

## A corpus of Andi field recordings

Samira Verhees<sup>1</sup>

Aigul Zakirova<sup>2</sup>

George Moroz<sup>2</sup>

Elena Sokur<sup>2</sup>

<sup>1</sup>Independent researcher, the Netherlands; <sup>2</sup>HSE University, Russia

20 September 2022

Tbilisi, VI International Conference: Language and Modern Technologies

# Outline of the talk

- Corpora of Linguistic Convergence Laboratory
- An example of Meadow Mari corpus
- Andi language and its data
- Our path towards the corpus

# Corpora of Linguistic Convergence Laboratory

# Corpora of Linguistic Convergence Laboratory:

<http://lingconlab.ru/>

- Russian Dialect corpora (14)
- Corpora of Bilingual Russian (7)
- Corpora of minority languages of Russia (6)

# Corpora of Linguistic Convergence Laboratory:

<http://lingconlab.ru/>

Linguistic Convergence Laboratory eng ru  
Resources

All resources

Corpora ▾

Dictionaries ▾

Other Resources ▾

Khislavichi  
dialect

Tokens: 260 793

[View](#)

Keba dialect

Tokens: 54 535

[View](#)

Luzhnikovo  
dialect

Tokens: 68 666

[View](#)

Lukh and Teza  
river basins  
dialects

Tokens: 146 350

[View](#)

Malinino  
dialect

Tokens: 138 943

[View](#)

Nekhochi  
dialect

Tokens: 88 965

[View](#)

OPOCHETSKY  
dialects

Tokens: 68 741

[View](#)

Rogovatka  
dialect

Tokens: 100 047

[View](#)

Spiridonova  
Buda dialect

Tokens: 70 565

[View](#)

Shetnevo and  
Makeevo  
dialect

Tokens: 58 003

[View](#)

Tserkovnoe  
dialect

Tokens: 19 960

[View](#)

Upper Pinega  
and Vyga river  
basins dialect

Tokens: 70 803

[View](#)

Ustja River  
Basin dialects

Tokens: 959 782

[View](#)

Zvenigorod  
dialect

Tokens: 68 324

[View](#)

# Corpora of Linguistic Convergence Laboratory:

<http://lingconlab.ru/>

Linguistic Convergence Laboratory eng ru

All resourcesCorpora ▾Dictionaries ▾Other Resources ▾

Resources

Dialect corporaCorpora of bilingual RussianCorpora of minority languages of Russia

Bashkortostan Russian

Tokens: ND

View

Beserman Russian

Tokens: 97 216

View

Chuvash Russian

Tokens: 46 307

View

Daghestanian Russian

Tokens: 227 885

View

Karelian Russian

Tokens: 74 014

View

Romani Russian

Tokens: 41 767

View

Yakut Russian

Tokens: 15 139

View

5

# Corpora of Linguistic Convergence Laboratory: <http://lingconlab.ru/>

Linguistic Convergence Laboratory eng ru  
Resources

All resources

Corpora ▾

Dictionaries ▾

Other Resources ▾

Dialect corpora

Corpora of bilingual Russian

Corpora of minority languages of Russia

Abaza

Tokens: 3 636

View

Adyghe

Tokens: ND

View

Bashkir

Tokens: ~25 000

View

Kabardian

Tokens: ND

View

Khakas

Tokens: ~58 000

View

Meadow Mari

Tokens: ND

View

# Example of Meadow Mari corpus

[Meadow Mari Spoken Corpus](#)

EN | RU | ?



● Word #1

Word:

Lemma:

Grammar:

Gloss:

Language/tier: Meadow Mari



speaker: anf  
speaker\_name: Фёдорова Анна  
Николаевна  
dialect: моркинско-сернурский  
gender: f  
place of birth:  
year of birth: 1930  
education:  
year: 2018

Full-text search:

☒ Precise match

Search sentences

Search words / lemmata



Select subcorpus



Search result: 98 occurrences, 94 sentences found in approximately 10 documents.

**Биография А. Н. Ф. (часть 2)** 2018

а мый кум кече веле **коштына** ыле Реру-Ушчыл...

а мы за три дня доезжали

**Биография З. И. Е. (часть 2)**

<нрзб> **коштыт**...

Очень много ходило...

**Биография З. И. Е. (часть 2)** 2018

Эмланже гын пеш **коштына** ыле,

**коштына**  
кошташ V  
кошт-ына  
STEM-NPST.1PL  
gr: npst, 1, pl  
trans\_ru: ходить



1

2

3

4

...

10



## Example of Meadow Mari corpus

- translation
- glosses
- audio/video
- export to `.xlsx`
- sociolinguistic information

9

- hierarchy of tiers
- time alignment
- stem forms

Developing the corpus of Andi

# The Andi Language: a sociolinguistic background

- Andic < Avar-Andic < East Caucasian, glottocode [andi1255]
- spoken in several villages of the Botlikh district of Dagestan.
- more than 20,000 speakers of Andi [Aglarov (2002); All-Russian National Census 2010]
- Andi speakers are trilingual in Andi, Avar and Russian
- Avar serves as a lingua franca and is taught in school (Dobrushina et al. 2017)



	Andi	Rikvani	Gagatli	Zilo
materials	(Kibrik and Kodzasov 1988; Alekseev 1999)	(Sulejmanov 1957)	(Salimov 2010)	(Kaye et al. orth)
grammar sketch	+	+	+	+
dictionary	+	-	±	±
morphological parser	± <sup>1</sup>	-	-	-

<sup>1</sup>A pilot version of a morphological parser of Andi is presented in (Buntyakova 2022).

# Andic data

- How many hours do we have?
- How many is annotated?

- Andi has several dialects and no cross-dialectal standard;
  - no full-fledged dictionary
  - no full-fledged grammatical parser (though see a first attempt in (Buntyakova 2022))
- Our recorded data are heterogeneous
  - different dialects
  - different conventions
  - different file formats
- Due to our limited knowledge of the Andi dialects, sometimes we do not know what the correct analysis of a given word form is.

The material has to be converted to a singular format using `phonfieldwork` (Moroz 2020). For the Andi dialectal corpus the pipeline is as follows:

- we preprocess the annotation files, converting them to ELAN `.eaf` format (Wittenburg et al. 2006),
- align them with the sound
- gloss them manually or correct mistakes and ambiguities left by morphological parser
- publish online using the Tsakorpus platform (Arkhangelskiy 2019)
- repeat all previous steps



# Conclusions

## Conclusions:

Thank you for your attention!

- Aglarov, M. A. (2002). *Andijcy: Istoriko-etnografičeskoe issledovanie* [*The Andi people: a historical and ethnographic study*]. Jupiter, Makhachkala.
- Alekseev, M. E. (1999). Andijskie jazyki. In Alekseev, M. E., Starostin, S. A., Klimov, G. A., and Testeleets, J. G., editors, *Jazyki mira. Kavkazskie jazyki*, pages 220–228. Moskva.
- Arkhangelskiy, T. (2019). Corpora of social media in minority Uralic languages. In *Proceedings of the fifth Workshop on Computational Linguistics for Uralic Languages*. Association for Computational Linguistics, Tartu, Estonia.
- Buntyakova, V. A. (2022). Sozdanie morfologičeskogo parsera dla andijskogo jazyka v sisteme lexd i twol [morphological parser of andi in lexd and twol]. Term paper.

## References

- Dobrushina, N., Staferova, D., and Belokon, A. (2017). Atlas of multilingualism in dagestan online. <https://multidagestan.com>.
- Kaye, S., Moroz, G., Rochant, N., Verhees, S., and Zakirova, A. (forth). Andi (Zilo dialect). In Lander, Y., Maisak, T., and Koryakov, Y., editors, *The Caucasian Languages. An International Handbook*. De Gruyter Mouton, Berlin/New York.
- Kibrik, A. E. and Kodzasov, S. V. (1988). *Sopostavitelnoye izucheniye dagestanskikh yazykov* [Comparative study of Daghestanian languages]. Moscow State University, Moscow.
- Moroz, G. (2017). *lingtypology: easy mapping for Linguistic Typology*.
- Moroz, G. (2020). *Phonetic fieldwork and experiments with phonfieldwork package*.

- Salimov, H. S. (2010). *Gagatlinskij govor andijskogo jazyka* [*The Gagatli dialect of Andi*]. IJaLI, Makhachkala.
- Sulejmanov, J. G. (1957). Grammatičeskij očerk andijskogo jazyka. na materiale govora s. Rikvani [A sketch of the grammar of Andi. Based on material from the dialect of the village Rikvani].
- Wittenburg, P., Brugman, H., Russel, A., Klassmann, A., and Sloetjes, H. (2006). Elan: a professional framework for multimodality research. In *5th International Conference on Language Resources and Evaluation (LREC 2006)*, pages 1556–1559.