# CORPORA OF CONTEMPORARY LATGALIAN TEXTS AND SPEECH
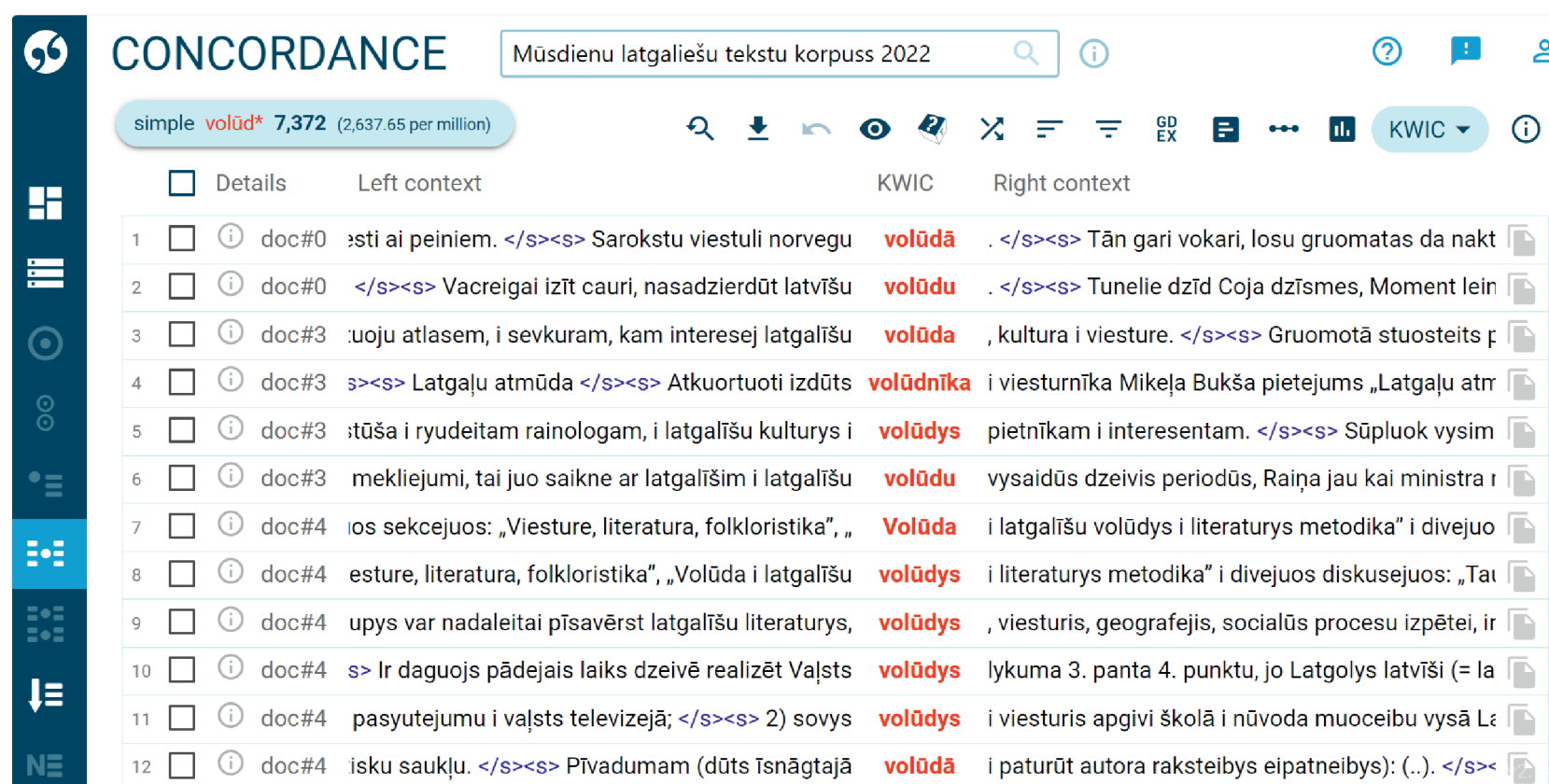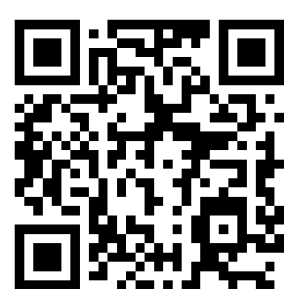
## MŪSDIENU LATGALIEŠU TEKSTU KORPUSS (MuLa)
### CORPUS OF CONTEMPORARY LATGALIAN TEXTS
**Website:** *https://korpuss.lv/id/MuLa2022*



- The current version **MuLa 2022** contains balanced samples of Latgalian texts published 1988–2021 with accompanying metadata about the author, place and year of the publication, as well as information about the type and genre of the text.
- Construction periods: 2011 – 2013; 2020 – 2022.
- Current size: 2,000,000 words, 2,800,000 tokens.
- Publishers: Rezekne Academy of Technologies, AiLab IMCS University of Latvia.
- Authors: Anna Briška, Ilze Ziņģe, Roberts Darģis, Krtistīne Pokratniece, Sanita Martena, Antra Kļavinska.
- Funding: National Research Programme project "Digital Resources for Humanities: Integration and Development" (No. VPP-IZM-DH-2020/1-0001).
- Corpus management platform: NoSketch Engine.
- The corpus is included in the Latvian National Corpora Collection (LNCC).

## MYUSU DĪNU LATGALĪŠU RUNYS KORPUSS
### CORPUS OF CONTEMPORARY LATGALIAN SPEECH
**Website:** *https://mularkorpuss.rta.lv*



- Construction periods: 202 – 2022, 202 – 2024.
- Current size: 23 hours of audio recordings, with transcribed text totalling 187296 tokens.
- Publisher: Rezekne Academy of Technologies.
- Authors: Sanita Martena, Nicole Nau, Antra Kļavinska, Angelika Juško-Štekele, Armands Kociņš-Kūceņš, Ausma Sprukte, Anna Briška, Ingars Gusāns, Laura Mazure.
- Funding: National Research Programme project "Digital Resources for Humanities: Integration and Development" (No. VPP-IZM-DH-2020/1-0001); National Research Programme project "Diversity of Latvian in Time and Space" (No. VPP-LETONIKA-2021/4-0003).
- Web Interface: SpoCo.

### Metadata:
- location and time of the recording,
- the duration of the audio segment,
- the speaker's gender and age.

**Content:**
- audio recordings:
  - interviews conducted during field work in Latgale and Siberia (2009–2021),
  - TV and radio broadcasts (2018–2023),
- transcriptions:
  - made using the ELAN software,
  - orthographic transcription that preserves dialect features.
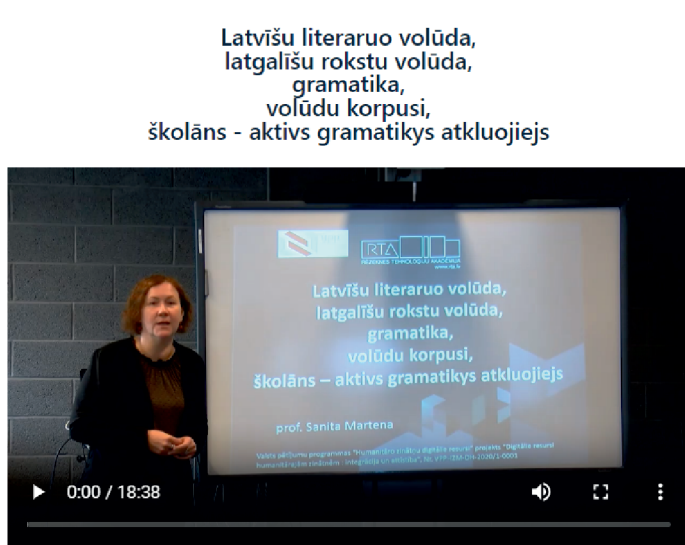
### DEVELOPMENT OF CORPUS LITERACY
- Study courses "Language corpora in the educational process", "Language technologies" for RTA teacher program students.
- Video lectures for Latgalian language teachers about the possibilities of using corpora.
  Website: *https://ltg.korpuss.rta.lv/video_lekcijas/*

### DEVELOPMENT OF CORPORA OF LESSER-USED LANGUAGE

**Opportunities and strengths:**
- new exciting ways of exploring Latgalian,
- long-term, stable cooperation with other European regional language research centers,
- adherence to FAIR data principles and inter-institutional cooperation in data management,
- raising awareness of the value of the Latgalian language and promoting awareness in the local community, especially among young people,
- fundraising, new projects, crowdsourcing.

**Challenges:**
- linguistic:
  - insufficient data (texts, audio recordings),
  - subjectivity factor in manual transcription of audio recordings,
  - variation as a challenge for the morphological tagging of corpora,
- extralinguistic:
  - lack of financial and human resources for the planned and sustainable digitization;
  - lack of interest / information among potential users,
  - lack of corpus literacy.

### Video lekcijas

Latvišu literaruo volūda, latgalīšu rokstu volūda, gramatika, volūdu korpusi, školāns - aktivs gramatikys atkluojiejs

Korpusprateiba i cytys digitaluos prasmis, ipazeistūt Latgolai rakstureigūs personvuordus (prīkšvuordus)