# Georgian
## spoken data

Stavros Skopeteas

University of Göttingen

Göttingen, October 25, 2021

# targets and contents

Goal

available Georgian corpus data

Contents

- resources with written texts

- Georgian spoken data

- recommendations

# written data

# Georgian National Corpus (GNC)

Georgian

rich collection of texts from various diachronic stages of Georgian
(Gippert & Tandashvili 2015)

http://gnc.gov.ge/gnc/page

| Corpus | Language(s) |
|---|---|
| **GNC Old Georgian** | 6 062 122 |
| **GNC Middle Georgian** | 1 432 262 |
| **GNC Modern Georgian** | 2 108 370 |
| **GRC** | 202 728 329 |
| **GDC** | 1 694 362 |
| **GNC Political texts** | 1 436 075 |
| **GNC Law texts** | 1 495 985 |
| **GNC Megrelian** | 89 404 |
| **GNC Svan** | 473 180 |

# Georgian Language Corpus, ILIA Univ. resources

## Georgian Language Corpus

http://corpora.iliauni.edu.ge/

texts of Modern Georgian (1832-2012), Old Georgian section, bilingual section with parallel corpora (Old Georgian – Old Armenian) such as the Georgian Chronicles (Doborjginidze & Lobzhanidze 2016).



## Epigraphic Corpus of Georgia

https://epigraphy.iliauni.edu.ge/en/

digitized ancient inscriptions in Greek, Aramaic, Armenian and Georgian (5 c. BC - 19 c. CE; Doborjginidze & Kalkhitashvili 2019)

## Wardrops' Collection Online

http://manuscript.iliauni.edu.ge/index.html

manuscript collection with an elaborate structure of textual metadata (Lobzhanidze 2018)

# Georgian dialect corpus

http://www.corpora.co/#/corpus

created at the TSU Arnold Chikobava Institute of Linguistics

- almost all dialects of Georgian

- Kartvelian languages

interface and contents in Georgian (also including a database of the Georgian linguistic terms)

http://www.corpora.co/#/corpus

# spoken data

# spoken data

Asatiani, Rusudan (recording/transcription,annotation) Stavros Skopeteas (design/supervision), Veronika Ries (recordings), Caroline Brokmann and Florian Fischer (revisions) 2019. Georgian spoken data corpus. *The Language Archive*, Corpus resource;

persistent identifier: https://hdl.handle.net/1839/00-0000-0000-0021-4DA3-5

# parallel corpus design

**Activity Description**
Please tell me how you are making a Chadzapuri. Do not worry if there are some details that clear description, such that another person can do the same.

**Ancestors**
Please tell me how do you imagine that the Ancient Georgians lived. It is not a problem if you Just tell me the story of your ancestors as far as you know it. If you do not know anything, ple that these people were living.

**Comparative Description**
Please tell me how you perceive the major differences between Georgian and Russian. If you just say to us how you perceive the major differences between these two languages. If you sp information about the difference in expressing yourself in both languages.

**Event Description**
Please tell me how did you enjoy the last New Year's feast: what did you prepare for the feas do, what did you think, what did you feel what happened.

**Path Description**
Please describe the path to go from Vake to Marjanishvili to me. Please give exact descriptior path that we have to follow (by telling me about all the important places on the way to Marjan houses, trees, crossroads, etc.). Please describe the path to go from Vake to Marjanishvili to descriptions, so that we can recognize the path that we have to follow (by telling me about all to Marjanishvili, e.g., characteristic houses, trees, crossroads, etc.).

GEO-TXT-PA-00000-15

GEO-TXT-PA-00000-16

GEO-TXT-PA-00000-17

GEO-TXT-PA-00000-18
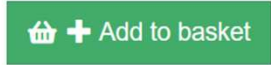
GEO-TXT-PA-00000-19

GEO-TXT-PA-00000-20

# download files

all data are accessible to registered users in the TLA (you need to create
an TLA account and login in order to download this data)

# open files

download ELAN

https://archive.mpi.nl/tla/elan

# process files

- phonological transcription (orth)

- glossing in English

- free translation (English)

# process files

- phonological transcription (orth)

- glossing in English

- free translation (English)

# exercise

- What do people describe on they way from Vake to Marjanishvili Square?

# ideas

This corpus is useful for examining phenomena that often occur in discourse, e.g.:

- phonetic realization of segments, syllable structure, intonation

- word order

- grammatical phenomena that frequently appear in this type of data (narratives/short dialogues).

This corpus does not have the necessary coverage for phenomena that are underrespresented in this type of texts, e.g., narratives/dialogues are not rich in commands. Constructions that are relevant for grammatical description and do not often occur in discourse may be poorly represented or not represented at all in a small-size corpus.

# references

Asatiani, R., S. Skopeteas, V. Ries, C. Brokmann, F. Fischer 2019. Georgian spoken data corpus. *The Language Archive*, Corpus resource; persistent identifier: https://hdl.handle.net/1839/00-0000-0000-0021-4DA3-5.

Doborjginidze, N. & I. Lobzhanidze. 2016. Corpus of the Georgian Language. *EURALEX 17*, 328-334.

Doborjginidze, N. & T. Kalkhitashvili 2019. The Epigraphic Corpus of Georgia. *Kadmos* 9, 222-233.

Gippert, J. & M. Tandashvili 2015. Structuring a diachronic corpus. The Georgian National Corpus project. In J. Gippert & R. Gehrke (ed.), *Historical Corpora. Challenges and Perspectives.* 305-322. Tübingen: Narr.

Lobzhanidze, I. 2018. Computational Model of the Modern Georgian Language and Search Patterns for an Online Dictionary of Idioms. *TbiLLC 2018*, 187-208.

# this lecture

is part of the series *Glottothèque: Languages of the Anatolia, Caucasus, Iran, Mesopotamia; grammatical snippets online*, ed. by. C. Bulut, A. Donabédian-Demopoulos, G. Haig, G. Khan, P. Samvelian, S. Skopeteas, N. Sumbatova. Bamberg/Cambridge/Göttingen/Moskow/Nicosia/Paris: LACIM network.



You may find related lectures and further information at the Glottothèque website at: https://spw.uni-goettingen.de/projects/lacim/