

Корпус вепсского языка: проблемы и перспективы

Veps corpus



¹Institute of Language, Literature and History

²Institute of Applied Mathematical Research
of the Karelian Research Centre of the RAS

¹Nina Zaitseva

²Andrew.Krizhanovsky

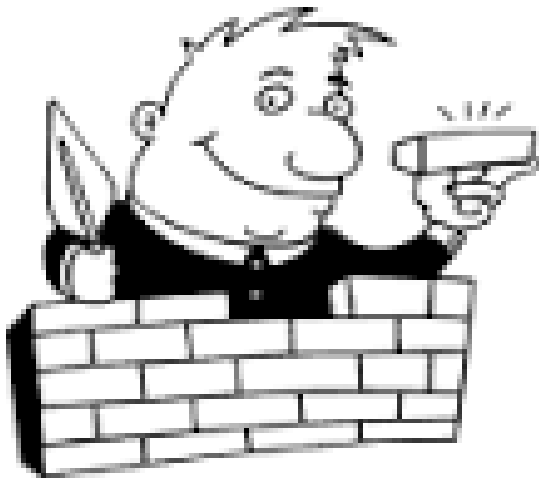


gmail.com



Goal

To create spoken
and written corpus



First task

Develop software
(corpus as a web-site).



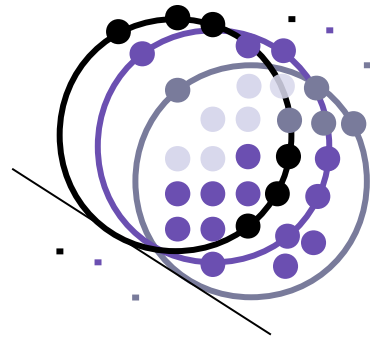
Second task

Fill data:
add texts, add words.





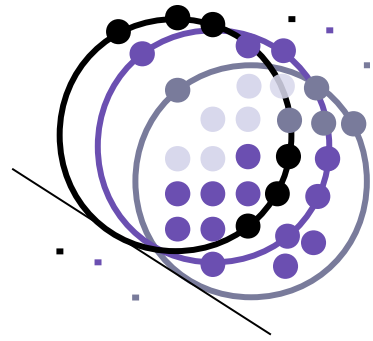
What is Corpus Linguistics?



Corpus Linguistics is the study of language/linguistic phenomena through the analysis of data obtained from a corpus.



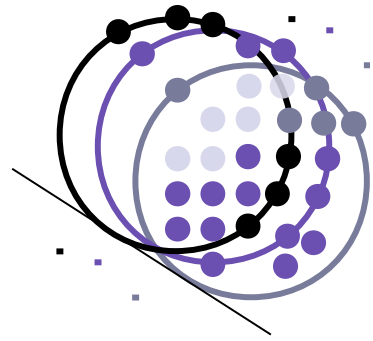
Why to use a *corpus*?



- Intuition alone is not enough
 - Is “*starting*” always replaceable by “*beginning*”?
 - Is it only “*time*” that is “*immemorial*”?
 - “*think of*” vs. “*think about*”
- Native speaker intuition is unreliable
 - provides no information on frequency of occurrence
 - “**head**” => body part - Is this the most used sense?
- Create corpus once, use for different linguistics tasks



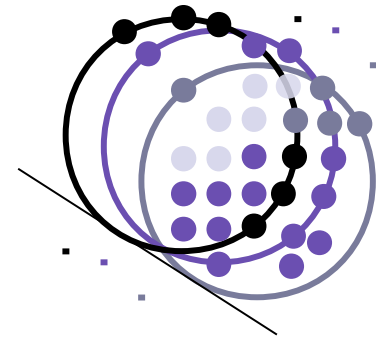
Тексты, входящие в состав корпусов



- отобраны исходя из определенных принципов,
- специально подготовлены и размечены,
 - машиночитаемый формат + разметка
- с помощью специальных программ в них можно искать необходимые фрагменты текста по заданным параметрам.



Types of corpora



- spoken vs. written
- monolingual vs. bi/multilingual
- parallel vs. comparable corpora (translation corpora)
- general language purpose vs. specialised language purpose
- diachronic vs. synchronic
- plain text vs. annotated (tagged) text

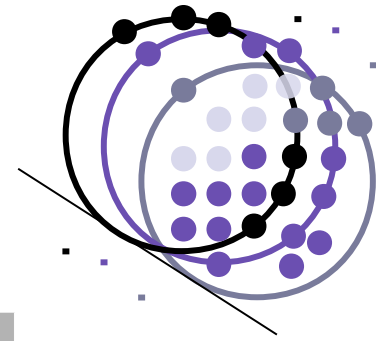


Web corpora

Корпусы(а) в Интернете



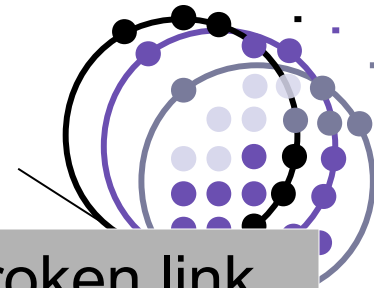
Russian corpora



| | |
|-------------------------------------------------------------------------------------------------------------|---------------|
| Национальный корпус русского языка http://ruscorpora.ru | >500 млн слов |
| Открытый корпус русского языка http://opencorpora.org | > 1 млн слов |
| | |



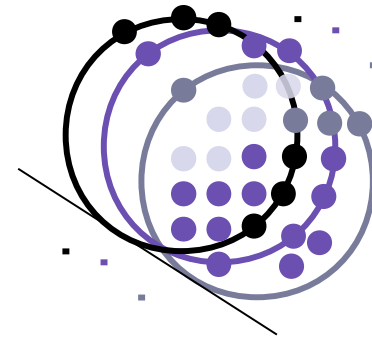
Finno-Ugric corpora



| | |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------|
| Языковой банк Финляндии http://www.csc.fi/tutkimus/alat/kielitiede | Broken link |
| Справочный корпус эстонского языка http://www.keeletehnoloogia.ee/projects-1/the-reference-corpus-of-the-estonian-language/comprehensive-corpus-of-estonian?set_language=en | Broken link |
| Фонетический корпус спонтанной эстонской речи http://www.murre.ut.ee/phonetic-corpus/ | 20 hours, 20 speakers |
| Венгерский национальный корпус http://mnsz.nytud.hu/index_eng.html | > 187 million words |



Veps corpus



Veps dictionary + corpus are available at:

<http://vepsian.krc.karelia.ru/>

Корпус вепсского языка

[О проекте](#) [Поиск в словаре](#) [Поиск по текстам](#) [Тексты](#) [Подкорпус причитаний](#) [Ссылки](#)



Поиск в словаре

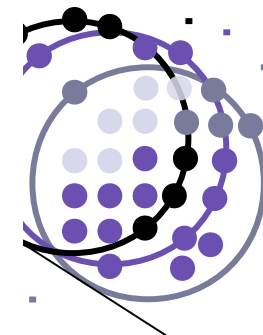
Statistics 1

Статистика

Новые материалы: леммы, словоформы, переводы, **омонимы**

Число лемм в словаре: 2750

Число словоформ в словаре: 8408



Wild card 2: %

везде в леммах в словоформах

Найти

ä ö ü č š ž '

Umlaut 3

Ваш запрос: __lud%

лемма: **kuluda**, POS: **V** (глагол)

| № | словоформы | Грамматические признаки |
|----|------------|-------------------------|
| 1. | kulub | |
| 2. | kulu | |
| 3. | kului | |

Перевод
русский: слышаться
английский: to be heard

лемма: **kuludas**, POS: **V** (глагол)

| № | словоформы | Грамматические признаки |
|----|--------------|-------------------------|
| 1. | kuluškanzihe | |

Перевод
русский: послышаться
английский: to be heard

лемма: **polut**, POS: **N** (существительное)

| № | словоформы | Грамматические признаки |
|---|------------|-------------------------|
| | | |

Поиск по текстам

Статистика

Новые материалы: тексты, источники

Число текстов в корпусе: 1097

Число библиографических источников: 871

Пример: vasan, minä или ленточка, или %o_o_ни%

В текстах

ä ö ü č š ž '

- приближённый поиск лемм и словоформ в текстах, то есть поиск подстрок (по умолчанию выполняется точный поиск)

*Ваш запрос: **vasan***

Леммы в словаре: vas

Перевод лемм: vas (квас || bread juice | kvass | quass)

Словоформы в словаре: vasan, vas-se, vas

По запросу найдено текстов: 2

1. [Kut tegiba vasan](#)

Язык и народ: тексты и комментарии, (2002), с. 67-68

Карпов Никанор Леонтьевич, г.р. 1897, место записи: Ладва (Ladv), Подпорожский р-н, Ленинградская обл.

1. [Kut tegiba vasan](#)

Язык и народ: тексты и комментарии, (2002), с. 67-68

Карпов Никанор Леонтьевич, г.р. 1897, место записи: Ладва (Ladv), Подпорожский р-н, Ленинградская обл.

1. Potom liiban paned sin'n'a da kipekkoizuu valad viluu, da vuu sepid paned dei **vas** tegese čoma.
2. Ende tegiba **vasan**, ka rugihen ligotadas da id'atadas da päče pandas.
3. Tegžihe imen, nu da sit' kuivatas da jouhtas da pašttas liibežiks, nu i **vasan** pandas.
4. Otad, **vasan** valad, lukhiineižid da ret'kašt da babuš.
5. Uuged ezmäks paned mišto ii jokseiž liib, riigun kohtha, riig mi jokseb **vas-se**.

2. [Van'ka-vor](#)

Вепские народные сказки, (1996), с. 122-124; ф/архив ИЯЛИ КарНЦ РАН: № 2624/2, НА КарНЦ, кол.83, ед.хр.146

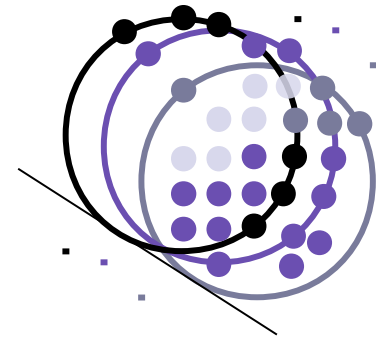
Микшина Марфа Захарьевна, г.р. 1910, место записи: Ладва (Ladv), Подпорожский р-н, Ленинградская обл., г. записи: 1980, записали: Онегина Нина Федоровна

По запросу найдено текстов: 2,

из этих текстов найдено цитат, содержащих искомые слова из словаря: 5



Veps corpus



User

- Search in dictionary
- Search in corpus

Editor (admin)

- Edit dictionary
- Edit texts in corpus

Все тексты

Выберите подкорпус

- диалектные тексты (199)
- библейские тексты (переводные) (431)
- младописьменный подкорпус (89)
- подкорпус вепских причитаний (47)
- подкорпус вепских сказок (55)

свойство текста

- северновепский диалект (27)
- средневепский диалект (214)
- южновепский диалект (57)
- восточные говоры (32)
- западные говоры (31)
- свадебные причитания (29)
- похоронные и поминальные причитания (18)
- публицистические тексты (51)
- художественные тексты (21)
- тексты для детей (17)

выбрать

Подкорпус вепских причитаний

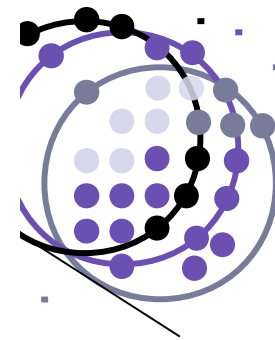
Список текстов только для одного языка: вепский

Северновепский диалект

Найдено текстов: 5.

1. [Nevestan voik svad'ban päiväl](#)

северновепский диалект, свадебные причитания

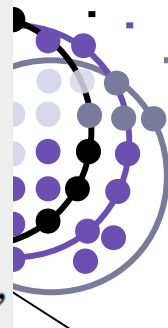


Nevestan voik svad'ban päiväl

подкорпус вепских причитаний

северновепский диалект, свадебные причитания

*Информант: Ермолаева Анастасия, г.р. 1899, урожд. Шелтозеро (Šoutarv),
Прионежский р-н, Республика Карелия, место записи: Шелтозеро (Šoutarv),
Прионежский р-н, Республика Карелия, г. записи: 1962, записали: Лонин
Рюрик Петрович*



Nevestan voik svad'ban päiväl

Плач невесты в день свадьбы

Tänambeižel päivaižel om milei-d'o
svad'baine.

В сегодняшний денечек у меня уже
свадьба.

Kogozihe kaik minun rodn'aine.

Собралась вся моя родня.

Milei vaise üht ii täudu bat'uškod-se
rodnijad.

Не хватает мне одного батюшки
родного.

Armhad tö minun kaik läheližed,

Любимые вы мои родственники,

armhad tö minun susedaižed,

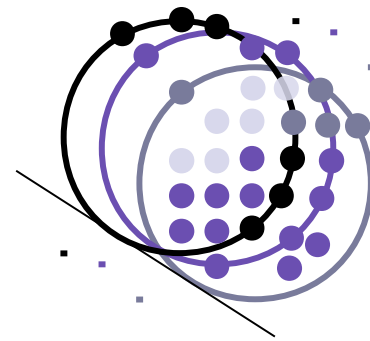
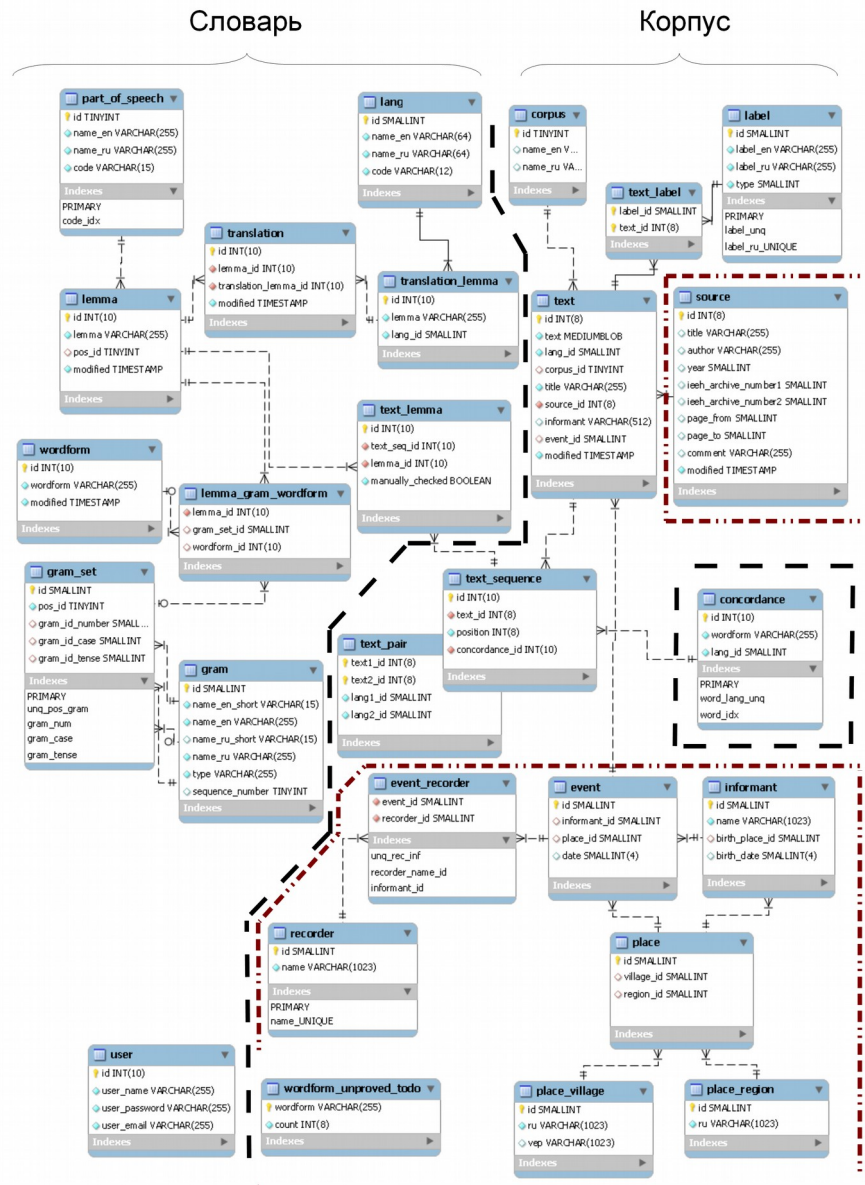
любимые вы мои соседушки,

avaikat tö mili dorogaine,

откройте вы мне дороженьку,



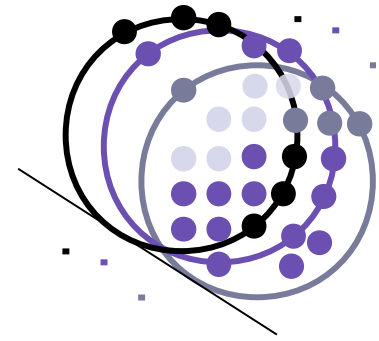
dictorpus



Паспорт, источники в корпусе текстов



Future work



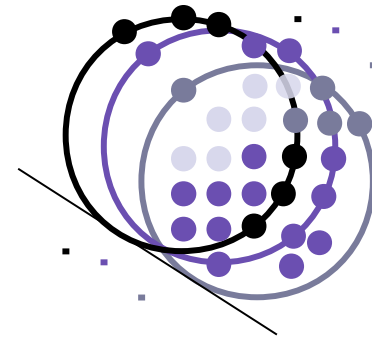
Open corpus of Veps and Karelian languages

(VerKar) (Открытый корпус вепсского и карельского языков)

- Open software license,
- Open (Creative-commons cc-by) data



Thank you



Veps corpus:

<http://vepsian.krc.karelia.ru/>

The paper “Корпус вепсского языка”:

<http://scipeople.com/publication/121149/>

Корпус вепсского языка

[О проекте](#) [Поиск в словаре](#) [Поиск по текстам](#) [Тексты](#) [Подкорпус причитаний](#) [Ссылки](#)

