | jāll | vār'-es't |
|------|------|-----------|
| n | v | n |
| deer | live\3SG:PRS | fjell(WK)-LOC.SG |

'the reindeer lives on the fjell'

# ANNOTATION DEPTH IN KSDP

## OUR ANSWER...

пуаз я̄лл ва̄ресьт

'олень живет в тундре'

'the reindeer lives on the fjell'

- Orthography and translation is adequate for the documentary linguist
- The theoretical linguist can do further analyses, annotations, descriptions, etc.
  - because sufficient phonological/grammatical and lexical description exists as the result of our and earlier research

# DOCUMENTARY LINGUISTICS AND APPLIED RESEARCH

For the *documentation of Kola Saami* an orthographic representation is even preferable, because

- faster work provides more text annotations for
  —revitalization *and* research
- merging our spoken language data with written text corpora results in a relatively large corpus, which helps producing tools for
  — revitalization *and* research, e.g.
  - Digital infrastructure: spell-checkers, machine translation, teaching aids, etc.
  - Corpus linguistic tools: morphological and syntactic analysers, etc.
- North Saami (15,000 speakers, corpus of 500,000 words)
- *Boazu eallá váris*

## DOCUMENTARY LINGUISTICS

- We see *documentary linguistics* as a discipline concerned with language documentation (i.e. "a comprehensive, multi-facetted and multi-purpose record of linguistic practices characteristic of the investigated speech community")
- Two important purposes of language documentations are
  - being a data pool for theoretical research
    (in typology, anthropology, etc.)
  - being a data a pool for practical research
    (on revitalization, pedagogy, language technology, etc.)
- Documentary linguistics as primarily an applied discipline

## RESEARCH ETHICS

- When working with speakers of endangered languages, you should always consider making available your data (and research results)
  - Archiving!

- If you believe in linguistics as a *science*, you should also consider making available your data (and research results) – and making your analyses provable (and falsifyable)
  - Archiving!

## LEGACY DATA

- The empirical linguist (of whatever framework) often also records data useful for the documentation of endangered languages.
- Important points to consider when creating legacy data:
  - Accessability (physically and intellectually)
  - Long-term storage and preservation (data formats and carriers)

# METADATA

- Minimal cataloging metadata
  - Session label
  - Actors (consultant, collector, annotator, etc.)
  - Date and place
  - Research background
  - …
- Annotations
  - Elicitation materials (questionnaire)
  - Transcription, translation
  - Research notes
  - …

## DATA FORMATS

- Audio
  - Preferably .wav
- Text (Metadata, Annotations)
  - Preferably open text formats (e.g. simple .txt)
  - Alternatively .pdf versions of your word document

## EQUIPMENT

- Audio is normally most important
- Example of mid-range semi-professional equipment
  - Microphone Sony ECM-MS 957 (approx. € 250)
  - Recorder Edirol R-09 (approx. € 350)
- Software
  - ELAN (freeware)
  - Texteditor (freeware)