# Applying spatial data in linguistics

Kristel Uiboaed, Siim Antso, Liina Lindström,

Maarja-Liisa Pilvik, Mirjam Ruutma

University of Tartu

# Overview

- Project
- Resources and data
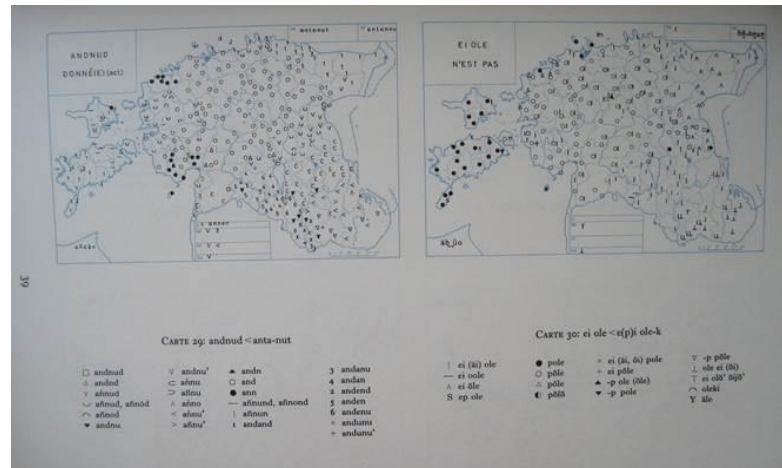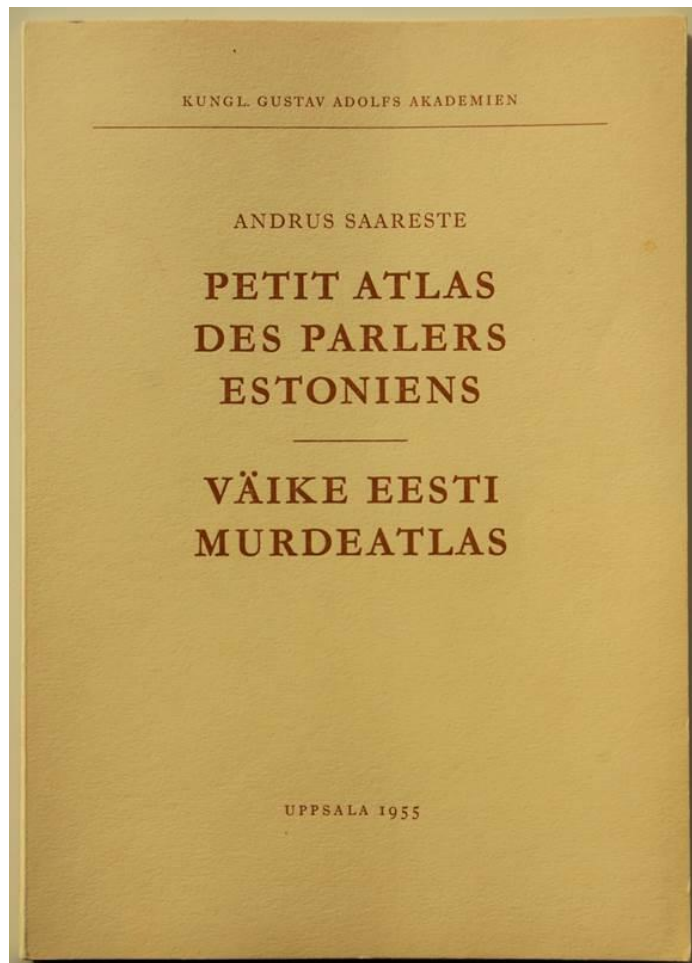- Workflow
- Conclusions

# Project

- Applying spatial data in linguistic research (01.04.2015–31.03.2017)
- Objectives:
  - digitalizing (scanning and georeferencing) data in Andrus Saareste's dialect atlases and his unpublished manuscript maps
  - integrating map applications with data from the Corpus of Estonian Dialects (CED)
  - representing linguistic distances in CED by building interactive applications for multidimensional analyses
  - making the above-mentioned resources freely available
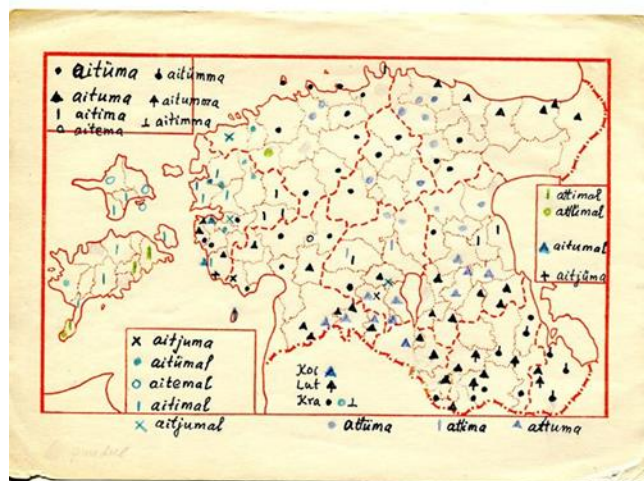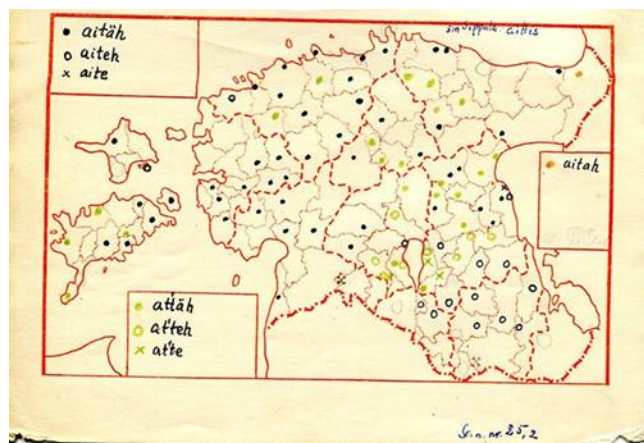
# Andrus Saareste's maps

- Mainly based on the data A. Saareste and his colleagues collected in Estonia during 1915–1944 (and that he managed to take with him when he emigrated to Sweden), to a lesser extent also on the material he collected from over 200 Estonian refugees while staying in Uppsala in 1945–1954.

# Andrus Saareste's maps

- Atlases
  - Petit atlas des parlers estoniens. Väike eesti murdeatlas (1955)
  - Eesti murdeatlas. I vihik (Atlas des parlers estoniens) (1938)
  - Eesti murdeatlas. II vihik (Atlas des parlers estoniens) (1941)
  - Altogether 215 maps
  - Based on Estonia's map during 1920–1940ies
- Manuscripts
  - Maps in Uppsala University archives
  - 2500 maps scanned and to be digitalized

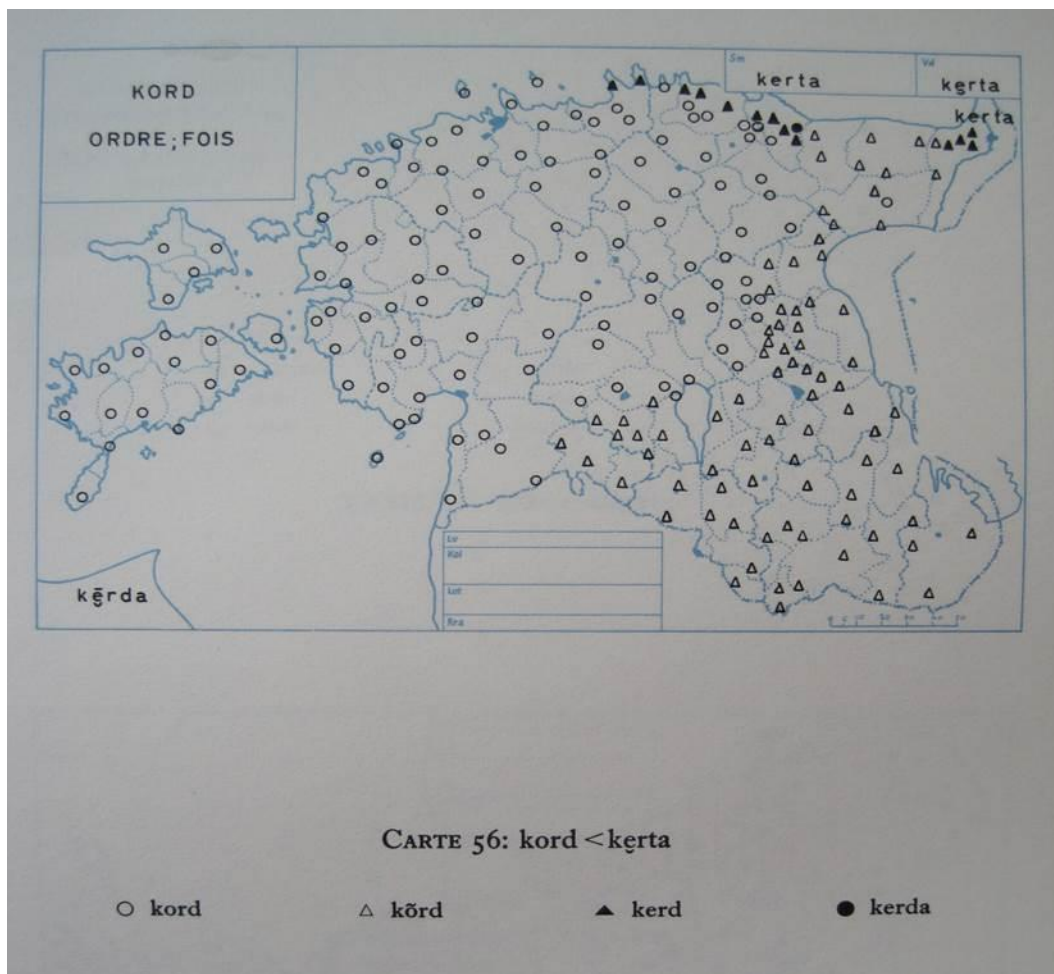**Digital Humanities in Estonia A° 2015**
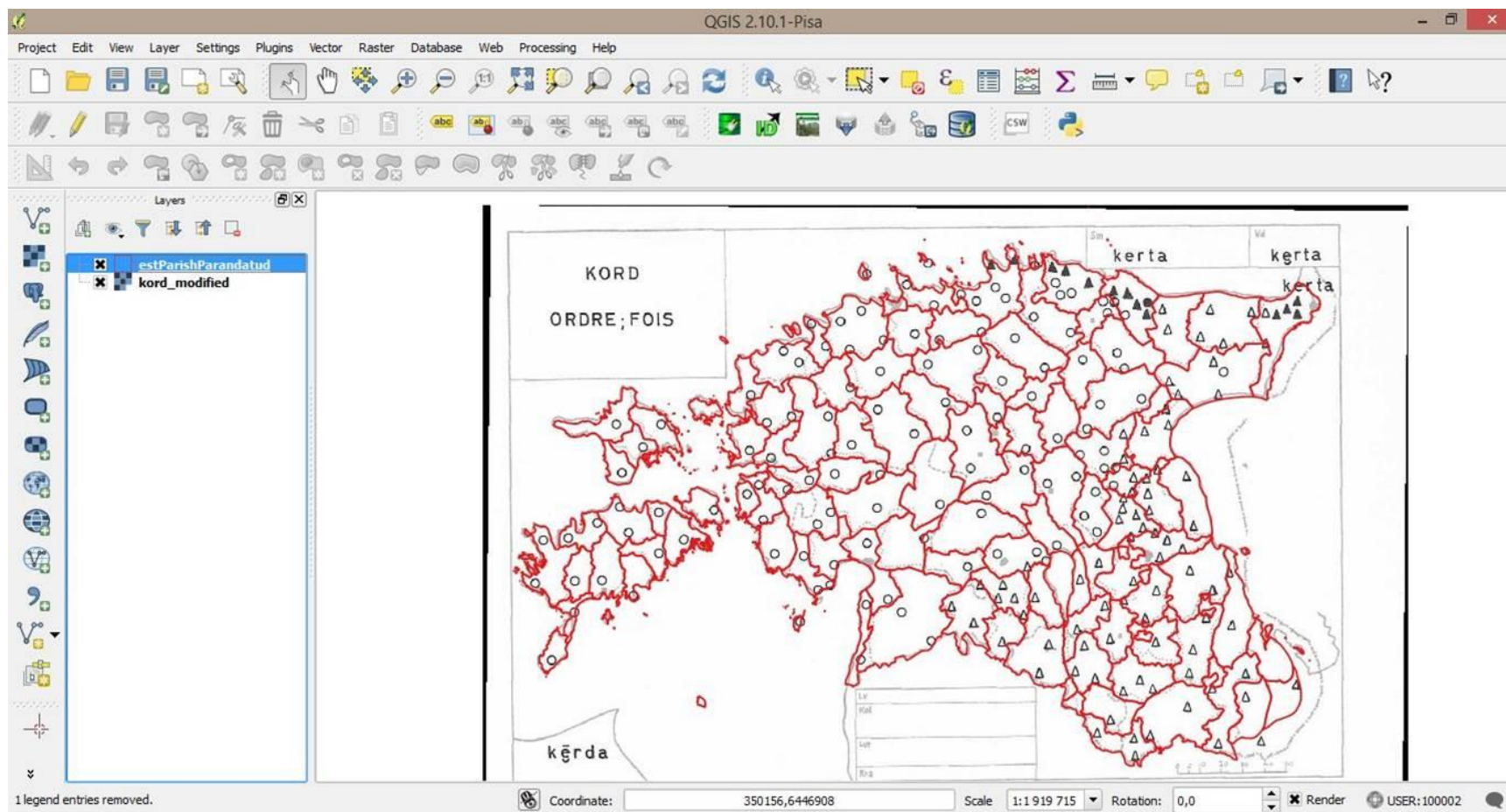
# Workflow: resources

- GIS for map digitalization
- R for analysis and visualization
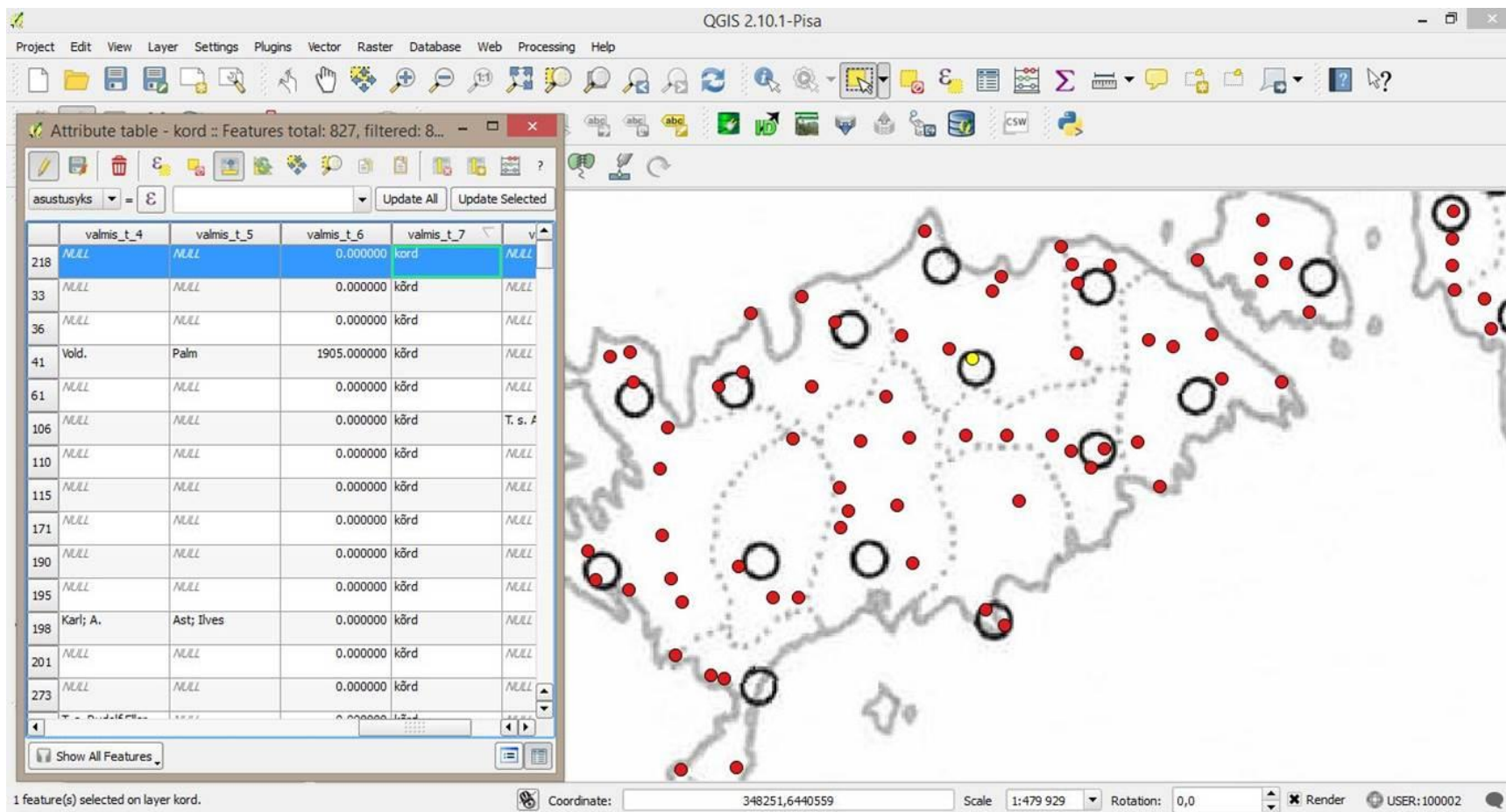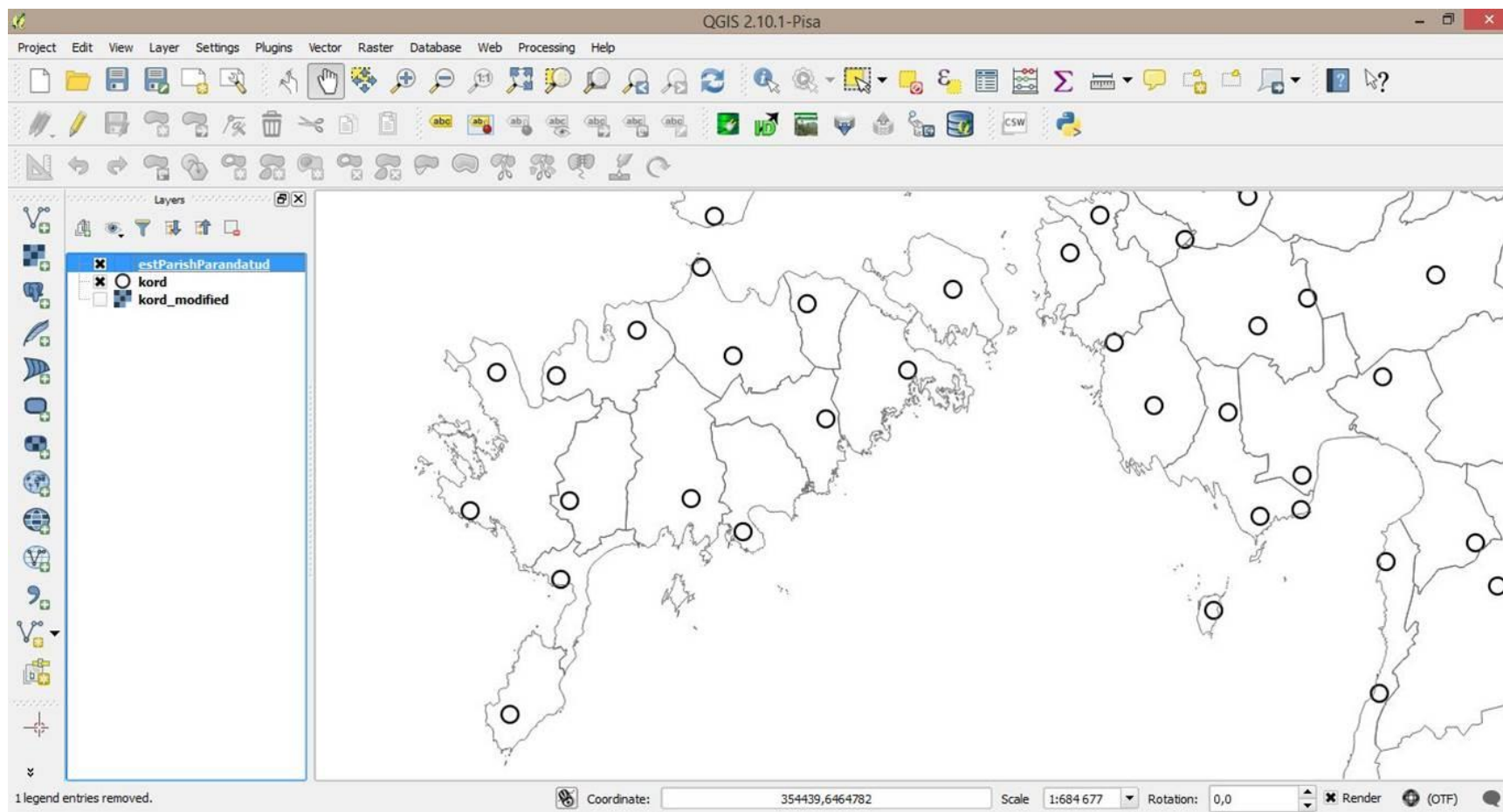- Free map resources (Maa-amet Geoportal, Google maps)

# Data insertion

- Geographical data collection for old data, aligning and verifying.
- Construct a shapefile map with the data collection points using the Maa-amet settlements base map (contemporary administrative division).
- Several villages in older data no longer exist → nearest settlement to mark that data collection point.
- Scan the maps from the *Petit Atlas des parlers estoniens* and georeference them.
- Insert the data on the map into the shapefile's attribute table.
- Design the final map (choose matching symbology, correct the position of the labels).

CARTE 56: kord < kęrta

○ kord     △ kõrd     ▲ kerd     ● kerda

# Corpus of Estonian Dialects

- www.murre.ut.ee/mkweb
- CED consists of dialect recordings (mainly from 1960-1970s), phonetically transcribed texts, dialect texts in simplified transcription, morphologically annotated texts, syntactically parsed texts, and a database with metadata about all the recordings and texts
- Older, non-mobile, less educated speakers
- 10 traditional dialect areas
- Morphologically annotated texts in XML format, which enables automated data extraction

# CED user interface

# CED map data

# Combining data sources

# Applications and conclusions

- CED data: frequencies and typicalities
- Atlas data: categorical
- Digitalized geographical information enables to compare this data
- Aggregating the data
- Linguistic distances based on both data sources
- Comparison of old data and new research findings
- Measuring linguistic and geographic distances

# References

CED = Corpus of Estonian Dialects.
*<http://www.murre.ut.ee/estonian-dialect-corpus/>*

Google Maps (2014) Google.

Kahle, David & Wickham, Hadley (2013) ggmap: A package for spatial visualization with Google Maps and OpenStreetMap. R package version 2.3.

*<http://CRAN.R-project.org/package=ggmap>*

EKI (2014) = Eesti Keele Instituudi kohanimeandmebaasi kihelkonnapiiride andmestik [Map of the Place Name Database created by the Institute of Estonian Language].

# References

Maa-amet Geoportal. *<http://geoportaal.maaamet.ee/eng/>*

QGIS Development Team (2015). QGIS Geographic Information System (*version* 2.10.1). Open Source Geospatial Foundation.

Saareste, Andrus (1938). Eesti murdeatlas. I vihik. Tartu: Eesti Kirjanduse Selts.

Saareste, Andrus (1941). Eesti murdeatlas. II vihik. Tartu: Teaduslik Kirjandus.

Saareste, Andrus (1955). Petit Atlas des parlers estoniens. Väike eesti murdeatlas. – Skrifter utgivna av Kungl. Gustav Adolfs Akademien 28. Uppsala: Almqvist & Wiksell.