

# Comparing phonological systems and syllable structure of Botlikh and Zilo Andi: a data-driven analysis

G. Moroz

Linguistic Convergence Laboratory, NRU HSE, Moscow, Russia

25 February 2020, MPI-SHH, Jena

Presentation is available here: [tinyurl.com/rvpqdaa](https://tinyurl.com/rvpqdaa)



## Phonological description: data-driven analysis

	Traditional analysis	Data-driven analysis
1.	Done by trained linguist	Evaluated by trained linguist
2.	Can be done from scratch	Previous description needed (or at least prior expectations)
3.	Doesn't care about amount of data	Care more about amount of data
4.	Less reproducible	More reproducible
5.	Can not be automated	Can be automated

## Phonological description: data-driven analysis

	Traditional analysis	Data-driven analysis
1.	Done by trained linguist	Evaluated by trained linguist
2.	Can be done from scratch	Previous description needed (or at least prior expectations)
3.	Doesn't care about amount of data	Care more about amount of data
4.	Less reproducible	More reproducible
5.	Can not be automated	Can be automated

Data-driven approach to phonological description and syllable structure analysis:

- was proposed in ([Moroz 2018](#))
- was applied to syllable structure in ([Moroz 2019](#)) to Adyghe data
- was applied to syllable structure in ([Romanova 2019](#)) to Russian and Macedonian data

# Phonological description: data-driven analysis

The main steps:

- start with some language corpus (corpus, dictionary etc.)
- automatically extract phonological units (segments, features, syllable structure etc.)

# Phonological description: data-driven analysis

The main steps:

- start with some language corpus (corpus, dictionary etc.)
- automatically extract phonological units (segments, features, syllable structure etc.)

Further steps:

- obtained frequencies and inventories could be compared with the same units from other languages

# Phonological description: data-driven analysis

The main steps:

- start with some language corpus (corpus, dictionary etc.)
- automatically extract phonological units (segments, features, syllable structure etc.)

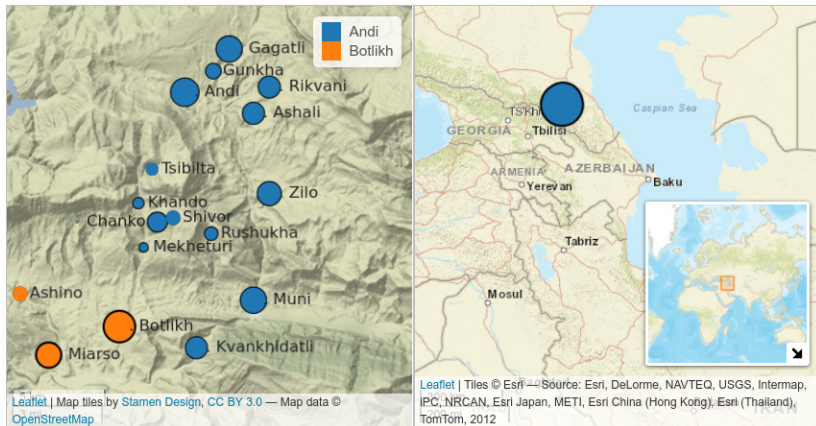
Further steps:

- obtained frequencies and inventories could be compared with the same units from other languages

Advantages:

- more reproducible
- could be updated with new data, see ([Moroz 2019](#)) slides from SLE on Bayesian typological research
- answers the question *‘How often is X present in language(s)?’* rather than *‘Is X present in language(s)?’*

# Andi and Botlikh villages



- Size of the dot corresponds with number of villages' inhabitants
- All villages except Botlikh are monoethnic
- Created with lingtypology package ([Moroz 2017](#))

---

**Botlikh < Andic group < EC**

---

Unwritten (can be written with extended Cyrillic script for Avar)

~5,000–8,000 speakers

Mostly spoken in 3 villages in northwestern Daghestan (Russian Federation): Botlikh, Miarso, Ashino, (Ankho); minor dialectal differences

One full reference grammar in Georgian ([Gudava 1962](#))

Two dictionaries: ([Saidova and Abusov 2012](#)), ([Alekseev and Azaev 2019](#))

---

---

**Andi < Andic group < EC**

---

Unwritten (can be written with extended Cyrillic script for Avar)

~16,500 speakers

About 14 villages; There are two main dialect groups: Lower Andi (Muni, Kvankhidatli) and Upper Andi (the rest)

Several reference grammars ([Suleymanov 1957](#)) (Rikvani), ([Salimov 1968](#)) (Gagatli), ([Tsertsvadze 1965](#)) (Andi)

No dictionary except ([Kibrik and Kodzasov 1988](#))

---



# Comparing two Botlikh dictionaries

(Saidova and Abusov 2012)

- Compiled in the 2000s by a native speaker (M. G. Abusov) and an experienced linguist (P. A. Saidova)
- Mostly Botlikh with some notes on Miarso

# Comparing two Botlikh dictionaries

(Saidova and Abusov 2012)

- Compiled in the 2000s by a native speaker (M. G. Abusov) and an experienced linguist (P. A. Saidova)
- Mostly Botlikh with some notes on Miarso

(Alekseev and Azaev 2019)

- Compiled in the 1960s / 1970s by a native speaker / philologist (X. G. Azaev) and later (in the 2000s) systematized by an experienced linguist (M. E. Alekseev)
- Subsequently edited by T. A. Maisak and scheduled for posthumous publication last year
- Botlikh only

# Comparing two Botlikh dictionaries

(Saidova and Abusov 2012)

- Compiled in the 2000s by a native speaker (M. G. Abusov) and an experienced linguist (P. A. Saidova)
- Mostly Botlikh with some notes on Miarso

(Alekseev and Azaev 2019)

- Compiled in the 1960s / 1970s by a native speaker / philologist (X. G. Azaev) and later (in the 2000s) systematized by an experienced linguist (M. E. Alekseev)
- Subsequently edited by T. A. Maisak and scheduled for posthumous publication last year
- Botlikh only

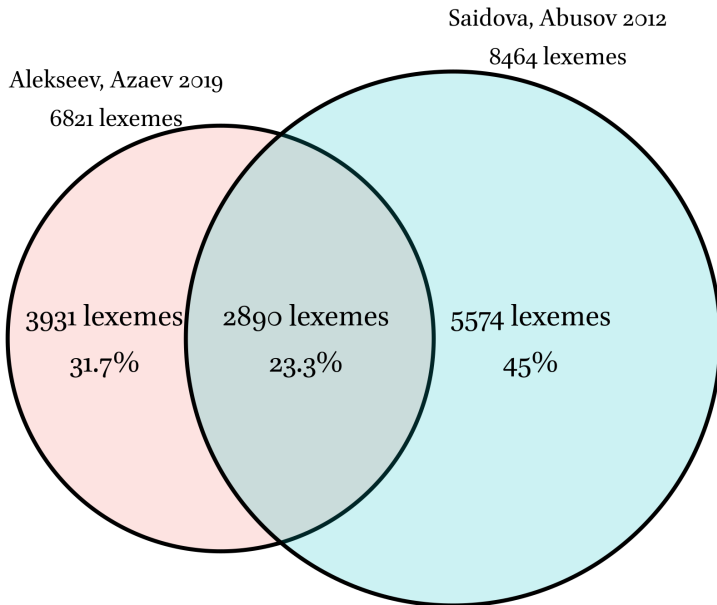
## Summary:

- Dictionaries were compiled independently of each other
- with no metadata on the speakers consulted
- data collection was separated with several decades break

## Comparing two Botlikh dictionaries: data preparation

- Automatically merge two `.doc` file into one unified `.xls` file
- Manually check for similarities (S. Verhees, C. Naccarato and me)

# Comparing two Botlikh dictionaries: data preparation



## Comparing two Botlikh dictionaries: results

- There are 1996 lexemes which look phonetically the same, and 909 are different (31%)
- If we remove the stress sign, there are 2449 lexemes which look phonetically the same, and 456 are different (16%)
- $\Rightarrow$  15% of lexemes have different stress pattern?..

## Comparing two Botlikh dictionaries: results

- There are 1996 lexemes which look phonetically the same, and 909 are different (31%)
- If we remove the stress sign, there are 2449 lexemes which look phonetically the same, and 456 are different (16%)
- $\Rightarrow$  15% of lexemes have different stress pattern?.. Yes, but including 265 (9%) cases where stress is present in one dictionary and absent in the other.

## Comparing two Botlikh dictionaries: results

- There are 1996 lexemes which look phonetically the same, and 909 are different (31%)
- If we remove the stress sign, there are 2449 lexemes which look phonetically the same, and 456 are different (16%)
- $\Rightarrow$  15% of lexemes have different stress pattern?.. Yes, but including 265 (9%) cases where stress is present in one dictionary and absent in the other.
- What causes the difference between dictionaries?
  - Stress pattern differences in 188 lexemes (about 6%)
  - Multiple cases where there is a small difference that could be explained either as a typo or in terms phonological variation: *čuhí* ‘to run’ [aa] vs. *čũhí* [sa], *kusu* ‘cherry plum’ [aa] vs. *kus:u* [sa]
  - Multiple cases where Russian borrowings were adopted differently: *awtobus* ‘bus’ [aa] vs. *abtabus* [sa], *biton* ‘milk can’ [aa] vs. *bitun* [sa], *apteka* ‘pharmacy’ [aa] vs. *abteka* [sa]
  - Morphological preferences: *dinija=w* ‘pious’ [aa] vs. *dinija=b* [sa]



## References

- Alekseev, M. and Azaev, X. (2019). *Botlixsko-russkij slovar'* [*Botlikh-Russian dictionary*]. Moscow: Academia.
- Gudava, Togo, E. (1962). *Botlixuri ena* [*The Botlikh language*]. Tbilisi: Mecniereba.
- Kibrik, A. E. and Kodzasov, S. V. (1988). *Sopostavitelnoye izucheniye dagestanskikh yazykov* [*Comparative study of Daghestanian languages*]. Moscow State University, Moscow.
- Moroz, G., A. (2018). lingphonology: automatic phonological description. R package draft.
- Moroz, G., A. (2019). Slogovaya struktura adygeyskogo yazika: ot dannyx k obosheniyam [*Adyghe syllable structure: From empirical data to generalizations*]. *Voprosy Jazykoznanija*, 2:82–95.
- Moroz, G. (2017). *lingtypology: easy mapping for Linguistic Typology*.

## References

- Romanova, K., I. (2019). Automatic Syllable Structure Extracting From Dictionaries: Slavic Data. Term paper.
- Saidova, P. A. and Abusov, M. G. (2012). *Botlixsko-russkij slovar'* [*Botlikh-Russian dictionary*]. Makhachkala: IJaLI.
- Salimov, X. S. (2010 (1968)). *Gagatlinskij govor andijskogo jazyka* [*The Gagatli dialect of the Andi language*]. Makhachkala.
- Suleymanov, J. G. (1957). *Grammatičeskij očerk andijskogo jazyka (po dannim govora s. Rikvani)* [*Grammar sketch of the Andi language (based on material from the dialect of the village Rikvani)*]. PhD thesis, Institut Jazykoznanja AN SSSR.
- Tsertsvadze, I. I. (1965). *Andiuri Ena*. Tbilisi: Metsniereba.