

Полевые данные и компьютерные инструменты

Г. Мороз

Международная лаборатория языковой конвергенции, НИУ ВШЭ,
Москва

Иркутский государственный университет

14 – 16 ноября 2022 г.

«Цифра» в социально-гуманитарных исследованиях: метод, поле,
реальность?

План презентации

Малые языки в большой лингвистике

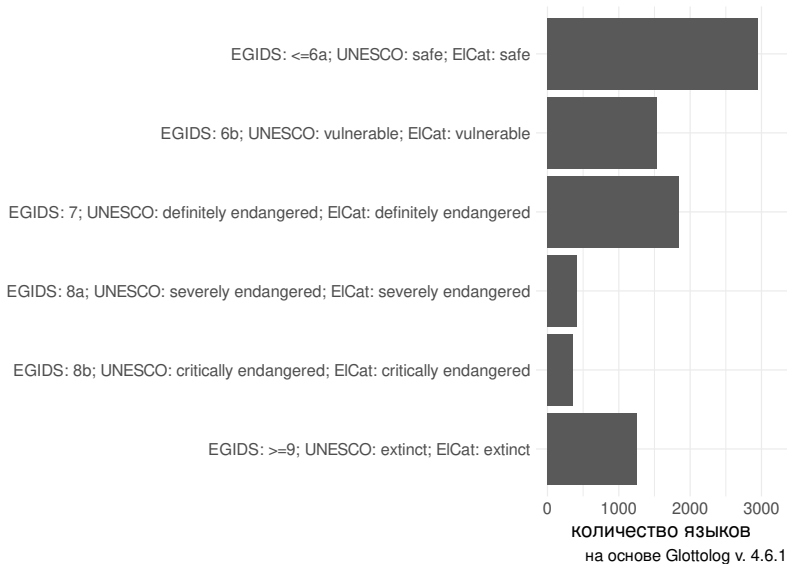
Устные корпуса

Морфологический анализаторы

Какие бывают языки с точки зрения компьютерных инструментов?

- огромные (шкала EGIDS ([Lewis and Simons 2010](#)) = 1)
 - доступно очень много текстовых данных, сформировалась литература, есть здоровый интернет
 - доступны даже исторические данные, скажем, на целый век
- средние (шкала EGIDS от 2 до 5)
 - доступны грамматические описания, двуязычные словари
 - литературная традиция часто ограничена одним веком, и обычно немногочисленна
- малые (шкала EGIDS больше 5)
 - острая нехватка материалов

Какие бывают языки с точки зрения компьютерных инструментов?



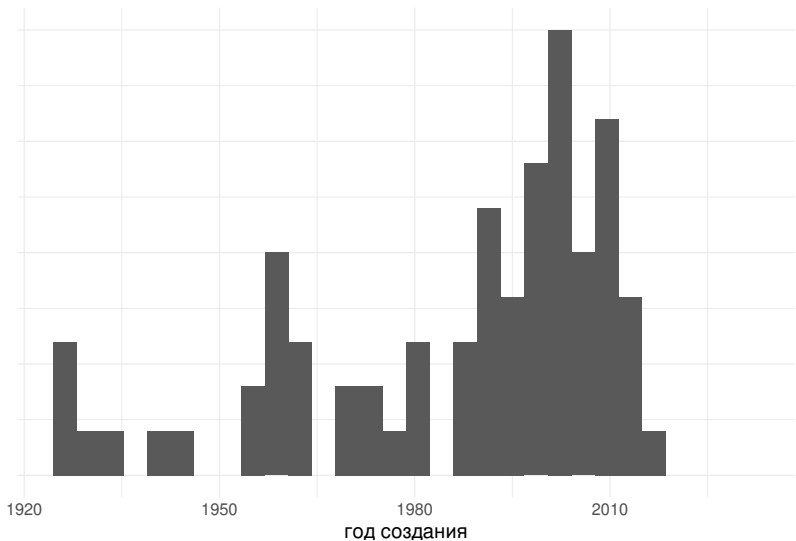
Как это влияет на инструменты компьютерной лингвистики для малых языков?

- нет данных или их очень мало, так что все что строится на нейросетях не работает
 - морфологический парсер?
 - синтаксический парсер?
 - распознавание/синтез речи?

Как это влияет на инструменты компьютерной лингвистики для малых языков?

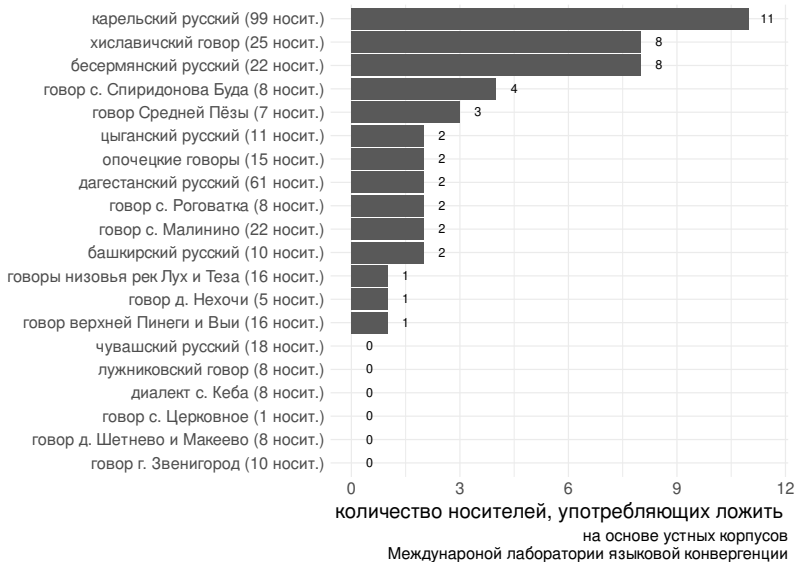
- нет данных или их очень мало, так что все что строится на нейросетях не работает
 - морфологический парсер?
 - синтаксический парсер?
 - распознавание/синтез речи?
- как следствие, не всегда можно вслепую копировать методы разработанные для больших языков

Лексема *ложить* в НКРЯ



77 уникальных авторов отфильтрованных из 141 примера

Лексема *ложить* в устных корусах



Как это влияет на инструменты компьютерной лингвистики для малых языков?

- нет данных или их очень мало, так что все что строится на нейросетях не работает
 - морфологический парсер?
 - синтаксический парсер?
 - распознавание/синтез речи?
- как следствие, не всегда можно вслепую копировать методы разработанные для больших языков

В результате, для малых языков необходимы особые инструменты и методы, которые бы позволяли облегчать или делать возможным лингвистический анализ и инструменты для языкового комьюнити (спеллчекеры, предективный набор и т. п.).

Инструменты для полевых лингвистов

- устные корпуса
 - без морфологической разметки
 - с морфологической разметкой
- морфологические анализаторы (трансдюсеры lexd и twol)
- синтаксические парсеры (проект Universal Dependencies)

План презентации

Малые языки в большой лингвистике

Устные корпуса

Морфологический анализаторы

SpoCo (von Waldenfels and Woźniak 2016)

- устьянские говоры (959 782 сл.-у.);
- хиславичский говор (260 793 сл.-у.);
- говоры низовья рек Лух и Теза (146 350 сл.-у.);
- говор с. Малинино (138 943 сл.-у.);
- говор с. Роговатка (100 047 сл.-у.);
- говор д. Нехочи (88 965 сл.-у.);
- говор Средней Пёзы (79 566 сл.-у.);
- говор верхней Пинегы и Выи (70 803 сл.-у.);
- говор с. Спиридонова Буда (70 565 сл.-у.);
- опочецкие говоры (68 741 словоупотребление);
- лужниковский говор (68 666 сл.-у.);
- говор г. Звенигород (68 324 сл.-у.);
- говор д. Шетнево и Макеево (58 003 сл.-у.);
- говор с. Кеба (54 535 сл.-у.);
- говор с. Церковное (19 960 сл.-у.).

SpoCo (von Waldenfels and Woźniak 2016)

- дагестанский русский (227 885 сл.-у.);
- бесермянский русский (97 216 сл.-у.);
- башкирский русский (93 127 сл.-у.);
- карельский русский (74 014 сл.-у.);
- чувашский русский (46 307 сл.-у.);
- цыганский русский (41 767 сл.-у.);
- якутский русский (15 139 сл.-у.).

SpoCo (von Waldenfels and Woźniak 2016)

Malinino Corpus

Corpus ▾

Statistics ▾

About ▾

Help

[rus en]

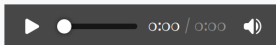
Query: [lemma='рука'%c]

Number of results: 35

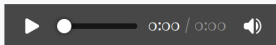


Recording

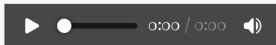
Transcription



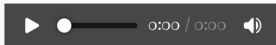
Так , дай **руку** пожалуйста .



А теперь всё , я говорю , до того износилось , вот как это
работаешь **рукой** . . . < смеется >



Тоже пил , повалился и **руки** отморозил .



Без **рук** остался , всё равно пил .



201008e_mgd

[*Слушать в одном файле \(mp3\)*](#)

First transcription: 0:0:0 Last transcription: 0:18:10

Interviewer: <>ⁱ ↓

mgd: <Рисовала>... Русалка.ⁱ ↓

Interviewer: Ага.ⁱ ↓ А масленицу делали у вас?ⁱ ↓

mgd: А?ⁱ ↓

Interviewer: А на масленицу не было?ⁱ ↓

mgd: О, а как же! И масленица бывает.ⁱ ↓ Тогда, бывало, каталися с горы.ⁱ ↓

Interviewer: С горы?ⁱ ↓

Tsakorpus (Arkhangelskiy 2019)

- Абазинский
- Адыгейский
- Башкирский
- Кабардинский
- Хакасский
- Луговой Марийский

Tsakorpus (Arkhangelskiy 2019)

Устный корпус башкирского языка дер. Рахметово и с. Баимово

RU | EN | ?

Возврат к поиску



Результат поиска: найдены 28202 словоформы, 2983 предложения примерно в 96 документах.

110717_gaj_Dialog_o_muzhe 2011

[S2_unknown] **jeläk** **jəjəyəðmə** ?
jeläk n jəjaw v
jeläk jəj-a-yəð-mə
berry collect-IPFV-2PL-Q
unmarked ipfv, 2pl, q
ягода собирать

[S2_unknown] Ягоды собираете?
ягода собирать



140708_ksg_Rozhdenie_detey 2014

но только пятница, пятница или понедельник.

но только пятница, пятница или понедельник.
но только пятница пятница или понедельник



160817_rmm_Vengrija 2016

beð beð pron beð we мы	qurqə qurqəw v be.afraid-IPFV ipfv be.afraid бояться	inek ine v ine-k be.PST-1PL pst, 1pl be.PST быть	samolyotta, läkin samolyot n samolyot-ta plane-LOC loc plane самолет	šul läkin conj šul pron šul that that но	tiklem tiklem post tiklem up.to up.to до	həjbät həjbät adj həjbät good good хороший	barəp barəw v bar-əp go-CV cv go идти	jettek jetew v jet-te-k be.enough-PST-1PL pst, 1pl be.enough хватать
---	--	---	---	---	--	--	--	---

Мы боялись в самолете, но очень хорошо доехали.
мы бояться в самолет но очень хорошо доезжать



Все наши корпуса доступны здесь

<http://lingconlab.ru/>

План презентации

Малые языки в большой лингвистике

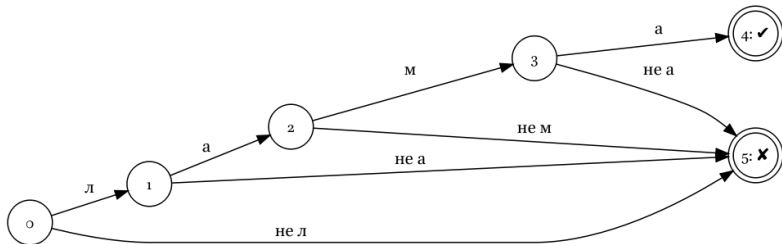
Устные корпуса

Морфологический анализаторы

Что такое трансдьюсер?

Трансдьюсер (конечный автомат с выходом) — это вид конечного автомата с двумя лентами памяти.

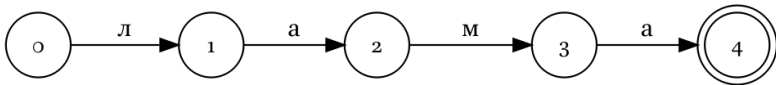
- Вот пример конечного автомата с одной лентой памяти. Он проверяет является ли поданное на вход слово словом *лама*:



Что такое трансдьюсер?

Трансдьюсер (конечный автомат с выходом) — это вид конечного автомата с двумя лентами памяти.

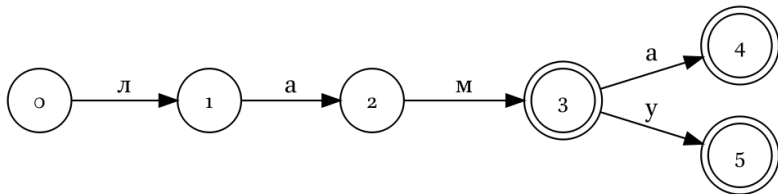
- Вот пример конечного автомата с одной лентой памяти. Он проверяет является ли поданное на вход слово словом *лама*.
- Обычно эти ветви “не X” не пишут:



Что такое трансдюсер?

Трансдюсер (конечный автомат с выходом) — это вид конечного автомата с двумя лентами памяти.

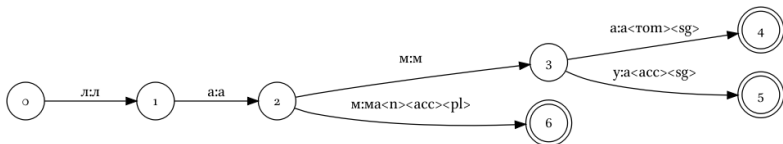
- Вот пример конечного автомата с одной лентой памяти. Он проверяет является ли поданное на вход слово словом *лама*.
- Обычно эти ветви “не X” не пишут.
- Можно закодировать больше слов (*лама, ламу, лам*):



Что такое трансдьюсер?

У трансдьюсеров две ленты памяти, что можно воспринимать как переписывание:

- лама переходит в лама<n><nom><sg>
- ламу переходит в лама<n><acc><sg>
- лам переходит в лама<n><acc><pl>
- все остальное – ошибка



Зачем использовать трансдюсер?

- они обратимы, так что анализ (ламу \rightarrow лама<acc><sg>) и генерация (лама<acc><sg> \rightarrow ламу) могут быть сделаны одним и тем же трансдюсером

Зачем использовать трансдюсер?

- они обратимы, так что анализ (ламу \rightarrow лама<acc><sg>) и генерация (лама<acc><sg> \rightarrow ламу) могут быть сделаны одним и тем же трансдюсером
- их можно оптимизировать для быстрого поиска

Зачем использовать трансдюсер?

- они обратимы, так что анализ (ламу → лама<acc><sg>) и генерация (лама<acc><sg> → ламу) могут быть сделаны одним и тем же трансдюсером
- их можно оптимизировать для быстрого поиска
- их можно соединять с другими трансдюсерами (например, транслитерация или даже перевод)

Зачем использовать трансдьюсер?

- они обратимы, так что анализ (ламу → лама<acc><sg>) и генерация (лама<acc><sg> → ламу) могут быть сделаны одним и тем же трансдьюсером
- их можно оптимизировать для быстрого поиска
- их можно соединять с другими трансдьюсерами (например, транслитерация или даже перевод)
- формализм трансдьюсеров позволяет описывать языковой материал в приближенном к лингвистическому описанию виде

Как использовать трансдюсер?

- можно почитать ([Beesley and Karttunen 2003](#); [Karttunen and Beesley 1992](#))
- `lexd` — компилятор для морфологии ([Swanson and Howell 2021](#))
- `twol` — компилятор для морфонологии

lexd пример (зиловский андийский)

PATTERNS

Numerals NumearalMarker

LEXICON Numerals

иЧшду	# пять; five
ойлИи	# шесть; six
гьокьу	# семь; seven
бейкьи	# восемь; eight
гьочIo	# девять; nine

LEXICON NumearalMarker

<num>:гy

lexd пример (зиловский андийский)

PATTERNS

Numerals NumearalMarker

LEXICON Numerals

и"шду	# пять; five
ойлИи	# шесть; six
гьокьу	# семь; seven
бейкьи	# восемь; eight
гьочIo	# девять; nine

LEXICON NumearalMarker

<num>:гу

и"шдугу:и"шдү<num>

ойлИигу:ойлИи<num>

гьокьугу:гьокьү<num>

бейкьигу:бейкьи<num>

гьочIoгу:гьочIo<num>

Как разработать морфологический трансдьюсер?

- опишите морфологию и морфонологию используя доступные ресурсы

Как разработать морфологический трансдьюсер?

- опишите морфологию и морфонологию используя доступные ресурсы
- составьте словарь с словоизменительной аннотацией

Как разработать морфологический трансдьюсер?

- опишите морфологию и морфонологию используя доступные ресурсы
- составьте словарь с словоизменительной аннотацией
- опционально можно сделать список тестовых форм, которые трансдьюсер должен разбирать

Как разработать морфологический трансдьюсер?

- опишите морфологию и морфонологию используя доступные ресурсы
- составьте словарь с словоизменительной аннотацией
- опционально можно сделать список тестовых форм, которые трансдьюсер должен разбирать
- чIe<NUM><num><obl.m><epent.m><an.sg><aff> чIeryшyбo (Zilo Andi)
- проверьте ваш трансдьюсер на аннотированном (или даже неаннотированном) корпусе

Спасибо за внимание!

agricolamz@gmail.com

References

- T. Arkhangelskiy. 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, 2019.
- K. R. Beesley and L. Karttunen. *Finite-state morphology: Xerox tools and techniques*. CSLI, Stanford, 2003.
- L. Karttunen and K. R. Beesley. *Two-level rule compiler*. Xerox Corporation, Palo Alto Research Center, 1992.
- M. P. Lewis and G. F. Simons. Assessing endangerment: expanding fishman's gids. 2010.
- D. Swanson and N. Howell. Lexd: a Finite-State lexicon compiler for non-suffixational morphologies. In M. Härmäläinen, N. Partanen, and K. Alnajjar, editors, *Multilingual facilitation*. 2021.

Ruprecht von Waldenfels and Michał Woźniak. Spoco-a simple and adaptable web interface for dialect corpora. *Journal for language technology and computational linguistics*, 31:155–170, 2016.