

Правительство Российской Федерации
Федеральное государственное автономное образовательное учреждение высшего
образования
Национальный исследовательский университет
«Высшая школа экономики»

Факультет гуманитарных наук
Образовательная программа
«Фундаментальная и компьютерная лингвистика»

Картина Элен Геннадьевич

**Правиловый морфологический парсер для шугнанского языка: существительные,
глаголы и прилагательные**

Выпускная квалификационная работа студента 4 курса бакалавриата группы БКЛ211

Академический руководитель образовательной программы
канд. филологических наук, доц.
Ю.А. Ландер

«_____» _____ 2025 г.

Научный руководитель
канд. филологических наук, доц.
Г.А. Мороз

Научный консультант
Стажёр-исследователь
М.Г. Мельченко

Москва 2025

Abstract

In this work I present a rule-based morphological analysis tool based on Helsinki Finite-State Technology (HFST) for the Shughni language (ISO: sgh; glottocode: shug1248), a language of the Iranian branch of the Indo-European family, a member of ‘Pamiri’ areal language group. While one existing rule-based parser exists for Shughni (Melenchenko, 2021), it does not utilize finite-state transducer technology. This work proposes the first HFST-based morphological parser implementation for Shughni, offering the advantages of this well-established framework for morphological analysis. The parser is presented in two variations: a morphological parser that breaks each word-form into stem and morphemes and assigns morphological tags to each one of them; a morphological generator that outputs word-forms taking a stem and morphological tags as an input. **TODO: prev sentence is questionable** This is a continuation my previous work, where nouns, pronouns, prepositions and numerals were implemented (Osorgin, 2024). This project covers **TODO: what**

TODO: Review abstract after finishing the work

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 1.1 | Shughni | 1 |
| 1.2 | Morphology modeling | 2 |
| 2 | Existing methods | 2 |
| 2.1 | Machine learning methods | 2 |
| 2.2 | Rule-based methods | 3 |
| 2.2.1 | Finite-state transducers | 3 |
| 2.2.2 | FST formalisms | 4 |
| 2.2.3 | Helsinki finite-state technology | 5 |
| 2.3 | Existing morphology models for Shughni | 5 |
| 3 | Data | 5 |
| 3.1 | Grammar descriptions | 5 |
| 3.2 | Dictionaries | 6 |
| 3.3 | Text corpora | 6 |
| 4 | Methods | 7 |
| 4.1 | Input and output format | 7 |
| 4.2 | FST compilation pipeline | 9 |
| 4.3 | lexd rule declaration | 12 |
| 4.4 | twol phonology | 13 |
| 4.5 | Stem lexicons processing | 14 |
| 4.6 | Transliteration | 15 |
| 4.7 | Russian lemmas TODO: rethink heading | 16 |
| 4.8 | Testing | 17 |
| 4.9 | Metrics | 19 |
| 5 | Results | 20 |
| 6 | Conclusion | 20 |
| | References | 21 |

1 Introduction

1.1 Shughni

The Shughni language (ISO: sgh; Glottolog: shug1248) is a language of the Iranian branch of the Indo-European family (Plungian, 2022, p. 12). As of June 1997, it was estimated to be spoken by approximately 100,000 people (Edelman & Yusufbekov, 1999, p. 225) in the territories of Tajikistan and Afghanistan. Both countries have a subregion where Shughni is the most widely spoken native language. The Shughni-speaking subregion of Tajikistan is called ‘Shughnon’ and it belongs to the ‘Gorno-Badakhshan Autonomous’ province. In Afghanistan, the Shughni-speaking region is called ‘Shughnan’ and it lies within the territory of ‘Badakhshan’ province (Parker, 2023, p. 2). Shughni belongs to ‘Pamiri’ areal language group, which is spoken along the Panj river in Pamir Mountains area.



Figure 1: Mountainous Badakhshan Autonomous Province of Tajikistan and Badakhshan Province of Afghanistan, (Parker, 2023, Fig 1.1)

There are three alphabets for Shughni that were derived from Cyrillic, Arabic and Latin scripts. Geographically the usage of said scripts correspond to the dominant script of each country where Shughni is spoken. In Tajikistan both official languages (Tajik and Russian) use Cyrillic script, so does Shughni on territory of Tajikistan. In Afghanistan Arabic script is used in Shughni, matching official languages (Pashto and Dari).

Latin script was developed and used in Tajikistan in 1930s (Edelman & Yusufbekov, 1999, p. 226)

(Edelman & Dodykhudoeva, 2009, p. 788), but according to Edelman and Yusufbekov (1999) was not widely adapted. Later around 1980s a Cyrillic script gained popularity in Tajikistan, having some poetic literature and school materials based on Tajik’s alphabet, which is Cyrillic (Edelman & Yusufbekov, 1999). Today, Latin script is mostly used by researchers in scientific works.

The morphological parser developed in this work is based on materials that focus on Shughni spoken in Tajikistan. All the base lexicon is Cyrillic and comes from dictionaries that cover Shughni in ‘Gorno-Badakhshan Autonomus Province’. Latin script is supported with the help of transliteration.

1.2 Morphology modeling

Today there are two general approaches to the task of morphology modeling. The deep learning (DL) approach and the rule-based approach.

The DL approach today typically makes use of training transformer models like BERT (Devlin et al., 2019) on vast amounts of marked-up data. This task becomes challenging, considering that Shughni is a low-resource language, meaning it lacks digital textual data. Although, DL approach was not utilized in this work, some existing DL approaches for low-resource languages are covered in section 2.1.

With the rule-based approach, morphological model is being built by writing grammar rules using some formalism language and by listing base lexicon. In this work, rule-based approach was utilized, as it does not depend on the amount of available marked-up data as the DL approach does. It requires lexicons and morphological grammar descriptions, which exist for Shughni and which are discussed in Section 3.

2 Existing methods

2.1 Machine learning methods

There are a variety of LLM (Large language model) architectures that were applied to the task of language modeling. One significant example is LSTM (Long short-term memory) model, that was introduced by Hochreiter and Schmidhuber (1997). LSTM is a variation of RNN (Recurrent neural network), and it was widely applied to language modeling, including morphology modeling. Another more recent significant example is the transformer architecture presented by Vaswani et al. (2017), off which two years later BERT model was based (Devlin et al., 2019).

One of the biggest downsides of ML methods is that its quality depends on training data quantity, which makes it challenging to apply to low-resource languages such as Shughni. However, with introduction of LLMs this problem was shown to be solvable, for example, as shown by developers of UDify model (Kondratyuk & Straka, 2019), which is a BERT-based model. In their work authors show, that their model pretrained on a large corpus of 104 languages can be fine-tuned on very little amounts of other languages’ data and still show decent results. For an example, they report that for Be-

larusian, UDify model achieved $UFeats = 89.36\%$ (accuracy of tagging Universal Features) after training on only 261 sentences from ‘Belarusian HSE’ Universal Dependencies treebank (Kondratyuk & Straka, 2019, Table 7).

However, working with LLM models is a highly resource-demanding task. The authors of UDify state, that the fine-tuning of their model for a new language would require at least 16 Gigabytes of RAM and at least 12 Gigabytes of GPU video memory, and the training process would take at least 20 days depending on the GPU model. While a deep learning approach would be interesting to explore, such computational resources are not available for this project. The neural approach is not the main target of this work and is implemented.

2.2 Rule-based methods

2.2.1 Finite-state transducers

The Rule-based approach historically is usually applied with the help of Finite-state transducers (FST), which is a variation of Finite-state machine, a mathematical abstract computational model. Following the terminology of Turing machines (Turing, 1937), a FST has two tapes: the input tape and the output tape. At any point it can read a next symbol from the input tape and then write a symbol to the output tape. Once a symbol was read from the input tape, it can not be read again, as the input tape shifts one symbol forward.

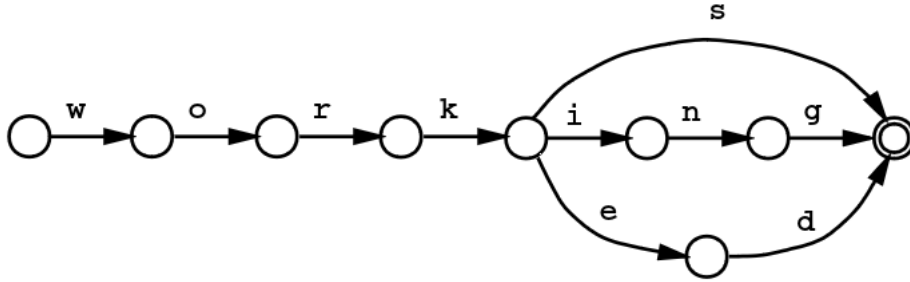


Figure 2: An example of FST with a single initial state (most left node) and a single final state (most right node) for a language where only three words exist: *works*, *working* and *worked*. The word *worker*, for example will not be recognized as a valid word by this FST, since there is no ‘d’ transition at state *work*. The only way from *work* state is via ‘d’ transition, which corresponds to the *worked* word. (Beesley & Karttunen, 2002)

The inner structure of FST can be illustrated as a directed graph with a set of all *states* (represented by graph’s nodes), a set of *transitions* (represented by graph’s edges), a set of *initial states* (a subset of all the states, these are states where FST can start reading from the input tape) and a set of *final states* (a subset of all the states, these are the states where FST can stop reading from the input tape). A simplified FST is shown on Figure 2. The letters above the graph’s edges denote *transition rules*, for an example *transition* ‘w’ means ‘read *w* from the input tape THEN write *w* to the output tape’.

While working, FST will only make transitions that are possible from the current state. If there

are no valid transitions then FST fails to process the input, and the input is considered to be impossible in the current language model. The measure of the amount of language’s grammatical wordforms that successfully pass through the FST from an *initial state* to a *final state* will be called *Coverage* from now and on. The measure of the amount of language’s ungrammatical wordforms that successfully pass through the FST from an *initial state* to a *final state* will be called *Overgeneration* from now and on. The ideal FST model of a language has a maximized 100% *Coverage* and a zero *Overgeneration*.

The model from the Figure 2 works effectively as a wordform paradigm dictionary, echoing back input wordforms that are grammatical and failing to output the whole ungrammatical wordforms. Now we can slightly adjust the transition rules in our example to make a morphological analysis tool that can be seen on the Figure 3. The notation of the *transition* ‘w:w’ dictates to read the left symbol from the input tape and write the right symbol to the output tape. If the spot on the right side is left empty, it means ‘write nothing to the output tape’. An important note to remember is that FST can output only one symbol to the output while making a single transition. In this example ‘<inf>’, ‘<pst>’ and ‘<prs><2sg>’ are ‘multichar’ symbols, meaning they are treated as three individual symbols by a FST, it will be covered in more detail in the Section 4.

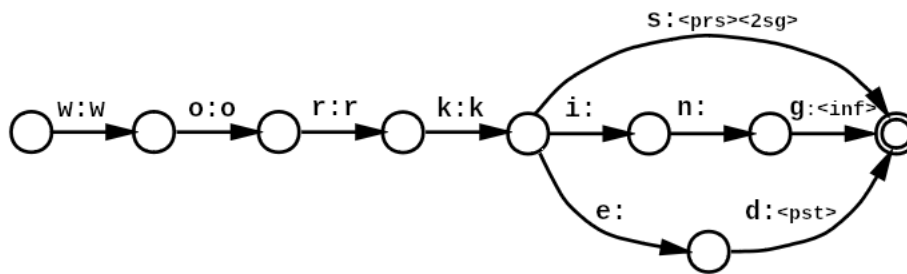


Figure 3: A modified version of Figure 2 which takes as input *works*, *working*, *worked* and outputs *work<prs><2sg>*, *work<inf>*, *work<pst>* respectively.

2.2.2 FST formalisms

By FST formalism I mean a human-readable formal language that can be compiled into a static FST, or from what a FST behavior can be emulated in runtime. A FST formalism usually includes a way to list lexicons and/or list lexicon combination rules and/or list phonological rules.

One of the first major fundamental advances was when (Koskenniemi, 1983) created a model, which introduced a FST formalism named Two-level morphology (TWOL) for describing morphological and morphonological paradigms. Its novice was in the addition of the phonology level of rules, which made it much easier and intuitive to implement cases like ‘*eye(-s)/*box(-es)*’. This model was capable of word-form recognition and production, but it was not yet compilable into static FSTs, it was working at runtime and was known for being slow. Then Karttunen et al. (1987) at Xerox Research Center developed a Two-level rule Compiler (`twolc`), which compiled TWOL rules into static FSTs. Later a separate compiler for lexicon definitions was introduced named `lexc` (**Lexicon Compiler**)

(Karttunen, 1993), it came with its own formalism language for describing lexicon and morphotactics. The standard approach to modeling a language at that point was using `lexc` to describe lexicon and morphology and `twolc` to describe morphonology, which stayed almost the same to this day.

One of the latest released tools was `lexd` lexicon compiler (Swanson & Howell, 2021). It is presented as a `lexc` alternative and is claimed to be much faster in the compilation time. It also introduced a tag system, which allows FST developers to specify how different lexicons combine with each other more precisely. **TODO: anything else?**

2.2.3 Helsinki finite-state technology

Helsinki finite-state technology (HFST) is a set of tools for creating and working with languages' morphology models in form of transducers (Lindén et al., 2009). It includes implementations of both `hfst-lexc` and `hfst-twolc` compilers, as well as command line interface commands for mathematical and other miscellaneous operations with transducers like FST combination and format conversion. Also, it comes with a file format `.hfst` designed to store compiled FSTs.

HFST is widely applied when it comes to creating rule-based morphological models. Some of the latest examples of HFST-based morphological tools are: morphological parser for the Tamil language by Sarveswaran et al. (2021), a morphological transducer for Kyrgyz by Washington et al. (2012), a morphological parser for Andi by Buntyakova (2023) and a morphological parser for the Chamalal language by Budilova (2023).

2.3 Existing morphology models for Shughni

At this time only one morphological parser exists for Shughni. It was developed by Melenchenko (2021) and was later included in 'Digital Resources for the Shughni Language' project (Makarov et al., 2022). It is a rule-based parser implemented in Python which shows good coverage and accuracy results. The main difference from the parser presented in this work is that Melenchenko's parser is not based on FST technology.

3 Data

3.1 Grammar descriptions

Several Shughni grammar descriptions were written throughout the years, starting from basic grammar description done by D. L. Ivanov (Salemman, 1895, pp. 274–281). An important mention is a work by Karamshoev (1963), which was the most detailed Shughni grammar description of its time. Latest significant works were 'Shughni language' (Edelman & Yusufbekov, 1999, pp. 225–242), 'Comparative Grammar of Eastern Iranian Languages' (Edelman, 2009) and 'A grammar of the Shughhi language' by Parker (2023), which is the biggest existing grammar, the most detailed and the most recent one.

For this work the main reference for compiling Shughni grammar rules was the work of Parker (2023). It was picked as it is the most recent, the most detailed and the biggest one. Other grammar description works were used too, but only as a secondary reference. The second most used grammar description was a work by Edelman and Yusufbekov (1999).

3.2 Dictionaries

There are two main dictionaries of the Shughni language: one by Zarubin (1960) and one by Karamshoev (1988–1999), both are written using Cyrillic script and include Russian translations. Some early dictionaries are ‘Brief grammar and dictionary of Shughni’ (Tumanovich, 1906), that is also using Cyrillic and translates to Russian, and ‘Shughni dictionary by D. L. Ivanov’ (Salemman, 1895), that translates to Russian but uses Arabic script alongside Cyrillic transcriptions for Shughni word-forms.

An important lexical source for this work was the ‘Digital Resources for the Shughni Language’ project (Makarov et al., 2022). As a part of their work, authors compiled a digital dictionary for Shughni, where they digitalized both major Shughni dictionaries by Karamshoev (1988–1999) and Zarubin (1960). The digital dictionary is available at their website via a web-interface, but I was given access by the authors to a copy of the underlying database, which simplified the process of exporting lexicons for this project. All the lexicons for FST compilation were taken from their database.

3.3 Text corpora

I was given access to unpublished native texts that were gathered during HSE expeditions to Tajikistan in 2019-2024. It is not a large corpus of texts. It consists of a translation of ‘The Gospel of Luke’ in Shughni, of a ‘Pear Story’ text, which is a spoken text that was written down from a retelling of the ‘Pear Story’ movie during an expedition, and it also includes a group of miscellaneous untitled small texts. Texts size can be seen on Table 1

| Text name | Total tokens | Unique tokens |
|-----------------------------|--------------|---------------|
| <i>‘The Gospel of Luke’</i> | 2978 | 1001 |
| <i>‘Pear Story’</i> | 1117 | 438 |
| Miscellaneous texts | 164 | 106 |
| All texts | 4259 | 1393 |

Table 1: A list of native Shughni texts and their sizes gathered during HSE expeditions to Tajikistan in 2019-2024

The database provided by Makarov et al. (2022) also contained a lot of different useful data parsed from dictionaries including dictionary entries’ usage examples. Such data is not as valuable as native texts, as sometimes it might not come from a native speaker but from a researcher. I would argue that for *Coverage* evaluation it might still be quite useful.

From materials of HSE expeditions to Tajikistan I also acquired manually glossed texts in .eaf (ELAN) format. These texts were utilized for *Accuracy* evaluation, which will be discussed in Section 4.9 along with *Coverage*. A full list of text sources can be seen on Table 2.

| Text name | Total tokens | Unique tokens | Native | Glossed |
|-----------------------------|--------------|---------------|-----------|---------|
| Dictionary examples | 164 225 | 29 013 | Uncertain | No |
| <i>‘The Gospel of Luke’</i> | 2 978 | 1 001 | Yes | No |
| <i>‘Pear Story’</i> | 1 117 | 438 | Yes | No |
| Miscellaneous texts | 164 | 106 | Yes | No |
| <i>‘The Gospel of Luke’</i> | 2 942 | 635 | Yes | Yes |
| <i>‘Pear Story’</i> | 228 | 83 | Yes | Yes |
| <i>‘Mama’</i> | 267 | 123 | Yes | Yes |

Table 2: A list of all available digital textual data

4 Methods

4.1 Input and output format

Every FST converts between two types of strings: wordforms (e.g. ‘дарйойен’ \approx ‘rivers’) and glossed strings (e.g. ‘дарйо<n>><pl>’ = ‘дарйо.N-PL’). Glossed strings format is based on Apertium format **TODO: citation needed?**, which is a standardized format for HFST-based transducers. The choice of this format is motivated by the fact that this is the standard, and I want to keep it consistent with existing tools and practices in the field. Important keys of the format are:

- Grammatical tags are enclosed with angular brackets and are lowercase.
stone.V = stone<v>
- Part of speech tags (POS) are obligatory and stand next to the stem or lemma.
stone.SG = flood<n>><sg>
- Different morpheme tags are separated with a single right angular bracket ‘>’.
stone-PL = stone<n>><pl>
- Multiple stems in a single word is possible.
lemma1<adj>><morph>>lemma2<adj>

The Shughni morphological parser in this work is presented in a variety of input and output formats. A full list of .hfst files with their formats of input and output is shown in Table 3. Motivation for this many variations of the morphological parser was support of different formats. Parser comes with four format variables:

- FST directionality: analyzer or generator. This simply shows the direction of a FST, analyzers take wordforms as input and return glossed strings, generators take glossed strings and return wordforms
- Glossed string stem glosses: Shughni stems or Russian lemmas. Notated in file names as `stem` and `rulem` respectively.
дарйо<n>><pl> (Shughni stem); река<n>><pl> (Russian lemma)
- Wordform morpheme borders: plain word or morphemes are separated with ('>') delimiter symbol. Notated in file names as `word` and `morph` respectively.
дарйойен (plain word); дарйо>йен (morphemes separated)
- Wordform script: Latin or Cyrillic.
дарйойен (Latin); daryoyen (Cyrillic)

| Transducer file name | Input example | Output example |
|----------------------------------|---------------|----------------|
| sgh_gen_stem_morph_cyr.hfst | дарйо<n>><pl> | дарйо>йен |
| sgh_gen_stem_word_cyr.hfst | дарйо<n>><pl> | дарйойен |
| sgh_gen_rulem_morph_cyr.hfst | река<n>><pl> | дарйо>йен |
| sgh_gen_rulem_word_cyr.hfst | река<n>><pl> | дарйойен |
| sgh_gen_stem_morph_lat.hfst | дарйо<n>><pl> | daryo>yen |
| sgh_gen_stem_word_lat.hfst | дарйо<n>><pl> | daryoyen |
| sgh_gen_rulem_morph_lat.hfst | река<n>><pl> | daryo>yen |
| sgh_gen_rulem_word_lat.hfst | река<n>><pl> | daryoyen |
| sgh_analyze_stem_morph_cyr.hfst | дарйо>йен | дарйо<n>><pl> |
| sgh_analyze_stem_word_cyr.hfst | дарйойен | дарйо<n>><pl> |
| sgh_analyze_rulem_morph_cyr.hfst | дарйо>йен | река<n>><pl> |
| sgh_analyze_rulem_word_cyr.hfst | дарйойен | река<n>><pl> |
| sgh_analyze_stem_morph_lat.hfst | daryo>yen | дарйо<n>><pl> |
| sgh_analyze_stem_word_lat.hfst | daryoyen | дарйо<n>><pl> |
| sgh_analyze_rulem_morph_lat.hfst | daryo>yen | река<n>><pl> |
| sgh_analyze_rulem_word_lat.hfst | daryoyen | река<n>><pl> |

Table 3: A full list of available HFST transducers

Four binary variables result in 16 ($= 2^4$) FST variations. Every FST listed in Table 3 can be built with `make` command as shown in Code block 1. Regular `.hfst` binary FSTs are not recommended using in production environments, as they are not optimized. For production use an optimized format called `.hfstol` (HFST **O**ptimized **L**ookup), which can also be compiled automatically for every listed FST using `make`.

```
$ make sgh_gen_stem_morph_cyr.hfst
$ make sgh_gen_stem_morph_cyr.hfstol
```

Code 1: Example of FST compilation with make.

4.2 FST compilation pipeline

Analyzers and generators

First, there are two main types of FSTs: generators and analyzers. They differ only in the directionality of a FST. Analyzers take wordforms as input and return glossed strings as output. Generators work in reverse, as shown in Code block 2.

```
$ echo "дарйойен" | hfst-lookup -q sgh_analyze_stem_word_cyr.hfst
дарйойен      дарйо<n>><3pl>  0.000000
дарйойен      дарйо<n>><pl>   0.000000
$ echo "дарйо<n>><pl>" | hfst-lookup -q sgh_gen_stem_word_cyr.hfst
дарйо<n>><pl>  дарйо-йен      0.000000
дарйо<n>><pl>  дарйойен      0.000000
```

Code 2: FST analyzer vs generator output formats.

The `lexd` source code is written as a generator, meaning by default, compiled FST takes glossed stem or lemma as input and returns a wordform. To compile any analyzer, a corresponding generator is inverted, as shown in Code block 3.

```
$ hfst-invert sgh_gen_stem_word_cyr.hfst \
-o sgh_analyze_stem_word_cyr.hfst
```

Code 3: FST analyzer creation from a FST generator.

Shughni stems and Russian lemmas

The next format only applies to the glossed side of FSTs (to the analyzers' output and to the generators' input). It sets whether a Cyrillic Shughni stem or a Russian translated lemma will be used as a stem's gloss, as shown in Code block 4. Shughni stems can have multiple Russian candidates (*дарйо* can be translated as *пека* or *мопе*). This leads to composed transducer having more output candidates. This works both ways, meaning Russian lemmas can translate as multiple Shughni stems (*пека* can be translated as *дарйо* or *хау*).

```

$ echo "дарйо<n>><3pl>" | hfst-lookup -q sgh_gen_stem_word_cyr.hfst
дарйо<n>><3pl>   дарйо-йен      0.000000
дарйо<n>><3pl>   дарйойен      0.000000
$ echo "река<n>><3pl>" | hfst-lookup -q sgh_gen_rulem_word_cyr.hfst
река<n>><3pl>     дарйо-йен      0.000000
река<n>><3pl>     дарйойен      0.000000
река<n>><3pl>     ҳац-ен        0.000000
река<n>><3pl>     ҳацен         0.000000
$ echo "дарйойен" | hfst-lookup -q sgh_analyze_rulem_word_cyr.hfst
дарйойен         море<n>><3pl>  0.000000
дарйойен         море<n>><pl>   0.000000
дарйойен         река<n>><3pl>  0.000000
дарйойен         река<n>><pl>   0.000000

```

Code 4: Shughni stem vs Russian lemma versions of FST.

The `lexd` source code contains lexicons with Shughni stems on the glossed side, meaning default compiled FST contains only Shughni stems. The process of creating a FST that works with Russian lemmas on the glossed side is more complicated. It is achieved with the help of a second FST `rulem2sgh.hfst`, its only purpose is translating stems to lemmas. It is attached to the input of a generator FST, creating a pipeline

`‘река<n>><pl>’ → rulem2sgh → ‘дарйо<n>><pl>’ → generator → ‘дарйойен’`

It can be done with ‘compose’ transducer operation (see Code block 5), which takes two FSTs, directs first’s output to the second’s input and returns the resulting composed FST. Details of the translator FST’s development are described in Section 4.7.

```

$ hfst-compose rulem2sgh.hfst sgh_gen_stem_word_cyr.hfst
-o sgh_gen_rulem_word_cyr.hfst

```

Code 5: Shughni stem translator composition.

Latin script support

The source code contains only Cyrillic lexicons. The support of Latin script comes with the help of a separate transliterator FST, as it was the case for Russian lemmas support. Transliteration is only applied to the wordform side of FSTs. The glossed side’s stems and lemmas are left Cyrillic. The

pipeline for transliteration is below:

‘дарйо<n>><pl>’ → generator → ‘дарйойен’ → cyr2lat → ‘daryoyen’

The compilation process is the same as it is for Russian lemmas shown in Code block 5. The only difference is that generator comes first and transliterator comes second, as it is applied to the wordform instead of the glossed side. See Section 4.6 for more details about the transliterator development.

Morpheme borders

Theoretical linguists sometimes need wordforms to have morphemes separated by a special symbol that does not appear naturally in a language. In this morphological parser I choose the right angular bracket symbol ‘>’ for this role to keep it consistent with glossed strings morpheme separator. This presented a slight challenge to the development, as regular wordforms often may contain hyphens (‘-’) between some morphemes. For an example, ‘дарйойен’ can also be spelled as ‘дарйо-йен’. But for the morpheme separated FST wordform must contain only ‘дарйо>йен’, with no extra optional hyphens.

This was solved with the help of `twol` rules. The base `lexd` FST (`sgl_base_stem.hfst`) contains both hyphens and morpheme separators, which looks like ‘дарйо>-йен’ and ‘дарйо>йен’ on the wordform side. Then the filtering is done with help of two FSTs that contain `twol` rules shown in Code block 6. The `sep.hfst` removes all separator symbols (‘>’) from the wordform, leaving ‘дарйо-йен’ and ‘дарйойен’. And `hyphen.hfst` removes all hyphens from the wordform, leaving ‘дарйо>йен’ and ‘дарйо>йен’, which then fold into a single FST path after optimization.

```
# hyphen.twol
Rules
"Hyphen removal to avoid '>-'/'->' situations"
% -:0 <=> _ ;

# sep.twol
Rules
"Morpeme separator removal"
%>:0 <=> _ ;
```

Code 6: `twol` rules that filter morpheme border.

For this case compose-intersect operation is applied, that allows `twol`-compiled FSTs to be composed with regular lexicon `lexd` FSTs. Compilation command is shown in Code block 7.

This way, every FST listed in Table 3 can be compiled using a specific combination of bash HFST tools shown in this section. For more details explore the `Makefile` in the repository (Kartina, 2025).

TODO: How to cite properly? Work vs Repo

```
$ hfst-compose-intersect sgh_base_stem.hfst twol/sep.hfst  
-o sgh_gen_stem_word_cyr.hfst  
$ hfst-compose-intersect sgh_base_stem.hfst twol/hyphen.hfst  
-o sgh_gen_stem_morph_cyr.hfst
```

Code 7: Plain text and morpheme separated text FST versions compilation.

4.3 `lexd` rule declaration

The choice of lexicon compiler was made in favor of `lexd` as it provides everything that `lexc` does and in addition has some extra useful functional in form of the tag system, which will be taken advantage of.

The `lexd` source code is stored in the `lexd/` directory. I decided to go with a modular file structure for `lexd` source code, as it helps to keep the source code organized. The `lexd/` directory contains `.lexd` source code files with morpheme lexicons (suffixes, clitics, prefixes, etc.) and lexicon combination patterns. Stem lexicons are stored separately in the `lexd/lexicons/` directory. For the most part `lexd/lexicons/` directory contains lexicons obtained from database dumps provided by Makarov et al. (2022). The stem lexicon processing is described in Section 4.7.

There is no module import feature in `lexd`. So in order to be able to make a modular `.lexd` source file structure compilable into a single `.hfst` file we can concatenate every `.lexd` module into a single large temporary `.lexd` file and feed it to the compiler. This is achieved with `bash` command shown in Code block 8. The `lexd` compiler outputs FST in AT&T format and `hfst-txt2fst` converts it to a binary `.hfst` file.

```
$ cat lexd/*.lexd lexd/lexicons/*.lexd > sgh.lexd  
$ lexd sgh.lexd | hfst-txt2fst -o sgh_base_stem.hfst
```

Code 8: Bash command pipeline compiling multiple `.lexd` files into a single FST.

TODO: Finish when done!!

Nouns

Verbs

Adjectives

Pronouns

Numerals

Anything else(?)

4.4 twol phonology

Shughni is not rich for morphonological rules. The only global phonological rule is ‘й’ (Latin: ‘y’; IPA: /j/) drop (**TODO: citation needed; ‘drop’ is a questionable name for this**) on some morphemes’ borders after a consonant:

$$j \rightarrow \emptyset / [+consonant] _$$

TODO: как показать границу морфемы, нужна ли вообще нотация правила если я не могу на неё сослаться

A special multichar symbol can be introduced to stand in place of ‘й’ letters that are affected by this phonological rule: ‘{й}’. We capitalize it and frame it with curly brackets so it is never confused with plain letters. A FST treats multichar symbols as single unique symbols, even though they visually consist of multiple characters. This feature allows us to mark which lexicon entries are affected by this phonological rule and which are not. An example of lexicon definition with this symbol is shown in Code block 9.

```
LEXICON PCS # Personal clitics
<1sg>:{й}ум
<2sg>:{й}ат
<2sg>:т
<3sg>:{й}и
<1pl>:{й}ам
<1pl>:{й}ам
<2pl>:{й}ет
<3pl>:{й}ен

LEXICON DIM # Diminutive clitic
<dim>:{й}ик
<dim>:{й}ак
```

Code 9: A real lexd example of ‘{й}’ multichar usage.

Without twol rules these morphemes will remain as they are specified in lexd, meaning that

feeding ‘вирод<n>><dim>’ (=brother.N-DIM) to the input of a generator will output literally ‘вирод{й}ик’. After applying a `twol` rule (shown in Code block 10) it generates wordforms ‘виродик’ (=вирод.N-DIM/brother.N-DIM) and ‘туйик’ (=ту.PERS.2SG-DIM/you.PERS.2SG-DIM). The composition of `twol` rules with the main FST is the same as it was for morpheme border filtration (Code block 7). This rule applies to all morphemes that can start with ‘j’. Including, but not limited to personal clitics shown in Code block 9, noun adjectivator suffix ‘-jin’ (Parker, 2023, p. 168), noun plural suffix ‘-yen’

```
"j drop after consonant"
%{Й%}:0 <=> Consonant (%>) (%-) (%>) _ ;
```

Code 10: A `twol` rule for ‘j’ drop depending on the previous morpheme’s segment. Symbol ‘%’ in `twol` is used to escape a character.

There are other phonological rules that do not apply to such variety of contexts:

- Two clitics drop a vowel ‘a’ when following a vowel: ‘=(a)нд’ (locative) and ‘=(a)рд’ (dative)

TODO: citation needed. `twol` rule looks like this:

```
%{А%}:0 <=> Vowel (%>) (%-) (%>) _ ;
```

- A verb PRS suffix ‘-д/-т’ becomes ‘-д’ following voiced obstruents and vowels and ‘-т’ in all other contexts (Parker, 2023, p. 262). `twol` rule looks like this:

```
%{ДТ%}:д <=> DGroup _ ;
```

Where `DGroup` stands for all vowels and voiced obstruents.

4.5 Stem lexicons processing

Describing morphological rules and morphemes is not sufficient to make a morphological parser. An important part of the development was gathering stem lexicons for various parts of speech. For some POS, like conjunctions or numerals, lexicons were listed by hand. But for others, like verbs, manual lexicon construction is an extremely time-consuming task. Such lexicons were extracted from the `SQLite` database provided by Makarov et al. (2022).

The database in question contains digital versions of Shughni main dictionaries: Karamshoev (1988–1999) and Zarubin (1960). It has a decently organized structure, which allows selecting dictionary entries by parsed metadata like POS or glosses. To transfer stems from the database to a `lexd`-formatted file it follows the following pipeline:

```
Database → noun.csv → form_lexd.py → noun.lexd
```

An example of real `csv` and resulted `lexd` fragments are shown in Code block 11. I decided to keep meanings from the source as `lexd` comments. Later it makes it easier to debug lexicons manually, this should not influence the compilation process as comments are ignored. The

Exported noun.csv fragment

```
cyrillic,meaning
аббуст,рукавица
аббустеџ,материал
абед,обед; время обеда
абйн,ненавистница
абкѹмпарт,обком партии
```

Generated noun.lexd fragment

```
LEXICON LexiconNoun[pl_all,sg]
аббуст      # рукавица
аббустеџ    # материал
абед        # обед; время обеда
абйн        # ненавистница
абкѹмпарт   # обком партии
```

Code 11: An example of exported noun stems from the database converted to `lexd` source code.

`form_lexd.py` script is located in `scripts/lexicons/` directory. Besides formatting it also preprocesses the Shughni stems: removes stresses, unifies Cyrillic script and filters out affix entries (source dictionaries have entries for affixes alongside with regular stems). For convenience the script also sorts stems alphabetically.

The `form_lexd.py` script does not read command line arguments. The configuration is ‘hard-coded’ to read all the `.csv` files from `scripts/db_dumps/` directory and save the generated `.lexd` files to `lexd/lexicons/`.

4.6 Transliteration

The transliteration in this work is implemented for Latin and Cyrillic scripts of the Shughni language. It is developed as a pair of separate FSTs: `translit/lat2cyr.hfst` (Latin to Cyrillic transliteration) and `translit/cyr2lat.hfst` (Cyrillic to Latin transliteration). I found it more convenient to compile transducers with `lexd` (opposed to raw ‘strings’ format). `Lexd` allows splitting characters and special symbols into separate lexicons and writing comments, which makes the source file much more readable. Also, as a side bonus, `lexd` syntax allows repeating (making it infinitely cyclic) the transducer without the need to repeat it via HFST tools in the `Makefile`.

The transliterator only works for wordforms (e.g. ‘daryoyen’=*river.N-PL*). It does not work for glossed strings (e.g. ‘daryo<n>><pl>’). The reason for this is that in case of Latin to Cyrillic transliteration the FST would also transliterate grammatical tags, outputting ‘дарйо<н>><пл>’ given ‘daryo<n>><pl>’ as input. This could be solved by listing all the possible grammatical tags as multichar symbols in the transliterator’s lexicon, but this seems unnecessarily complicated to me for transliteration purposes.

As shown in Code block 12, transliteration for wordforms works for both plain text wordforms and morpheme-separated wordforms (the difference was explained in Section 4.2: *Morpheme borders*). This ensures that the transliteration FST can be successfully composed with any version of wordform

side of FSTs.

```
$ cat lat.txt | hfst-lookup -q translit/lat2cyr.hfst
na>čis>en    на>чис>ен    0.000000
na-čis-en    на-чис-ен    0.000000
načisen      начисен      0.000000
$ cat cyr.txt | hfst-lookup -q translit/cyr2lat.hfst
на>чис>ен    na>čis>en    0.000000
на-чис-ен    na-čis-en    0.000000
начисен      načisen      0.000000
```

Code 12: An example of transliteration FST work. The mentioned word ‘na>čis>en’ can be glossed as *NEG-look. V.PRS=3PL*

4.7 Russian lemmas **TODO: rethink heading**

The Shughni stem translator is a separate FST which translates Shughni stems to Russian lemmas. It is supposed to be applied to the glossed side of transducers as shown in Code block 13. Shughni has some cases of compound words which consist of multiple stems, e.g. ‘*дапоз-зѣѣ*’ (=‘long-sleeved’) (Parker, 2023, p. 173). This was taken into account, so the translator FST translates any amount of stems in a single word. This FST is blind to morphology rules and will translate anything that contains a set of existing stems and tags, even if the combination is nonsensical. But after composition with a morphology model FST ungrammatical and nonsensical strings are left outside the of resulted FST’s paradigm.

The `lexd` source code pattern is shown in Code block 14. `RuLemmasBase` pattern contains lexicons grouped by POS. This ensures that homographs (wordforms with same graphical form but different meanings) with different POS will not share contexts. An example of a homograph in Shughni is ‘*ду*’, which can stand for a verb (*hit.V.PRS/IMP*), a demonstrative (*D2.M.SG*) or a conjunction (=‘*when.CONJ*’). An example of FST with and without POS-grouping is shown in Code block 15. A wordform ‘*ди*’, when fed into `analyze_pos_ignored.hfst`, where homographs are not taken into account, returns ungrammatical glosses like ‘*когда<v>*’(‘when<v>’) and ‘*этом<v>*’(‘this_one<v>’), which is an unwanted result. But when fed into `analyze_pos_fixed.hfst`, where translator has POS tags are glued to stems, the result contains only verbal Russian lemma ‘*бумь<v>*’(‘hit<v>’).

Translation lexicon generation

Translation lexicons (`RuLemmasAdj`, `RuLemmasConj`, etc.) are compiled using a Python script `scripts/ru_lemmas/form_lexd.py`. This FST module (`sgn2rulem.hfst`) does not share source code with the main morphology FST, so its lexicons are stored in a separate `translate/lexd/`

```

$ cat glossed.txt | hfst-lookup -q translate/sgh2rulem.hfst
зѣѣ<n>          pyкaв<n>          0.000000
дapоз<adj>>зѣѣ<n>    длинный<adj>>pyкaв<n>    0.000000
дapоз<adj>><dim>      длинный<adj>><dim>      0.000000
$ cat nonsense.txt | hfst-lookup -q translate/sgh2rulem.hfst
<adj><v>>зѣѣ<n><pl><sg>    <adj><v>>pyкaв<n><pl><sg>    0.000000
<dim><pl><sg><1sg><2sg>    <dim><pl><sg><1sg><2sg>    0.000000
$ hfst-compose translate/rulem2sgh.hfst sgh_gen_stem_word_cyr.hfst
-o sgh_gen_rulem_word_cyr.hfst
$ cat nonsense.txt | hfst-lookup -q sgh_gen_rulem_word_cyr.hfst
<adj><v>>зѣѣ<n><pl><sg>    <adj><v>>зѣѣ<n><pl><sg>+?    inf
<dim><pl><sg><1sg><2sg>    <dim><pl><sg><1sg><2sg>+?    inf

```

Code 13: An example of Shughni stem translator work. ‘+?’ at the end of second column’s string and ‘inf’ in the third column means that FST did not find a valid path for the current word.

‘зѣѣ’=a sleeve; ‘дapоз’=long

Note: output is edited: contextually insignificant lines are removed to keep Code section small.

directory. For convenience every POS lexicon is stored in a separate file under this directory. Compilation process is done by the same principle as for the main FST (see Code block 8).

The biggest challenge of this process was the extraction of Russian lemmas from the database dictionary (provided by Makarov et al. (2022)). Gloss lemmas are preferred to be concise. They usually consist of a single word like (e.g. *swim*) or maximum of 2-3 words (e.g. *swim_deep*) if lemma language can not describe semantics in a single word, and it is important to mention this semantic nuance. Dictionary entries in the database are parsed from real dictionaries, that were written by hand and therefore do not have a strict format. This fact makes it challenging to consistently extract perfect lemmas, some examples are presented in Table 4. In addition, as shown by ‘чи’ entry, sometimes a real lemma (in this case ‘кто’) can be hidden in the middle of a text.

Fortunately, most of the lemmas can be extracted automatically without a problem. Russian lemma size minimal statistics are provided in Table 5. Most of the Russian lemmas consist of a couple of words: 76.5% consist of 1 word; 88.9% of ≤ 2 words; 94.2% of ≤ 3 words; 96.6% of ≤ 4 words. Nevertheless, there are still many long lemmas similar to the ones shown in the Tables 4 and 5. Considering that a single Shughni stem often has multiple Russian lemmas, the probability of getting at least one long lemma rises.

4.8 Testing

Whoa so much testing

| Word | Dictionary entry | Extracted lemma (by script) |
|--------|---|---|
| цирӯвд | горевать, тосковать, печалиться; скорбеть | горевать |
| ҳйвдов | сбивать палкой орехи, бить по орешнику палкой | сбивать_палкой_орехи |
| цӯн | подчинительный союз: сколько ни, как ни | сколько_ни |
| чи | косвенная форма к прямой форме вопросительного местоимения ЧАЙ кто (см.), употребляемая в различных косвенных позициях | косвенная_форма_к_прямой_форме_ _вопросительного_местоимения_чй_ _кто |

Table 4: Examples of extracted Russian lemmas from dictionary entries.

| Type | Mean | Median | Max |
|---------------------|--------|--------|---|
| Length (characters) | 10.737 | 8 | 105 ‘указывает_на_косвенное_дополнение_со_ _значением_адресата_действия_при_ряде_ _глаголов_и_глагольных_сочетаний’ |
| Length (words) | 1.487 | 1 | 16 ‘косв_форма_мн_ч_указ_мест_дальн_ст_к_ _прямой_форме_w_них_и_т_п’ |

Table 5: Examples of extracted Russian lemmas from dictionary entries.

```

PATTERNS
(RuLemmasBase|RuLemmasTags) +

PATTERN RuLemmasBase
RuLemmasAdj [<adj>]
RuLemmasConj [<conj>]
RuLemmasN [<n>]
RuLemmasNum [<num>]
RuLemmasPrep [<prep>]
RuLemmasPro [<pro>]| [<pers>]| [<dem>]
RuLemmasV [<v>]

```

Code 14: A fragment of translator FST's `lexd` source code. The fragment contains only pattern rules. **TODO: update if set of POS changes!**

```

$ hfst-compose translate/rulem2sgh_no_pos.hfst sgh_base_stem.hfst |
  hfst-invert -o analyze_pos_ignored.hfst
$ echo "мā>ди" | hfst-lookup -q analyze_pos_ignored.hfst
ди      битъ<v><imp>          0.000000
ди      когда<v><imp>         0.000000
ди      этот<v><imp>          0.000000
$ hfst-compose translate/rulem2sgh.hfst sgh_base_stem.hfst |
  hfst-invert -o analyze_pos_fixed.hfst
$ echo "мā>ди" | hfst-lookup -q analyze_pos_fixed.hfst
ди      битъ<v><imp>          0.000000

```

Code 15: A comparison of FST analyzers. The first one called `analyze_pos_ignored.hfst` treats all homographs independently of their part of speech, which results in incorrect output. The second one (`analyze_pos_fixed.hfst`) returns only verb Russian lemmas.

Note: output is edited: contextually insignificant lines are removed to keep Code section small.

4.9 Metrics

Whoa so much metrics

5 Results

6 Conclusion

References

- Beesley, e. R., & Karttunen, L. (2002). *Finite-state morphology: Xerox tools and techniques*. CSLI, Stanford.
- Budilova, Z. A. (2023). *Создание морфологического парсера для чамалинского языка в системе lexd u twol [Morphological parser of Chamalal in lexd and twol]*, NRU HSE. <https://www.hse.ru/edu/vkr/837214661>
- Buntuyakova, V. A. (2023). *Создание морфологического парсера андийского языка в системе lexd u twol [Morphological parser of Andi in lexd and twol]*, NRU HSE. <https://www.hse.ru/edu/vkr/837214826>
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In J. Burstein, C. Doran, & T. Solorio (Eds.), *Proceedings of the 2019 conference of the north American chapter of the association for computational linguistics: Human language technologies, volume 1 (long and short papers)* (pp. 4171–4186). Association for Computational Linguistics. <https://doi.org/10.18653/v1/N19-1423>
- Edelman, D. I. (2009). Сравнительная грамматика восточноиранских языков [Comparative Grammar of Eastern Iranian Languages].
- Edelman, D. I., & Yusufbekov, S. (1999). Шугнанский язык [Shughni language]. In *Языки мира: Иранские языки. III: Восточноиранские языки [Languages of the world: Iranian languages. III. Eastern Iranian languages]*. <https://iling-ran.ru/web/ru/publications/langworld/volumes/7>
- Edelman, D. I., & Dodykhudoeva, L. R. (2009). Shughni. In G. Windfuhr (Ed.), *The iranian languages* (pp. 787–824). London & New York: Routledge.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9, 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- Karamshoev, D. (1963). *Баджувский диалект шугнанского языка [Badzhuvskij dialect of the Shughni language]*. изд-во АН Тадж. ССР. <https://books.google.ru/books?id=8q1GXwAACAAJ>
- Karamshoev, D. (1988–1999). *Шугнанско-русский словарь [Shughni-Russian dictionary]*. Izd-vo Akademii nauk SSSR.
- Kartina, E. (2025). *Afmigoo/bachelor_diploma*. https://github.com/afmigoo/bachelor_diploma
- Karttunen, L. (1993). Finite-state lexicon compiler. Technical Report ISTL-NLTT-1993-04-02, Xerox Palo Alto Research Center, Palo Alto, CA.
- Karttunen, L., Koskenniemi, K., & Kaplan, R. (1987). A compiler for two-level phonological rules. *tools for morphological analysis*.
- Kondratyuk, D., & Straka, M. (2019). 75 languages, 1 model: Parsing universal dependencies universally.
- Koskenniemi, K. (1983). Two-level Morphology: A General Computational Model for Word-Form Recognition and Production.
- Lindén, K., Silfverberg, M., & Pirinen, T. (2009). HFST Tools for Morphology – An Efficient Open-Source Package for Construction of Morphological Analyzers. In C. Mahlow & M. Piotrowski (Eds.), *State of the Art in Computational Morphology* (pp. 28–47, Vol. 41). Springer Berlin Heidelberg. https://doi.org/10.1007/978-3-642-04131-0_3
- Makarov, Y., Melenchenko, M., & Novokshanov, D. (2022). Digital Resources for the Shughni Language. *Proceedings of The Workshop on Resources and Technologies for Indigenous, Endangered and Lesser-Resourced Languages in Eurasia within the 13th Language Resources and Evaluation Conference*, 61–64. <https://aclanthology.org/2022.euralli-1.9>
- Melenchenko, M. G. (2021). *Автоматический морфологический анализ шугнанского языка [Automatic full morphology analysis for Shughni]*, NRU HSE.

-
- Osorgin, I. G. (2024). *Создание морфологического парсера для шугнанского языка в системе lexd и twol [Creating a morphological parser for the Shughni language using lexd and twol systems]*, NRU HSE.
- Parker, C. (2023). *A grammar of the Shughni language* [Doctoral dissertation, Department of Linguistics, McGill University].
- Plungian, V. (2022). The study of shughni: The past and the future. *RSUH/RGGU Bulletin: "Literary Theory. Linguistics. Cultural Studies"*, Series. 2022;(5):11-22.
- Salemann, K. (1895). Шугнанский словарь Д.Л. Иванова [Shughni dictionary by D.L. Ivanov]. In *Восточные заметки [Eastern notes]*.
- Sarveswaran, K., Dias, G., & Butt, M. (2021). ThamizhiMorph: A morphological parser for the Tamil language. *Machine Translation* 35, 37–70. <https://link.springer.com/article/10.1007/s10590-021-09261-5>
- Swanson, D., & Howell, N. (2021). Lexd: A Finite-State Lexicon Compiler for Non-Suffixational Morphologies.
- Tumanovich. (1906). *Краткая грамматика и словарь шугнанского наречия [Brief grammar and dictionary of Shughni]*.
- Turing, A. M. (1937). On computable numbers, with an application to the entscheidungsproblem. <https://doi.org/10.1112/plms/s2-42.1.230>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 30). Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf
- Washington, J., Ipasov, M., & Tyers, F. (2012). A finite-state morphological transducer for Kyrgyz. In N. Calzolari, K. Choukri, T. Declerck, M. U. Doğan, B. Maegaard, J. Mariani, A. Moreno, J. Odijk, & S. Piperidis (Eds.), *Proceedings of the eighth international conference on language resources and evaluation (LREC'12)* (pp. 934–940). European Language Resources Association (ELRA). <https://aclanthology.org/L12-1642/>
- Zarubin, I. (1960). *Шугнанские тексты и словарь [Shughni texts and dictionary]*. Izd-vo Akademii nauk SSSR.