

# ВепКар: руководство для пользователей

Крижановская Н. Б., Новак И. П., Пеллинен Н. А.

Издание 1. Декабрь 2022

## Содержание

<b>Введение</b>	<b>4</b>
<b>1. Принципы работы в современных корпусах текстов</b>	<b>4</b>
1.1. Корпусная лингвистика и виды корпусов	4
1.1.1 Основные понятия корпусной лингвистики	4
1.1.2. Разметка корпусов	6
Понятие разметки	6
Лингвистическая разметка	6
Экстралингвистическая разметка	6
1.1.3. Типология корпусов	7
1.1.4. Обзор существующих корпусов различных типов	7
Зарубежные корпуса	8
Корпусы русского языка	8
Корпусы уральских языков	9
1.2. Описание корпуса ВепКар	11
1.2.1. История создания	11
1.2.2. Архитектура ВепКар	11
1.2.3. Страница с текстом в корпусе ВепКар	13
1.2.4. Словарная статья в корпусе ВепКар	15
1.2.5. Открытость данных	17
<b>2. Методические рекомендации по использованию интернет-ресурса ВепКар</b>	<b>17</b>
2.1. Применение интернет-ресурса ВепКар	17
2.2. Рекомендации по использованию интернет-ресурса ВепКар	18
2.2.1. Общие вопросы	18
Шаг 1. Как найти сайт Открытого корпуса вепского и карельского языков (ВепКар) в Интернете?	19

Шаг 2. Как переключить интерфейс сайта ВепКар на русский язык?	21
Шаг 3. Обязательно ли нужно регистрироваться на сайте ВепКар?	21
Шаг 4. Как ввести особые графические символы карельского и вепсского языков, если их нет на клавиатуре?	21
Стоит особо отметить, что в процессе поиска символ ü приравнен к символу у, w - к u / y / ü, а знак палатализации ' не учитывается.	21
Шаг 5. Как использовать простые шаблоны в текстовых полях?	21
Шаг 6. Как использовать специализированные шаблоны в текстовых полях?	21
2.2.2. Поиск в корпусе ВепКар	23
Шаг 7. Как осуществлять поиск текстов в корпусе ВепКар?	23
Шаг 8. Как найти текст на вепсском / карельском языке?	23
Шаг 9. Как осуществлять расширенный поиск по текстам ВепКар?	24
Шаг 10. Как найти озвученный текст на вепсском / карельском языке?	25
Шаг 11. Как найти вепские сказки, записанные до 1950 года?	25
Шаг 12. Как осуществить поиск примеров употребления слова в корпусе?	26
Шаг 13. Как найти примеры употребления в корпусе людиковских имен в форме аппроксиматива?	26
Шаг 14. Как найти примеры употребления в корпусе ливвиковских глаголов в форме перфекта и плюсквамперфекта кондиционала?	27
Шаг 16. Как найти примеры употребления конкретной словоформы в вепсских текстах 30х годов?	30
Шаг 17. Как найти частотные словари корпуса?	31
Шаг 18. Как найти самые частотные слова в текстах?	31
Шаг 19. Как найти «Речевой корпус»?	32
Шаг 20. Как найти и прослушать образцы речи интересующего населённого пункта?	33
Шаг 21. Как ссылаться на материалы корпуса ВепКар?	34
2.2.3. Поиск в словаре ВепКар	34
Шаг 22. Как осуществлять поиск в словаре ВепКар?	34
Шаг 23. Как найти вепское / карельское слово по начальной форме?	34
Шаг 24. Как указать язык/наречие при поиске леммы?	35
Шаг 25. Как осуществлять расширенный поиск в словаре ВепКар?	35
Шаг 26. Как найти вепское (карельское) слово по русскому толкованию?	36
Шаг 27. Как найти вепское (карельское) слово по словоформе?	36

Шаг 28. Как найти все глаголы новописьменного севернокарельского языка, начинающиеся с буквы 'm', у которых начальная форма (инфинитив) оканчивается на 'uo' или 'uö', а форма 1 лица ед.ч. настоящего времени оканчивается на 'un' или 'un'?	37
Шаг 29. Как найти обратный словарь?	37
<b>Заключение</b>	<b>38</b>
<b>Тезаурус</b>	<b>39</b>
<b>Список использованной литературы</b>	<b>39</b>

# Введение

Сохранение уникальной языковой культуры имеет серьезное научное и социальное значение. В России и во всем мире уменьшается число носителей коренных языков и стоит задача сохранения языков малых народов. Из 250 языков России около 100 являются миноритарными языками<sup>1</sup>. Стоит отметить, что даже языки, не относящиеся к меньшинствам, например, якутский (саха), могут находиться под угрозой исчезновения. Для большинства языков России, кроме русского, цифровых ресурсов (электронных словарей, лингвистических корпусов и текстов онлайн) либо нет, либо их относительно мало [Klyachko, 2019].

Для сохранения языкового богатства и последующего изучения языков создаются лингвистические корпусы текстов.

Электронные языковые корпусы в настоящее время являются одним из самых важных инструментов сохранения языков малых народов. Они представляют собой базу для решения миллионов задач прикладной лингвистики. Научное направление лингвистики, изучающее построение, анализ и использование языковых корпусов, называется *корпусной лингвистикой*.

Разработка открытых ресурсов, позволяющих хранить, систематизировать, осуществлять поиск информации по заданным параметрам, анализировать и интегрировать языковой материал в необходимый формат (печатный, электронный), представляется перспективным направлением в соответствии с требованиями современной науки.

## 1. Принципы работы в современных корпусах текстов

### 1.1. Корпусная лингвистика и виды корпусов

#### 1.1.1 Основные понятия корпусной лингвистики

Корпусная лингвистика занимается разработкой общих принципов построения и использованием лингвистических корпусов. Существует много определений того, что составляет лингвистический корпус, одно из них гласит, что корпус представляет собой «совокупность текстов или частей текстов, по которым можно провести общий лингвистический анализ» [Meuser, 2002, xi].

Наличие корпуса для языка, объём корпуса и различные его характеристики крайне важны, поскольку «любое лингвистическое исследование в той или иной мере опирается на анализ языкового материала, языковых данных» [Баранов, 2001].

Стоит также отметить, что благодаря корпусам многократно повысились не только эффективность и скорость обработки языковых данных, но и достоверность результатов, поскольку на корпусе данных легче проверить гипотезу, предположения, выводы [Горина, 2014, с. 4].

Для анализа языкового явления и проведения лингвистического исследования необходим языковой материал в электронной форме, то есть *лингвистический корпус*,

---

<sup>1</sup> Миноритарный язык - это язык национального (этнического) меньшинства [ССТ: 129].

представляющий собой коллекцию текстов, специально отобранных, размеченных по различным лингвистическим параметрам и обеспеченных системой поиска, которую называют *корпусным менеджером*.

### **Основные характеристики корпусов**

Исходя из определений можно представить минимальные требования к корпусу текстов, которые выражаются в следующих параметрах [Копотев, 2014]:

- *репрезентативность* — свойство корпуса, заключающееся в статистически достоверном представлении языка или его части и достигаемое за счет необходимого объема и жанрового многообразия текстов;
- *сбалансированность* — параметр, определяющий насколько равномерно представлены тексты разных типов;
- *объем корпуса* — информация об общем объеме корпуса и о количестве извлеченных из текста примеров;
- *электронная форма* хранения обеспечивает быстрый поиск и извлечение материала;
- *разметка* (аннотация) — введенная вручную или автоматически лингвистическая или метатекстовая информация обо всех выбранных единицах текста (тексте, предложении, словосочетании, словоформе, морфеме).

### **Прагматическая ориентированность**

Практика разработки и применения электронных корпусов текстов показала, что невозможно создать универсальный корпус, обеспечивающий решение всех задач. Задачи и цели любого исследования определяют тип корпуса, правила отбора текстов, а также способ и степень их обработки. Корпусы всегда создаются для решения определенной задачи или круга задач. Это определяет как наполнение корпуса текстами (например русская драма XIX в., тексты языка охотников), так и разметку корпуса [Захаров, Богданова, 2020].

Практика показывает также, что корпусная лингвистика оперирует как минимум тремя разными типами корпусов текстов:

1. *универсальными*, то есть отражающими все многообразие речевой деятельности;
2. *специфичными*, или представляющими бытование некоторого языкового или культурного явления в общественной речевой практике, например, корпус пословиц или корпус политических метафор в газетной речи;
3. *специальными*, то есть создаваемыми для решения специальной задачи, например, для обучения, для задач социолингвистики, для отладки систем машинного перевода.

## 1.1.2. Разметка корпусов

### Понятие разметки

Разметка (аннотирование) корпуса заключается в приписывании текстам и их компонентам специальных атрибутов (тегов) [Захаров, Богданова, 2020]. Теги бывают двух типов:

- *лингвистические*, описывающие лексические, грамматические и прочие характеристики элементов текста;
- *экстралингвистические* или метатекстовые (сведения об авторе или исполнителе и о тексте: название, год и место издания, жанр, тематика).

### Лингвистическая разметка

Особое значение имеет лингвистическая разметка, т.е. любая аннотация, основанная на лингвистических характеристиках текста. Данные лингвистической разметки можно добавлять к текстовым элементам разных уровней. С точки зрения технологии, разметка может быть автоматической, ручной или автоматической с ручной правкой [Захаров, Богданова, 2020].

Среди лингвистических типов разметки можно выделить следующие типы:

- *Морфологическая разметка* отражает признаки части речи и грамматических категорий, свойственных данной части речи. Код привязан к отдельному слову или группе слов. Это основной тип разметки, поскольку он используется для синтаксической и семантической разметки.
- *Синтаксическая разметка* описывает синтаксические связи между лексическими единицами и/или различные синтаксические конструкции (придаточное предложение, именное сказуемое и т. п.). Синтаксический код может быть закреплен за предложением или за синтаксическим отношением.
- *Семантическая разметка* обозначает семантические категории, к которым относится данное слово или словосочетание, и более узкие подкатегории, определяющие его значение. Семантический код закрепляется за словом или словосочетанием.
- *Анафорическая* разметка фиксирует референтные связи, например, местоименные.
- *Просодическая* разметка отражает ударение и интонацию.
- *Дискурсная* разметка служит для обозначения пауз, повторов, оговорок и т. п.

### Экстралингвистическая разметка

Дополнительная информация о текстах корпуса называется внешней или экстралингвистической разметкой (метаразметкой). Экстралингвистическая разметка может описывать как общие свойства текста, так и дополнительные технические данные:

- сведения об авторе;
- информацию о событии записи текста (характерно для диалектных текстов): год, место записи, собиратель (человек, записавший текст), информант (имя, год и место рождения, род занятий, национальность, пол);
- данные о публикации: автор, название, год, источник, страницы;

- типологические данные: художественный текст (например, рассказ, роман, басня) или нехудожественный (например, публицистический, научный, учебный, бытовой);
- тематическая (жанровая) классификация: историческая проза, детектив, драматургия и другие;
- технические данные: даты обработки, исполнители, кодировка, лицензия, даты загрузки;
- и другие свойства.

### 1.1.3. Типология корпусов

При разделении корпусов на классы рассматривают два подхода [Захаров, 2020]:

1. Относится ли корпус к языку в целом (возможно, к языку заданного периода) или корпус содержит тексты определённого стиля, жанра, отдельного писателя, социальной или возрастной группы и т.д.
2. Деление корпусов по типу разметки. В основном корпуса содержат только морфологическую разметку. Отдельно выделяют корпуса с синтаксической разметкой (treebanks).

Классификация корпусов по различным подходам представлена на Рис. 1.

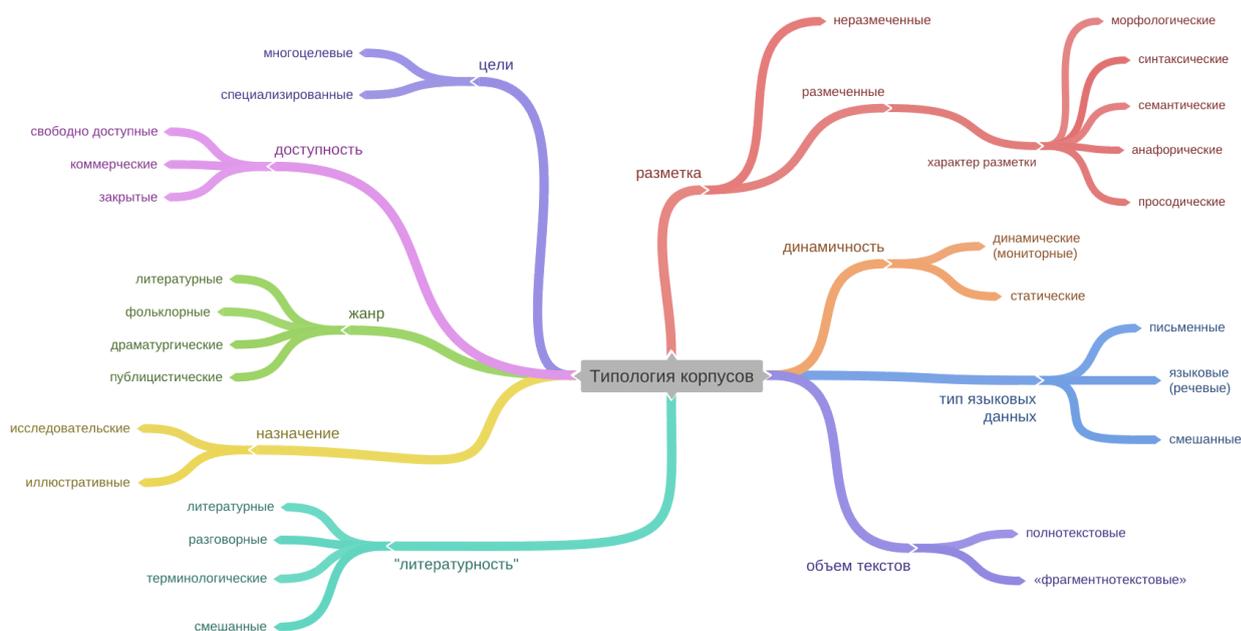


Рис. 1. Типы корпусов по классификации В.П. Захарова

### 1.1.4. Обзор существующих корпусов различных типов

В настоящее время для большинства основных языков мира существуют национальные общезыковые корпуса. Число лингвистических корпусов измеряется тысячами и постоянно растет. Два крупнейших специализированных каталога CLARIN<sup>2</sup> и ELRA<sup>3</sup> содержат информацию о более чем 3000 корпусах. Список лингвистических ресурсов,

<sup>2</sup> Common Language Resources and Technology Infrastructure, <https://www.clarin.eu/resource-families>

<sup>3</sup> European Language Resources Association, <http://www.elra.info/en/>

включающий и корпуса, можно найти на сайте Ассоциации компьютерной лингвистики<sup>4</sup>. Корпусы используются как для лингвистических исследований, так и в прикладных целях.

### **Зарубежные корпуса**

*Британский национальный корпус*<sup>5</sup> (кратко BNC) является большим (100 млн слов) корпусом британского варианта английского языка. Это *синхронический* корпус, он содержит актуальные для одного временного периода (конец XX века) примеры. Корпус разрабатывали два университета: Оксфордский и Ланкастерский в 1991-1994 годах.

Корпус содержит художественную и документальную прозу, тексты газет, периодических научных изданий и журналов, выпускаемых для различных возрастов, популярную научную фантастику, опубликованные и неопубликованные письма, школьные и университетские сочинения и др. 90% корпуса BNC составляют письменные тексты.

Устная речь в виде аудиозаписей с транскрипцией составляет 10% от объёма корпуса BNC. Авторы BNC наняли 124 добровольцев, мужчин и женщин разных возрастов, принадлежащих разным социальным группам, живущих в разных частях Великобритании, чтобы они записывали на плёнку свои обычные разговоры с друзьями, соседями, в семье, на работе. Была записана речь в самых разных контекстах — от разговоров на формальных деловых или правительственных встречах до радиошоу и телефонных разговоров. Было записано более 700 часов. К аудиозаписям были добавлены транскрипции (тексты с расшифровкой), тексты были аннотированы. Лаборатория фонетики Оксфордского университета создала звуковой подкорпус, для этого была выполнена оцифровка записей, каталогизация и анализ массива аудиокассет<sup>6</sup>, бывших недоступными в Звуковом архиве Британской библиотеки [Albert, 2012].

Одним из наиболее известных корпусов общего типа является *Чешский национальный корпус*<sup>7</sup>. Это синхронический морфологически размеченный корпус, представляющий современный чешский язык. Корпус был создан в конце 90-х годов в Карловом университете в Праге.

### **Корпусы русского языка**

Первый русскоязычный корпус был создан в 1980-е годы в Университете Упсалы (Швеция). В *Упсальский корпус* отбирались специальные тексты с 1985 по 1989 годы и художественные тексты с 1960 по 1988 годы с целью представить современный литературный язык, поэтому в массиве нет образцов разговорной речи. Упсальский корпус русского языка состоит из 600 текстов, его объём составляет 1 млн словоупотреблений. Корпус имеет морфологическую разметку.

Первый общедоступный, представительный и размеченный корпус русского языка, с которым могли бы работать лингвисты, стал доступен в 2004 году [Плунгян, Сичинава, 2004]. *Национальный корпус русского языка (НКРЯ)* содержит тексты на русском языке, корпус является диахроническим и содержит тексты, созданные на протяжении нескольких веков, включая современные тексты. Авторы НКРЯ использовали понятие

<sup>4</sup> Список ресурсов с разбивкой по языкам, [https://aclweb.org/aclwiki/List\\_of\\_resources\\_by\\_language](https://aclweb.org/aclwiki/List_of_resources_by_language)

<sup>5</sup> British National Corpus, BNC, <http://sara.natcorp.ox.ac.uk/> или <https://www.english-corpora.org/bnc/>

<sup>6</sup> См. примеры аудиозаписей BNC <http://www.phon.ox.ac.uk/SpokenBNCcontents>

<sup>7</sup> Český národní korpus, <http://korpus.cz/>

«подкорпуса» для разделения текстов на следующие большие группы: основной, поэтический, диалектный, параллельные тексты, газетный, церковнославянский, исторический, синтаксический, акцентологический, мультимедийный и обучающий. Как все современные лингвистические корпуса, он доступен онлайн<sup>8</sup>. Отметим такие важные составляющие НКРЯ, как богатые инструменты поиска и статистические сведения о корпусе, делящиеся на статическую и динамическую статистику [Гришина, Плунгян, 2003, с. 335]. Сложности создания большого корпуса русского языка при одновременном сосуществовании разных классификаций языковых уровней, ведущие к неизбежным трудностям при автоматическом аннотировании на основе этих классификаций, описаны в работе «Национальный корпус русского языка» [Копотев, Янда, 2006].

Ряд особенностей интерфейса НКРЯ и поисковых возможностей описан в работе «Национальный корпус русского языка» [Копотев, Янда, 2006, с. 150].

Table 1. **Примеры корпусов русского языка** (объем словоупотребления в млн.)

Корпусы	Адрес	Объем
Упсальский и Тюбингенский корпуса	<a href="http://www.sfb441.uni-tuebingen.de/b1/rus/korpora.html">http://www.sfb441.uni-tuebingen.de/b1/rus/korpora.html</a>	1
Машинный фонд русского языка	<a href="http://cfrl.ruslang.ru/">http://cfrl.ruslang.ru/</a>	100
Национальный корпус русского языка (НКРЯ)	<a href="https://ruscorpora.ru/new/index.html">https://ruscorpora.ru/new/index.html</a>	1 000
Хельсинкский аннотированный корпус русских текстов (ХАНКО)	<a href="http://h248.it.helsinki.fi/hanco/index.html">http://h248.it.helsinki.fi/hanco/index.html</a>	0,1
Открытый корпус русского языка	<a href="http://opencorpora.org/">http://opencorpora.org/</a>	1,5
Генеральный Интернет-корпус Русского Языка (ГИКРЯ)	<a href="http://www.webcorpora.ru/">http://www.webcorpora.ru/</a>	20 000
Корпус русского литературного языка	<a href="https://narusco.ru/">https://narusco.ru/</a>	1
Корпус русских учебных текстов (КРУТ)	<a href="http://web-corpora.net/learner_corpus">http://web-corpora.net/learner_corpus</a>	3,1
Русский учебный корпус	<a href="http://web-corpora.net/RLC">http://web-corpora.net/RLC</a>	2

### Корпусы уральских языков

В области уралоистики известны, прежде всего, электронные ресурсы наиболее крупных финно-угорских языков, такие как Языковой банк Финляндии, Венгерский национальный корпус, Текстовый корпус Института эстонского языка. В России ведется работа по созданию и наполнению корпусов мордовских, марийского, коми, удмуртского, мансийского языков. На сайте проекта “Малые языки Сибири” ([Siberian Lang | Малые](#)

<sup>8</sup> См. <https://ruscorpora.ru>

[языки Сибири: наше культурное наследие \(msu.ru\)](http://msu.ru)) представлены также аннотированные тексты с видео на селькупском самодийском языке.

Большой теоретический и практический интерес представляет собой корпус уральских языков Поволжья, разработанный Т. Архангельским. Для построения веб-корпуса пяти уральских языков (коми-зырянского, лугового марийского, мокшанского, удмуртского и эрзянского) обрабатываются тексты сети Интернет, корпус содержит только автоматическую разметку. В него включены современная художественная литература, научные статьи, переводы Библии, статьи Википедии, официальные тексты и публичные записи в социальных сетях [Arkhangelskiy, 2020, с. 58-61].

Table 2. **Примеры корпусов языков России** (объем словоупотребления в млн)

Корпусы	Адрес	Объем
Башкирский поэтический корпус	<a href="http://web-corpora.net/bashcorpus/search/index.php?interface_language=ru">http://web-corpora.net/bashcorpus/search/index.php?interface_language=ru</a>	1,8
Калмыцкий корпус	<a href="http://web-corpora.net/KalmykCorpus/search/index.php?interface_language=ru">http://web-corpora.net/KalmykCorpus/search/index.php?interface_language=ru</a>	0,8
Бурятский корпус	<a href="http://web-corpora.net/BuryatCorpus/search/index.php?interface_language=ru">http://web-corpora.net/BuryatCorpus/search/index.php?interface_language=ru</a>	2,2
Татарский национальный корпус «Туган тел»	<a href="http://web-corpora.net/TatarCorpus/search/index.php?interface_language=ru">http://web-corpora.net/TatarCorpus/search/index.php?interface_language=ru</a>	26

Table 3. **Корпусы уральских языков мира** (объем словоупотребления в млн.)

Корпусы	Адрес	Объем
Языковой банк Финляндии	<a href="https://sanat.csc.fi">https://sanat.csc.fi</a>	
Венгерский национальный корпус	<a href="http://mnsz.nytud.hu/index_hun.html">http://mnsz.nytud.hu/index_hun.html</a>	188
Текстовый корпус Института эстонского языка	<a href="https://portaal.eki.ee/corpus">https://portaal.eki.ee/corpus</a>	10,4
Корпус коми-зырянского языка	<a href="http://komi-zyrian.web-corpora.net/">http://komi-zyrian.web-corpora.net/</a>	1,76
Корпус лугового марийского языка	<a href="http://meadow-mari.web-corpora.net/">http://meadow-mari.web-corpora.net/</a>	5,53
Корпус мокшанского языка	<a href="http://moksha.web-corpora.net/">http://moksha.web-corpora.net/</a>	1,74
Корпус эрзянского языка	<a href="http://erzya.web-corpora.net/">http://erzya.web-corpora.net/</a>	2,3
Корпус удмуртского языка	<a href="http://udmurt.web-corpora.net/">http://udmurt.web-corpora.net/</a>	9,57

Национальный корпус удмуртского языка	<a href="http://udmcorpus.udman.ru/home">http://udmcorpus.udman.ru/home</a>	
Аннотированный корпус мансийских текстов	<a href="http://mansi.pro/corpus/">http://mansi.pro/corpus/</a>	0,05

## 1.2. Описание корпуса ВепКар

Открытый корпус вепского и карельского языков как результат многолетней междисциплинарной работы лингвистов и программистов Карельского научного центра РАН является уникальной источниковой базой для новых исследований.

### 1.2.1. История создания

Электронный ресурс ведет свою историю с 2009 года, когда под руководством Зайцевой Н. Г. был создан «Корпус вепского языка»<sup>9</sup> [Зайцева, 2012]. Стоит задача возрождения вепского и карельского языка и культуры. Для этого нужны большие электронные языковые ресурсы, которые позволят исследователям создавать и развивать, например, правила орфографии, вводить в научный и общественный оборот значительное количество материалов, востребованных авторами учебников по языку и истории родного края, учителями, писателями. В связи с этим, в целях сохранения, развития и популяризации не только вепского, но и карельского языка в 2016 году сотрудники Института языка, литературы и истории и Института прикладных математических исследований КарНЦ РАН приступили к созданию многоязычного корпуса. В 2016 году на базе «Вепского корпуса» была создан новый интернет-ресурс «Открытый корпус вепского и карельского языков» (ВепКар)<sup>10</sup>. В корпус помимо вепского вошли три карельских подкорпуса: собственно карельский, ливвиковский и людиковский. Объединённая лингвистическая платформа получила название «Открытый корпус вепского и карельского языков» (ВепКар). В 2022 году ВепКар «заговорил»: был разработан речевой модуль, позволяющий слышать тексты или леммы, озвученные носителями языка. Редакторами ВепКара осуществляется постепенное наполнение звукового корпуса аудиотекстами с их расшифровкой.

### 1.2.2. Архитектура ВепКар

Открытый корпус вепского и карельского языков — это многоязычный полнотекстовый лингвистический корпус, содержащий морфологическую, семантическую и метатекстовую разметку.

Корпус включает в себя хранящиеся в базе данных тексты и словари на вепском и карельском языках, часть из которых сопровождается аудиофайлами, а также компьютерную программу (корпусный менеджер), обеспечивающую поиск и обработку текстов. Корпусный менеджер написан на языке программирования PHP в системе разработки веб-сайтов Laravel. Данные хранятся в базе данных MySQL. Словари и тексты корпуса

<sup>9</sup> См. <http://vepsian.krc.karelia.ru/about/>

<sup>10</sup> См. <http://dictorpus.krc.karelia.ru>

вместе с поисковой системой доступны онлайн. Авторы проекта уделяют внимание популяризации корпуса ВепКар с помощью сайтов YouTube и Википедия.

Особенностью корпуса ВепКар является тесная взаимосвязь словарей и текстов. Многофункциональные словари вепского и карельского языков содержат толкование в русском и частично английском языках, перевод, диалектные пометы, семантические отношения (синонимы, антонимы и др.), примеры словоупотреблений со ссылкой на тексты, а также полные словоизменительные парадигмы. Все тексты автоматически размечаются и от слов в тексте идут отсылки на соответствующие значения в словарных статьях.

Разработчики добавляют в корпусный менеджер новые полезные функции, призванные облегчить работу редакторов. Например, были сформулированы и запрограммированы правила именного и глагольного словоизменения [Krizhanovskaya et al., 2022] для всех диалектов вепского языка и его младописьменного варианта, а также для ливвиковского, севернокарельского и тверского новописьменных вариантов карельского языка. Благодаря этому в системе ВепКар в полуавтоматическом режиме было сгенерировано 2.1 млн словоформ. Кроме семантической разметки, представленной в корпусе (2.1 млн связей между словами из текста и значениями лемм в словаре), была добавлена грамматическая разметка, позволившая автоматически установить 1.1 млн связей между словами из текста и грамматическими характеристиками словоформ из словаря.

Многоязычный корпус ВепКар делится на подкорпуса по языкам и наречиям, также есть стилевая и жанровая классификация текстов. В корпусе организована развитая система поиска с фильтрацией текстов по языковой, стилистической и диалектной принадлежности, по информанту, собирателю или автору, году записи или году публикации. Поиск лемм возможен по диалектам, частям речи, грамматическим признакам и даже по лексико-семантическим категориям.

Корпус ВепКар можно рассматривать как электронную библиотеку, предоставляющую пользователям доступ к полным текстам документов. Материалами для ВепКар могут служить тексты только с открытой лицензией<sup>11</sup>.

---

<sup>11</sup> См. страницу корпуса ВепКар с разрешениями авторов текстов: <http://dictorpus.krc.karelia.ru/ru/page/permission>



Рис. 2. Архитектура корпуса ВепКар. Связь словаря и корпуса текстов. Данные на ноябрь 2022 г.

Таким образом, интернет-ресурс ВепКар объединяет в себе три составляющие: размеченный корпус, многофункциональный словарь и библиотеку.

Подробнее об архитектуре (Рис. 2) и возможностях ВепКар читайте в работе [Бойко и др., 2021], в докладе на конференции [Крижановская, Крижановский, 2020] и смотрите видеозапись на YouTube<sup>12</sup>.

### 1.2.3. Страница с текстом в корпусе ВепКар

Страница с текстом (Рис. 3) в зависимости от информации о тексте (то есть от метаданных) содержит:

- информацию об авторе;
- заголовок;
- название подкорпуса / коллекции;
- языковую принадлежность;
- диалектную принадлежность;
- для диалектных текстов сведения о месте и времени записи, информанте, собирателе;
- сведения о публикации;
- информацию о хранении исходного материала в Научном архиве Карельского научного центра РАН или Фонограммархиве Института языка, литературы и истории КарНЦ;
- аудиофайл;
- размеченный текст;
- перевод текста (в основном на русский язык):

<sup>12</sup> См. видеозапись «Архитектура корпуса ВепКар и диалектные особенности» <https://www.youtube.com/watch?v=cUpqM97LXGs>

- автор перевода
- заголовок
- текст, выровненный по предложениям с исходным.

При наведении мышки на предложение исходного текста, жёлтым цветом выделяется соответствующее ему предложение-перевод. При наведении мышки на предложение-перевод, голубым цветом выделяется соответствующее ему предложение в исходном тексте.

Если кликнуть на размеченное слово в вепском/карельском тексте, откроется окно с леммой, значением и грамматическими свойствами. Каждое автоматически размеченное слово проверяется экспертом. Зеленый цвет означает, что слово проверено, синий и красный — не проверено. Если система нашла одно соответствие в словаре, то слово помечено синим, если несколько — красным. Чёрный цвет означает, что в словаре пока нет соответствия для слова в тексте.

### "L'ähtiimmö müä siäneh..."

подкорпус: диалектные тексты

информант(ы): Куршиева Анна Николаевна, 1910

Макарова Фекла Алексеевна, 1902

место записи: Самбатукса (Sammatus), Олонецкий район, Республика Карелия, г. записи: 1963

записали: Вяйзинен Т. И., Макаров Григорий Николаевич

источник: Г.Н. Макаров, В.Д. Рягоев, Образцы карельской речи. Говоры ливвиковского диалекта карельского языка, 1969, с. 61-63

ф/архив ИЯЛИ КарНЦ РАН: №221/3

### "L'ähtiimmö müä siäneh..."

Карельский: ливвиковское наречие  
Коткозерский

[Anni]: L'ähtiimmö müä siäneh jäl'les voinan  
lorpuw.

Da müä segoimmo sinne, vihmupäivü oli, müä  
emmo kodih voinnut puwttua.

Puwtuimmo sinne majakan alle, Leningradskoih  
mežah, dai müä üävüimmö, magaimmo üän.

N'eveskü se uinoi muate, virš  
a minä üän kaiken puhuin ha

Viižitoštu vuattu oli majakal luajittuw, niigoin  
[halgoloih] üän kaiken tuldu puhuin.

Sit huandesčura ku rodih, ruvettih siä hukat  
ulvomah, meččü oli ü'len suwri.

L'ähtiimmö müä (huandes tuli) poikki suas.

### «Отправились мы за рыжиками...»

Русский

[Анна]: Отправились мы за рыжиками после  
того, как война кончилась.

Мы там заблудились, дождливый день был,  
мы не смогли попасть домой.

Мы попали туда, к маяку ('под маяк'), где  
инградская межа, мы там и заночевали,  
провели.

тка-то уснула, корзиночку положила в  
овье, а я всю ночь поддерживала огонь  
(дула в огонь) в гнилом валежнике.

Пятнадцать лет прошло [с тех пор], как маяк  
тот в лесу был поставлен, дрова из этого  
маяка были, и я всю ночь поддерживала  
костёр.

Стало за полночь, начали там выть волки,

Рис. 3. Фрагмент ливвиковского диалектного текста "L'ähtiimmö müä siäneh..."<sup>13</sup>

<sup>13</sup> Полный текст см. <http://dictorpus.krc.karelia.ru/ru/corpus/text/1514>

## 1.2.4. Словарная статья в корпусе ВепКар

На сайте корпуса ВепКар организован словарь, по которому можно производить поиск (Рис. 33). Щелкнув на найденное слово, можно перейти к словарной статье. В словаре ВепКар словарные статьи называются “**леммами**” (лемма - начальная форма слова). Если слово из одного языка (наречия) употребляется в разных частях речи, то для каждой части речи имеется отдельная словарная статья.

После автоматической разметки каждому значению слова привязаны примеры из полнотекстового корпуса ВепКар. После ручной разметки (экспертной проверки) эти примеры помечены звёздочками с оценками (“лучший”, “хороший” и т.д.) (Рис. 4). Если у примера нет оценки, то этот пример возможно не проверен экспертом, а значит может относиться к другому значению или даже другой словарной статье.

### ahven

язык: карельский: ливвиковское наречие

часть речи: существительное



#### 1 значение

понятие: окунь

- русский: окунь
- английский: European perch, redfin perch, perch



перевод

вепский: ahven; ahn'

карельский: людиковское наречие: ahven

карельский: собственно карельское наречие: ahven

диалекты употребления: Ведлозерский, Видлицкий, Кондушский, Неккульский, Сямозерский, Тулмозерский

Примеры (17)

★ лучший    ★ отличный    ★ хороший    ★ плохой

- ★ **Ahvenet** ollah ei ylen suuret, vie poijat.  
Окуни ещё не крупные, ещё мальки. ("Minä sanon, kui müö provodiimmo kanikuluw")
- ★ Ahven kudou hätken, tedri kiimuiččou hätken.  
Окунь нерестится долго, тетерев токует долго. (Kalua pyvimmö ijän kaiken)
- ★ Vie sanotah: **ahven** ahaval kudou, a särgi siäl lämmäl.  
Еще говорят: окунь в холодную сухую погоду нерестится, а плотва в теплую. (Kalua pyvimmö ijän kaiken)
- ★ Sit nowzow, tämä joi, **ahven**: sit ahvenet kuvotetah kalakkahat järvet.  
Потом поднимается, так сказать, окунь, в рыбных озёрах во время нереста ловят окуня. (Minä olen rodinuh Čil'miele)
- ★ Eriastu sanondua on kalastajih näh, što vaiku pattii ei kergie **ahvenen** kudah da tedrin kiimah.  
Различные выражения есть о рыбаках, что только ленивый не успеает на окуневый нерест и тетеревиный ток. (Kalua pyvimmö ijän kaiken)

Рис. 4. Словарная статья для ливвиковского существительного *ahven* с тремя «лучшими», одним «отличным» и одним «хорошим», цитатами из корпуса ВепКар, проверенными лингвистами

## sanuo

язык: карельский: ливвиковское наречие

часть речи: глагол

фонетические варианты: [sanua](#) (Сямозерский, Видлицкий, Неккульский, Тулмозерский); [sania](#) (Кондушский)

### фразеологизмы

[sanuo silmiñ](#) - сказать в лицо (букв. сказать в глаза)

[sanuo tervehyöt](#) - передать привет

[sanuo valmistelemattah](#) - сказать экспромтом

[sanuo vastah](#) - перечить (букв. говорить наперекор)

### 1 значение

понятие: [сказать](#)

- русский: говорить, сказать

перевод

вепсский: [sanuda](#)

карельский: людиковское наречие:  
[sanuo](#); [sanuda](#); [sanoda](#)

карельский: собственно карельское  
наречие: [šanua](#); [sanuo](#); [šanuo](#)

диалекты употребления: Ведлозерский

### Примеры (3043)

★ лучший ★ отличный ★ хороший ★ плохой

1. ★ I kummoksie ei midä ole, hūö ei voija kohti [sanuo](#), što heil muudu ni midä ei pie, ku vai ottua valdu ruadajoil da krest'janoil käzis iäre.

*И неудивительно, они не могут прямо сказать нам, что им ничего другого не надо, кроме как отобрать власть из рук рабочих и крестьян. (Ruadorahvahale kaiken muan)*

2. ★ "Kai ruado rahvahale, ni midä ruadamattomile – vierahan ruavon süöjile" – vot kui [sanotah](#) bol'shevikat. "Все должно быть предоставлено трудящемуся народу, ничего не давать тем, кто не работает, кто сосёт чужую кровь" – вот как говорят большевики. (Ruadorahvahale kaiken muan)

3. ★ "Towta, – ma [sanon](#), – minul oli brihaččuiine kodij d'iänüh, ga tiiätgo kus on?"

*"Тетя, – говорю я, – остался у меня мальчик дома, не знаешь ли, где он?" ("Meil omii lapsii iellon")*

[еще примеры >>](#)

### 2 значение

- русский: рассказывать, рассказать

### Примеры (2959)

★ лучший ★ отличный ★ хороший ★ плохой

1. ★ Minä [sanon](#), kui müö provodiimmo kanikuluw.

*Я расскажу, как мы проводили каникулы. ("Minä sanon, kui müö provodiimmo kanikuluw")*

[еще примеры >>](#)

### 3 значение

- русский: называть, назвать

### Примеры (2962)

★ лучший ★ отличный ★ хороший ★ плохой

1. ★ Nu müö kävüimmo müömäh talolui müöte: sukija, piččie (nemme kruwživot – suomekse piččkse [sanotah](#)), rihmjä.

*Ну, мы ходили торговать по домам: чулки, кружева (по-фински кружева «питси» называют), нитки. (Torrun kelle Suomes)*

Рис. 5. Фрагмент словарной статьи ливвиковского глагола 'sanuo' с тремя значениями<sup>14</sup>

Словарная статья (Рис. 5) в зависимости от степени её проработки содержит:

- язык;
- грамматическую информацию
  - часть речи,
  - возвратность, безличность, переходность у глагола,
  - одушевленность у существительного,
  - степень сравнения у прилагательного и наречия,
  - тип наречия,
  - разряд местоимения и т.д.;
- фонетические варианты в других диалектах;
- фразеологизмы или устойчивые словосочетания с этим словом;
- значения с толкованиями на русском, английском, вепсском, наречиях карельского и финском языках;
- переводы значений на вышеперечисленные языки;
- семантические отношения с другими леммами в данном языке: синонимы, антонимы, гиперонимы, гипонимы, холонимы, миронимы и т.д.;
- диалекты, в которых употребляется данное значение;

<sup>14</sup> Полную статью см. <http://dictorpus.krc.karelia.ru/ru/dict/lemma/15190>

- примеры из корпуса с оценками эксперта;
- список словоформ (Рис. 6) со ссылками на тексты корпуса.

словоформы (150)

№	грамматические признаки	Новописьменный ливвиковский (149)	Коткозерский (1)
<b>Индикатив, презенс, положительные формы</b>			
1.	1 л., ед. ч.	<a href="#">sanon</a>	
2.	2 л., ед. ч.	<a href="#">sanot</a>	
3.	3 л., ед. ч.	<a href="#">sanou</a>	<a href="#">sanow</a>
4.	1 л., мн. ч.	<a href="#">sanommo</a>	
5.	2 л., мн. ч.	<a href="#">sanotto</a>	
6.	3 л., мн. ч.	<a href="#">sanotah</a>	
7.	ед. ч., коннегатив	<a href="#">sano</a>	
8.	мн. ч., коннегатив	<a href="#">sanota</a>	

Рис. 6. Фрагмент словарной статьи ливвиковского глагола 'saniu' с тремя значениями<sup>15</sup>

### 1.2.5. Открытость данных

ВепКар включает в себя только те тексты, на которые получено разрешение авторов на публикацию под открытой лицензией CC-BY 4.0. При такой политике не каждый текст можно включить в корпус, однако наличие разрешения даёт возможность публиковать тексты (Рис. 3) и свободно передавать третьим лицам для научных исследований и разработки коммерческих приложений.

## 2. Методические рекомендации по использованию интернет-ресурса ВепКар

### 2.1. Применение интернет-ресурса ВепКар

Лингвистические корпуса используются студентами и преподавателями разных дисциплин, связанных с анализом употреблений слов, с поиском типичных или необычных словоформ и оттенков значений.

С помощью ВепКар преподаватели могут

- подбирать примеры и формировать задания по грамматике, фразеологии, фольклору, литературе и разговорной практике по карельскому и вепсскому языкам;
- использовать текстовые материалы корпуса для проведения курса домашнего чтения;
- привлекать диалектные тексты в курсе диалектологии в целях выявления фонетических и грамматических особенностей карельской и вепсской диалектной речи;
- использовать материалы корпуса в процессе разработки новых учебных пособий;
- привлекать материалы параллельных словарей для обучения основам перевода.

<sup>15</sup> Полную статью см. <http://dictorpus.krc.karelia.ru/ru/dict/lemma/15190>

С помощью ВепКар студенты могут

- находить полные словоизменительные парадигмы изменяемых слов, т.е. использовать материалы словарного модуля корпуса в качестве электронного грамматико-орфографического словаря карельского и вепсского языков;
- использовать словарь корпуса в качестве переводного словаря карельского и вепсского языков;
- определять управление (глагольное управление, употребление послелогов и предлогов) с помощью контекстных примеров из корпуса текстов;
- выявлять фонетические и грамматические особенности карельской и вепсской диалектной речи на основе диалектных текстов;
- подбирать контекстные примеры в качестве источниковой базы для курсовых и дипломных работ.

С помощью ВепКар исследователи могут

- находить примеры, демонстрирующие различные фонетические явления, словоизменительные и словообразовательные особенности карельской и вепсской диалектной речи (расширенный поиск по заданным параметрам);
- получать выборки заданных словоформ из размеченных текстов с целью анализа грамматических категорий;
- анализировать и сравнивать тексты на новописьменных вариантах языка с целью выявления грамматических правил, нуждающихся в редактировании или доработке;
- проводить лингвостатистический анализ при помощи частотных словарей;
- отслеживать динамику изменения лексического состава анализируемых языков за вековой период;
- анализировать язык определенных жанров или конкретных авторов;
- использовать материалы корпуса в процессе разработки новых словарей.

С помощью ВепКар журналисты и писатели имеют возможность

- пользоваться большими объемами разнородных и разножанровых текстов для детального знакомства с творчеством коллег, разработки тем, образов и др.;
- уточнять и подбирать подходящие для текстов лексемы, словоформы и фразеологизмы, пользуясь линейкой словарей ВепКар (например, при подборе синонимов);
- привлекать данные корпуса для создания текстов и произведений на диалектах вепсского и карельского языков;
- использовать платформу ВепКар для популяризации своего творчества.

## **2.2. Рекомендации по использованию интернет-ресурса ВепКар**

### **2.2.1. Общие вопросы**

**Шаг 1. Как найти сайт Открытого корпуса вепского и карельского языков (ВепКар) в Интернете?**

Наберите адрес <http://dictorpus.krc.karelia.ru> и Вы попадете на стартовую страницу сайта корпуса ВепКар (Рис. 7).



Рис. 7. Стартовая страница сайта Открытого корпуса вепсского и карельского языков (VepKar)

## Шаг 2. Как переключить интерфейс сайта ВепКар на русский язык?

На главной странице сайта найдите иконку с флагом России и ссылку «Русский» и щелкните по ссылке мышкой (Рис. 8):

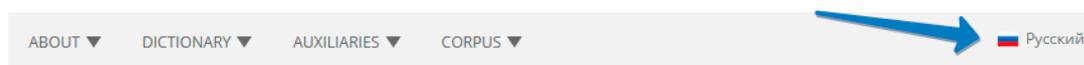


Рис. 8. Переключение на русский язык ВепКар.

## Шаг 3. Обязательно ли нужно регистрироваться на сайте ВепКар?

Регистрация не обязательна. Пока она необходима только редакторам ВепКар. Обычные пользователи на настоящий момент не получают никаких привилегий при регистрации. В будущем планируются дополнительные функции для зарегистрированных пользователей, например, сохранение типичных запросов, подобранных данных, организация «тяжёлых» запросов.

## Шаг 4. Как ввести особые графические символы карельского и вепского языков, если их нет на клавиатуре?

В тех полях, где могут понадобиться отличающиеся от английской раскладки буквы вепского или карельского алфавита (äöüčšž'), справа размещена иконка «специальные символы»: . Щелкните по этой иконке, иконка исчезнет, сверху откроется панель с буквами-кнопками (Рис. 9). Щелкните по нужному символу, символ вставится в поле на месте курсора.

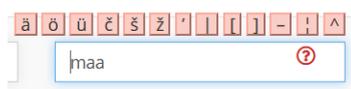


Рис. 9. Открытая панель «специальные символы»

Стоит особо отметить, что в процессе поиска символ  $\ddot{u}$  приравнен к символу  $u$ ,  $w$  - к  $u / y / \ddot{u}$ , а знак палатализации  $'$  не учитывается.

## Шаг 5. Как использовать простые шаблоны в текстовых полях?

В текстовых полях, если Вам нужен неточный поиск (по фрагменту слова), используйте процент  $\%$  для замены любого количества символов и подчеркивание  $\_$  для замены одного символа.

Например, по шаблону  $\%ta$  будут найдены все строки, оканчивающиеся на  $'ta'$ . А по шаблону  $\%änd\_$  будут найдены все строки, оканчивающиеся на  $'änd'$  и один любой символ. Например, будут найдены  $'eländy'$  или  $'vierändü'$ .

## Шаг 6. Как использовать специализированные шаблоны в текстовых полях?

В некоторых текстовых полях для построения более сложных запросов можно использовать специализированные шаблоны с регулярными выражениями (Regex и собственные). В table 4 приведены шаблоны и примеры их использования для поиска строк.

Table 4. Специализированные шаблоны для поиска текста

шаблон	описание	пример	пояснение
^	начало строки	^с	все слова, начинающиеся на букву <i>c</i> ( <a href="#">Поиск лемм по словоформам</a> )
\$	конец строки	in\$	все слова, оканчивающиеся на <i>in</i> (в словарной форме имен с.к. наречия: <a href="#">Поиск лемм по словоформам</a> )
.	любой один символ	h.\$	все глаголы на <i>ha/hä</i> в ливвиковском новописьменном варианте ( <a href="#">Поиск лемм по словоформам</a> )
[...]	любой один символ из перечисленных в квадратных скобках	läht[öe]m[äy]	все фонетические варианты понятия 'нетель' в словарной форме ( <a href="#">Поиск лемм по словоформам</a> )
[a-z]	любая латинская буква от <i>a</i> до <i>z</i>	[a-zA-Z]	любая латинская буква вне зависимости от регистра, например, нужно исключить символ палатализации
[^...]	любой один символ, <b>кроме</b> перечисленных в квадратных скобках	[^0-9]	любой символ, <b>кроме</b> цифры
?	идуший перед знаком вопроса символ может встретиться, а может и нет	^ikkuna?\$	<i>a</i> на конце может быть, а может не быть ( <a href="#">Поиск лемм по словоформам</a> )
*	ноль или более символов, идущих перед звездочкой	ab*c	<i>abc</i> ИЛИ <i>ac</i> ИЛИ <i>abbbc</i>
+	один или более символов, идущих перед плюсом	ab+c	<i>abc</i> ИЛИ <i>abbc</i> ИЛИ <i>abbbc</i>
{n}	<i>n</i> символов, стоящих перед скобками	b{3}	<i>bbb</i>
{m,n}	от <i>m</i> до <i>n</i> символов, стоящих перед скобками	a{1,3}	<i>a</i> ИЛИ <i>aa</i> ИЛИ <i>aaa</i>
{m,}	предыдущий символ может встретиться <i>m</i> и более раз	^pert{1,}i?\$	<i>pert</i> ИЛИ <i>perтт</i> ИЛИ <i>perti</i> ИЛИ <i>pertti</i> ( <a href="#">Поиск лемм по словоформам</a> )

(...)	круглые скобки задают группировку символов	(abc){1,3}	abc ИЛИ abcabc ИЛИ abcabcabc
p1 p2	p1 или p2	^rein ^rien ^ rejen ^regen	все варианты слабоступенной основы существительного regi ( <a href="#">Поиск лемм по словоформам</a> )
V	любая гласная буква	ttVV\$	все собственно карельские инфинитивы одноосновных глаголов, содержащих в основе чередование tt : t ( <a href="#">Поиск лемм по словоформам</a> )
C	любая согласная буква	^kaCrV\$	фонетические варианты с глухим и звонким согласным для kagra / kagru / kakra ( <a href="#">Поиск лемм по словоформам</a> )

## 2.2.2. Поиск в корпусе ВепКар

### Шаг 7. Как осуществлять поиск текстов в корпусе ВепКар?

- Выполните шаги 1 и 2.
- В верхнем меню найдите пункт «КОРПУС» и щелкните на него мышкой.
- В выпадающем меню щёлкните по ссылке «Тексты» (Рис. 10).

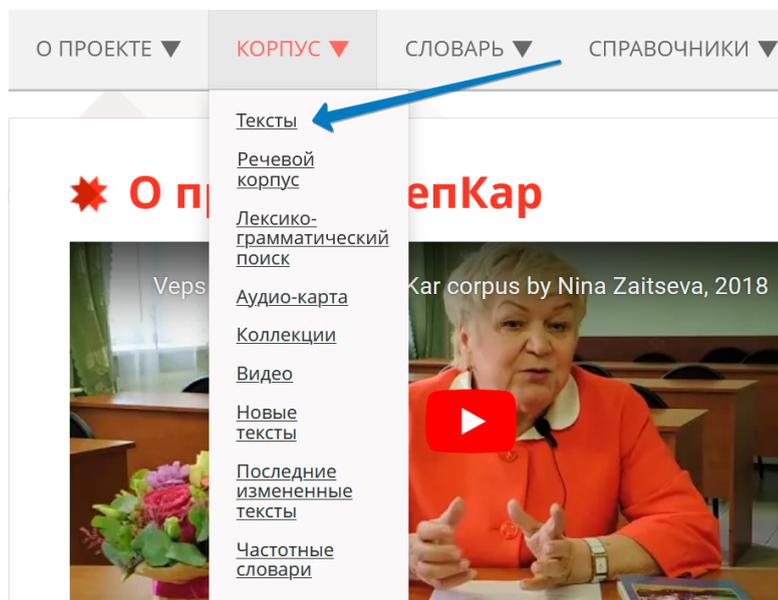


Рис. 10. Поиск в корпусе ВепКар

### Шаг 8. Как найти текст на вепском / карельском языке?

- Выполните шаг 7.
- В поле «Язык» выберите нужный язык/наречие.

Нажмите кнопку **ПОКАЗАТЬ**.

Загрузится таблица с текстами (Рис. 11).

Чтобы перейти на страницу с текстом, щелкните на заголовок текста.

## ★ Тексты

[Лексико-грамматический поиск](#) | [Создать новый](#) | [? Помощь](#)

[Расширенный поиск ↓](#)

Язык: \* карельский: людиковское наречие (255) | Подкорпус: | по 10 записей **ПОКАЗАТЬ**

Найдено 255 записей.

No	Язык	Диалект	Подкорпус	Жанр	Заголовок	Перевод
1	карельский: людиковское наречие	Новописьменный людиковский	публицистические тексты		<a href="#">ADVOIŠ</a>	В ГОСТЯХ
2	карельский: людиковское наречие	Южнолюдиковский (святозерский)	диалектные тексты	бытовой рассказ	<a href="#">Agđan St'opan</a>	Степан Агдян
3	карельский: людиковское наречие	Михайловский	диалектные тексты	бытовой рассказ	<a href="#">Äjj on dielot kalan suada</a>	
4	карельский: людиковское наречие	Южнолюдиковский (святозерский)	диалектные тексты	бытовой рассказ	<a href="#">Akimat</a>	Семейство Акиман
5	карельский: людиковское наречие	Новописьменный людиковский	публицистические тексты		<a href="#">Paušin Šan'uu. Aleksandr Barancev on meiden kuuluž heimolaine</a>	

Рис. 11. Поиск в корпусе *ВепКар*. Пример найденных людиковских текстов

### Шаг 9. Как осуществлять расширенный поиск по текстам *ВепКар*?

Выполните шаг 7.

Щелкните на ссылку «*расширенный поиск*» справа вверху над поисковой формой (Рис. 11).

В форме откроются дополнительные поля (Рис. 12). Под поисковой формой выводится количество найденных текстов и таблица с данными. Если сформирован поиск по языку, подкорпусу, диалекту или жанру, то соответствующие колонки в таблице не дублируются.

Язык	Диалект	Подкорпус
Жанр	Сюжет	Тема
Область, республика записи	Район записи	Населенный пункт записи
Область, республика рожд. информанта	Район рожд. информанта	Населенный пункт рожд. информанта
Информант	Собиратель	Автор текста или перевода
Заголовок	Год (с)	Год (по)
Источник	Источники	
Фрагмент текста	<input type="checkbox"/> с аудио-записями	
		по 10 записей <b>ПОКАЗАТЬ</b>

Найдено 4 129 записей.

Простой поиск ↑

№	Язык	Диалект	Подкорпус	Жанр	Заголовок	Перевод
1	карельский: собственно карельское наречие	Валдайский	диалектные тексты		Mar'joa möimmä imen'jah	Пунжина Александра Васильевна. Ягоды продавали в имение

Рис. 12. Расширенный поиск текстов в корпусе VenKar

## Шаг 10. Как найти озвученный текст на вепском / карельском языке?

- Выполните шаг 9.
- В поле «Выберите язык» выберите нужный язык/наречие.
- Кликните в поле «с аудио-записями».
- Нажмите кнопку **ПОКАЗАТЬ**.

Загрузится таблица с текстами (Рис. 13).

Язык	<input checked="" type="checkbox"/> с аудио-записями	по 10 записей <b>ПОКАЗАТЬ</b>
* карельский: ливвиковское наречие (959)		

Найдено 15 записей.

№	Диалект	Подкорпус	Жанр	Заголовок	Перевод
1	Рыпушкальский	диалектные тексты		Kaikin paištih hierus livvikse	Мичурова Надежда. Все говорили в деревне на ливвиковском
2	Некульский	диалектные тексты		Joga talois oli lehmy	Мичурова Надежда. В каждом доме была корова
3	Коткозерский	диалектные тексты		Vastaimmo mejän saldattoi	Мичурова Надежда. Встречали наших солдат
4	Рыпушкальский	диалектные тексты	бытовой рассказ	Poimiččulojn da viršilöjn azundu	Мичурова Надежда. Изготовление маленьких и больших корзин

Рис. 13. Поиск в корпусе VenKar. Пример найденных ливвиковских текстов с аудио-записями.

## Шаг 11. Как найти вепские сказки, записанные до 1950 года?

- Выполните шаг 9.
- В поле «Язык» выберите «вепский».
- В поле «Жанр» выберите «сказка».
- В поле «Год (по)» укажите '1950'.

Нажмите кнопку **ПОКАЗАТЬ**.

Расширенный поиск ↓

<b>Язык</b>	<b>Жанр</b>	<b>Год (по)</b> ?
» вепский (1 748)	» сказка	1950
по 10	записей	<b>ПОКАЗАТЬ</b>

Найдено 20 записей.

№	Диалект	Подкорпус	Заголовок	Перевод
1	Средневепский западный	фольклорные тексты	<a href="#">Pihkmüt da kivut</a>	Скатёрка и жерновок
2	Средневепский западный	фольклорные тексты	<a href="#">Eli akaine, oli hānou poig</a>	Жила женщина, был у нее сын
3	Северновепский	фольклорные тексты	<a href="#">Viikuško čapab sizarel käded</a>	Брат отрубает руки у сестры

Рис. 14. Поиск вепских сказок, записанных до 1950 года.

## Шаг 12. Как осуществить поиск примеров употребления слова в корпусе?

- В верхнем меню найдите пункт «КОПИУС» и щелкните на него мышкой.
- В выпадающем меню щёлкните по ссылке «Лексико-грамматический поиск» (Рис. 15).

### ✳ Лексико-грамматический поиск

Расширенный поиск ↓

<b>Язык</b>	<b>Подкорпус</b>	по 10	записей	<b>ИСКАТЬ</b>
<b>Слово 1</b> ? ã	<b>Часть речи</b> ? +	<b>Грамматические признаки</b> ? +		

Рис. 15. Поиск в корпусе заданной последовательности лемм и/или словоформ, обладающих определенными грамматическими характеристиками

В поле «Слово» можно задать начальную форму (лемму) и получить все формы этой лексемы, найденные в корпусе. Если нужно найти одну определенную форму слова, укажите ее в кавычках.

**Например**, для поисковой строки  $\wedge$ vuozis\$ будут найдены примеры вхождения *vuuvvel*, *vuodel*, и т.д. Если нужно найти только *vuuvvel*, задайте эту форму в кавычках: " $\wedge$ vuuvvel\$".

Иконки **+** открывают дополнительные окна для полей «Часть речи» и «Грамматические признаки». Галочками отметьте нужные значения.

## Шаг 13. Как найти примеры употребления в корпусе людиковских имен в форме аппроксиматива?

- Выполните шаг 12.
- В поле «Выберите язык» выберите «карельский: людиковское наречие».
- В поле «Часть речи» первого слова нажмите на иконку **+**.

Откроется дополнительное окно для выбора частей речи (Рис. 16).

Выберите часть речи x

<input type="checkbox"/> вспомогательный глагол	<input type="checkbox"/> глагол	<input type="checkbox"/> другое
<input type="checkbox"/> имя собственное	<input type="checkbox"/> междометие	<input type="checkbox"/> местоимение
<input type="checkbox"/> наречие	<input type="checkbox"/> послелог	<input type="checkbox"/> предикатив
<input type="checkbox"/> предлог	<input type="checkbox"/> прилагательное	<input type="checkbox"/> союз подчинительный
<input type="checkbox"/> союз сочинительный	<input checked="" type="checkbox"/> существительное	<input type="checkbox"/> фразеологизм
<input type="checkbox"/> частица	<input type="checkbox"/> числительное	

Рис. 16. Выбор части речи «существительное» в дополнительном окне.

- Отметьте галочкой «существительное».
- Нажмите кнопку **выбрать**.
- В поле «Грамматические признаки» первого слова нажмите на иконку **+**.

Откроется дополнительное окно для выбора грамматических признаков.

- Отметьте галочкой падеж «аппроксиматив».
- Нажмите кнопку **выбрать**.
- В основной форме (Рис. 17) нажмите кнопку **ИСКАТЬ**.

В новой вкладке загрузится результат поиска (Рис. 18). Стрелочки слева и справа предложения добавляют левый и правый контекст соответственно. Если щелкнуть на слово в предложении, откроется окошко с дополнительной грамматической и семантической информацией о слове.

Расширенный поиск ↓

Язык:  Подкорпус:  по  записей

Слово 1:  Часть речи:  Грамматические признаки:  

Рис. 17. Поиск в корпусе примеров употребления людиковских имен в форме аппроксиматива

## Результаты поиска

(язык: карельский: людиковское наречие) И

Слово 1: (существительное) И ((аппроксиматив))

Найдено 3 текста, 3 вхождения.

[Уточнить запрос](#)

1. Paušin Šan'uu. [Ken ei suvaiče hiiiriid, sil pidäy armastada kažid](#) (Oma mua. № 8, 2021, с. 10)  
 ← Konz hyö oldih kuhn'al **päçilluo**, a minä magazin kätkydes i kieras rubezin itkemäh, Anapka ylen teravas hypii minulluo i kaččeli minun silmiih, i minä hedi tyyništuin. →
2. Paušin Šan'uu. [Brus'a-kažiine](#) (Oma Mua. № 34 (1624), 2022, с. 9)  
 ← Konz hyö oldih kuhn'al **päçilluo**, a minä magazin kätkydes i kieras rubezin etkemäh, Anapka ylen teravas hypii minulluo i kaččeli minun silmiih, i minä hedi tyyništuin. →
3. [Lapsel oli yönit'kettai](#) / У ребёнка была ночная плачя (Спасская Губа (Mundärven laht), Кондопожский район, Республика Карелия, 2012)  
 ← Šie vot mina kakš yön magain sel **babuškalluo** tulin da oma rodn'a da vs'o. →

Рис. 18. Результаты поиска примеров людиковских имен в форме аппроксиматива

## Шаг 14. Как найти примеры употребления в корпусе ливвиковских глаголов в форме перфекта и плюсквамперфекта кондиционала?

- Выполните шаг 12.

- В поле «Язык» выберите «карельский: ливвиковское наречие».
- В поле «Слово 1» укажите лемму ‘olla’<sup>16</sup>.
- В поле «Часть речи» первого слова нажмите на иконку **+**.
- В дополнительном окне (Рис. 16) отметьте галочкой «глагол».
- Нажмите кнопку **выбрать**.
- В поле «Грамматические признаки» первого слова нажмите на иконку **+**.
- В дополнительном окне отметьте галочкой наклонение «кондиционал».
- Нажмите кнопку **выбрать**.
- Справа нажмите на иконку **+**.

На месте плюса появятся поля «Расстояние». Ниже откроются аналогичные поля для второго слова.

- В поле «Часть речи» второго слова нажмите на иконку **+**.
- В дополнительном окне отметьте галочкой «глагол».
- Нажмите кнопку **выбрать**.
- В поле «Грамматические признаки» второго слова нажмите на иконку **+**.
- В дополнительном окне отметьте галочками залог «актив» и «2-е причастие».
- Нажмите кнопку **выбрать**.
- В основной форме (Рис. 19) нажмите кнопку **ИСКАТЬ**.

В новой вкладке браузера загрузится результат поиска (Рис. 20). Если щелкнуть на слово в предложении, откроется окошко с дополнительной грамматической и семантической информацией о слове.

Рис. 19. Поиск в корпусе примеров употребления ливвиковских глаголов в форме перфекта и плюсквамперфекта кондиционала

<sup>16</sup> Чтобы отсеять варианты лемм, которые содержат подстроку ‘olla’, можно использовать шаблон ^olla\$.

## Результаты поиска

(язык: карельский: ливвиковское наречие) И

Слово 1: *olla* И (глагол) И ((кондиционал)) И

на расстоянии от 1 до 1 Слово 2: (глагол) И ((актив) И (2-е причастие))

Найдено 9 текстов, 9 вхождений.

1. [Bul'uu borkananke](#) (Irina Kudel'nikova, Oma mua. № 1, 2018, с. 11)

← Vie pidäs teis dengat ottua, ka olgah, prostin. Äjjän rahvastu **olluzin parandannuh**, a työ kallehen syömizen kaimaitto!

2. [Kargiet voinuaijat](#) (Fodorova Anni (nygöi Ivanova), Oma mua. № 24, 2018, с. 11)

← Kai kylä **ollus palanuh**, ku vahnembat ei ehtittys. →

3. [Ken da kui kalasti](#) (NIKOLAI FILATOV, Oma mua. № 12, 2020, с. 11)

← Tiettäväine, emmo **olis tundenuh** heidy, kerdu kolme vuottu olimmo

4. [Ken da kui kalasti](#) (NIKOLAI FILATOV, Oma mua. № 12, 2020, с. 11)

← Tiettäväine, emmo **olis tundenuh** heidy, kerdu kolme vuottu olimmo

5. [Kuldastu laduu ei umbua lumel](#) (Raisa Bogdanova, Oma mua. № 5, 2018, с. 11)

← Tanävuon Fodor Terentjev **olis täyttänh** 95 vuottu. →

parandannuh

parandua

глагол

1) **улучшать, улучшить;**

**совершенствовать,**

**усовершенствовать**

2) **лечить, вылечить,**

**исцелять, исцелить**

**- актив, 2-е причастие**

mmo eigo hyögi oldas ajamas. →

mmo eigo hyögi oldas ajamas. →

Рис. 20. Результаты поиска примеров ливвиковских глаголов в форме перфекта и плюсквамперфекта кондиционала

### Шаг 15. Как найти примеры использования глаголов с существительными в определённых падежах, например, в партитиве?

- Выполните шаг 12.
- В поле «Часть речи» первого слова нажмите на иконку **+**.
- В дополнительном окне (Рис. 16) отметьте галочкой «глагол».
- Нажмите кнопку **выбрать**.
- Справа нажмите на иконку **+**

На месте плюса появятся поля «Расстояние». Ниже откроются аналогичные поля для второго слова.

- В поле «Часть речи» второго слова нажмите на иконку **+**.
- В дополнительном окне отметьте галочкой «существительное».
- Нажмите кнопку **выбрать**.
- В поле «Грамматические признаки» второго слова нажмите на иконку **+**.
- В дополнительном окне отметьте галочкой падеж «партитив».
- Нажмите кнопку **выбрать**.
- В основной форме нажмите кнопку **ИСКАТЬ**.

Такой запрос является достаточно «тяжёлым» и если загрузка результата затянулась на несколько минут, попробуйте уточнить запрос. Выбор языка и других параметров тексту ускоряет процесс поиска примеров. Например, нас интересуют только ливвиковские фольклорные тексты.

- В поле «Язык» выберите «карельский: ливвиковское наречие».
- В поле «Подкорпус» выберите «фольклорные тексты».
- В основной форме (Рис. 21) нажмите кнопку **ИСКАТЬ**.

Расширенный поиск ↓

<b>Язык</b>	<b>Подкорпус</b>	по	10	записей	<b>ИСКАТЬ</b>
× карельский: ливвиковское наречие (954)	× фольклорные тексты (307)				

<b>Слово 1</b>	<b>Часть речи</b>	<b>Грамматические признаки</b>	<b>Расстояние</b>
	VERB		от 1 до 1
<b>Слово 2</b>	<b>Часть речи</b>	<b>Грамматические признаки</b>	
	NOUN	PRT	

Рис. 21. Поиск в корпусе примеров употребления глаголов с существительными в партитиве

В новой вкладке браузера загрузится результат поиска (Рис. 22).

## Результаты поиска

(язык: карельский: ливвиковское наречие) И (подкорпус: фольклорные тексты) И

**Слово 1:** (глагол) И

на расстоянии от 1 до 1 **Слово 2:** (существительное) И ((партитив))

Найдено 26 текстов, 124 вхождения.

[Уточнить запрос](#)

- [Варка da d'etka ištutettih nagrehen](#) / Бабка и дедка посадили репку (Карельская (Karjažet), Лодейнопольский р-н, Ленинградская обл., 1959)  
← Siiten koiru **kučui kaži:** veettih, veettih viijei – ei voittu ni kui vediä. →
- [\[ivanuška-boranuška\]](#) / [Иванушка-баранушка] (Самбатукса (Sammatus), Олонецкий район, Республика Карелия, 1959)  
← Brihačču sanowgi: "Oi, čidžoini, minä täs **juan vettü** lehmän kabjan jäl'les!" →  
← **Ruvetah händü** tappamah. →  
← Häi menöw, itköw-itköw: "Oi, sestrica, Ol'onuska, minuw ruvetah iškemäh; veiččii **tahkotah, kirvehii** hivotah, kattilat kiahutah, minuw iškemäh ruvetah!" →  
← Häi ku menöw vezirandah dai **sanow:** "Oi, sestrica Ol'onuska! →  
← Minuw iškemäh **ruvetah: veiččii hivotah, kirvehii** tahkotah, minuw ruvetah iškemäh!" →
- [Suaru hukkah niškoi](#) / Сказка про волка (Самбатукса (Sammatus), Олонецкий район, Республика Карелия, 1959)  
← "Na, täs, akku, **sanow, lihua,** keitä lapsil'e lihua, sanow, da piästäh n'al'1'äs hot' iäres". →

Рис. 22. Результаты поиска в корпусе примеров употребления глаголов с существительными в партитиве

## Шаг 16. Как найти примеры употребления конкретной словоформы в вепсских текстах 30х годов?

- Выполните шаг 12.
- В поле «Выберите язык» выберите «вепсский».
- Нажмите «Расширенный поиск» (Рис. 15).
- В поле «Год (с)» впишите '1930'.
- В поле «Год (по)» впишите '1940'.
- В поле «Слово 1» укажите лемму или шаблон, например '^ö\$'<sup>17</sup>.
- В поле «Грамматические признаки» нажмите на иконку **+**
- В дополнительном окне отметьте галочкой нужные грамматические категории, например, падеж «адессив».
- Нажмите кнопку **выбрать**.
- В основной форме (Рис. 23) нажмите кнопку **ИСКАТЬ**.

<sup>17</sup> По шаблону ^ö\$ будут найдены формы слова ö. Если написать просто 'ö', то будут найдены формы всех слов с входящей буквой ö: söda, hö, löuta и т.д. Подробнее о шаблонах в таблице 4.

Рис. 23. Поиск в корпусе примеров употребления вепской лексики 'ö' в форме адессива в текстах 30-х годов

В новой вкладке браузера загрузится результат поиска (Рис. 24).

Рис. 24. Поиск в корпусе примеров употребления вепской лексики 'ö' в форме адессива в текстах 30-х годов

### Шаг 17. Как найти частотные словари корпуса?

- Выполните шаги 1 и 2.
- В верхнем меню найдите пункт «КОРПУС» и щелкните на него мышкой.
- В выпадающем меню щёлкните по ссылке «Частотные словари».

Загрузится страница со списком частотных словарей (Рис. 25).

### ★ Частотные словари

- [Частотный словарь токенов](#)
- [Частота символов в текстах](#)
- [Частотный словарь лемм, размеченных в тексте](#)

Рис. 25. Страница с частотными словарями корпуса

### Шаг 18. Как найти самые частотные слова в текстах?

- Выполните шаг 17.
- Щелкните по ссылке «Частотный словарь лемм, размеченных в текстах».

Откроется страница с поисковой формой (Рис. 26).

Рис. 26. Частотный словарь лемм, размеченных в тексте

- В полях «Язык», «Диалект», «Части речи» выберите нужные значения (или оставьте пустыми поля).
- Нажмите кнопку **ПОКАЗАТЬ**.

Выведется список слов, встречающихся в текстах (Рис. 27).

Язык  
карельский: ливвиков... ▾

Диалект  
Новописьменный ливвиковск... x ▾

Часть речи  
существительное (552 ▾

**ПОКАЗАТЬ**

Показать 10 ▾ записей

Поиск:

No	⇕ Лемма	⇕ Часть речи	⇕ Частота	⇕
1763	<a href="#">oma</a>	существительное	<a href="#">1629</a>	
1055	<a href="#">sana</a>	существительное	<a href="#">1231</a>	
175	<a href="#">vuozǵi</a>	существительное	<a href="#">1195</a>	
2470	<a href="#">kylä</a> <span style="color: red;">сортировка по лемме</span>	существительное	<a href="#">1125</a>	
1274	<a href="#">rahvas</a>	существительное	<a href="#">1070</a>	
3260	<a href="#">hyvä</a>	существительное	<a href="#">987</a>	
3847	<a href="#">aigu</a>	существительное	<a href="#">979</a>	
2964	<a href="#">karjal</a>	существительное	<a href="#">897</a>	
2712	<a href="#">kodi</a>	существительное	<a href="#">890</a>	
96	<a href="#">äijy</a>	существительное	<a href="#">854</a>	

Записи с 1 до 10 из 3,876 записей

Предыдущая 1 2 3 4 5 ... 388 Следующая

Рис. 27. Частотный словарь ливвиковских существительных, размеченных в новописьменных текстах

По умолчанию список отсортирован по частоте встречаемости слова в текстах. Можно менять сортировку, нажимая стрелочки СПРАВА от заголовка колонки (Рис. 27).

Стоит отметить, что считаются все связи слова в тексте со словарной статьей (в том числе и непроверенные), кроме связей с меткой "совсем не подходит". Если слово в тексте автоматически было "привязано" к лемме с несколькими значениями и эти значения еще не проверены (не снята семантическая омонимия), то будут посчитаны все связи с этими значениями.

### Шаг 19. Как найти «Речевой корпус»?

- Выполните шаги 1 и 2.
- В верхнем меню найдите пункт «КОРПУС» и щелкните на него мышкой.
- В выпадающем меню щёлкните по ссылке «Речевой корпус».

Выведется список озвученных текстов (Рис. 28).

## ★ Речевой корпус

[Аудио-карта](#) | [Создать новый](#) | [? Помощь](#)

[Расширенный поиск ↓](#)

Язык  по  записей [ПОКАЗАТЬ](#)

Найдено 46 записей.

№	Диалект	Говор	Заголовок	Перевод	Прослушать
1	Вокнаволоцкий	<a href="#">Суднозеро</a>	<a href="#">Meijen talo šuuri oli</a>	Наш дом большой был	
2	Рыпушкальский	<a href="#">Алексала</a>	<a href="#">Kaikin paištih hierus livvikse</a>	Мичурова Надежда. Все говорили в деревне на ливвиковском	
3	Неккульский	<a href="#">Обжа</a>	<a href="#">Joga talois oli lehmy</a>	Мичурова Надежда. В каждом доме была корова	

Рис. 28. Речевой корпус прибалтийско-финских языков Карелии

В колонке «Говор» выводится населенный пункт - место рождения информанта.

### Шаг 20. Как найти и прослушать образцы речи интересующего населённого пункта?

- Выполните шаг 19.
- Нажмите на ссылку «*Расширенный поиск*» (Рис. 28).
- Если Вас интересует населенный пункт, где были собраны образцы, то выберите нужное значение в списке «*Населенный пункт записи*».
- Если Вас интересует населенный пункт, в котором родились информанты, то выберите нужное значение в списке «*Населенный пункт рожд. информанта*» (Рис. 29).

Язык  Диалект  Заголовок

Область, республика записи  Район записи  Населенный пункт записи

Область, республика рожд. информанта  Район рожд. информанта  Населенный пункт рожд. информанта 

- Виданы
- Видлица
- Виллала
- Винжа
- Воздвиженье
- Корпяярви

Информант  Собиратель

Источник  по  записей [ПОКАЗАТЬ](#)

[? ä](#) ↑

Рис. 29. Выбор населенного пункта рождения информанта

## Шаг 21. Как ссылаться на материалы корпуса ВепКар?

Если вы используете материалы корпуса в своей исследовательской работе, то вы можете выбрать следующие варианты.

1. В тексте или в сноске указать:

Исследование проведено на материале Открытого корпуса вепского и карельского языков ([dictorpus.krc.karelia.ru](http://dictorpus.krc.karelia.ru)).

2. В списке литературы в конце работы привести ссылку на Корпус:

Открытый корпус вепского и карельского языков (ВепКар). 2009—2022. Доступен по адресу: [dictorpus.krc.karelia.ru](http://dictorpus.krc.karelia.ru)

3. Ссылку на конкретный текст оформляйте, например, так:

Lehme № 2051 // Открытый корпус вепского и карельского языков ВепКар. URL: <http://dictorpus.krc.karelia.ru/ru/corpus/text/2051>

Отметим, что вместо ссылки на весь корпус, вы можете поставить ссылку на какую-либо из наших публикаций о корпусе. Этот список есть на сайте ВепКар: <http://dictorpus.krc.karelia.ru/ru/page/publ>

### 2.2.3. Поиск в словаре ВепКар

#### Шаг 22. Как осуществлять поиск в словаре ВепКар?

- Выполните шаги 1 и 2.
- В верхнем меню найдите пункт «СЛОВАРЬ» и щелкните на него мышкой.
- В выпадающем меню щёлкните по ссылке «Леммы» или «Словоформы» для поиска в словаре соответственно лемм или словоформ (Рис. 30).

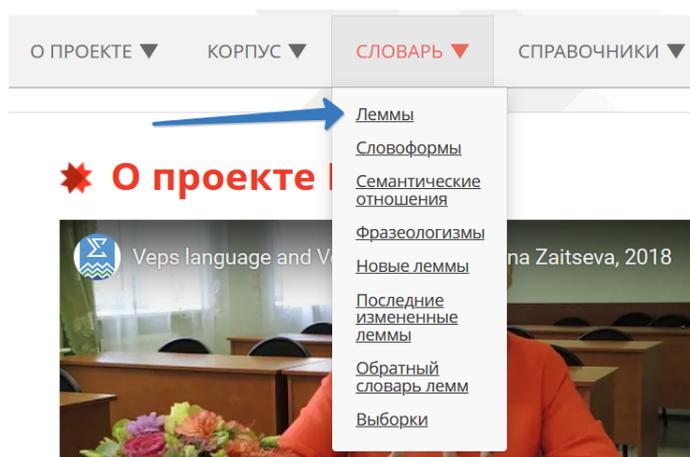


Рис. 30. Поиск лемм в словаре ВепКар

#### Шаг 23. Как найти вепское / карельское слово по начальной форме?

- Выполните шаг 22, щелкнув по ссылке «Леммы».
- Укажите вепское (карельское) слово в поле «лемма».
- Нажмите кнопку **ПОКАЗАТЬ** (Рис. 31).

При поиске Вы можете задавать фрагменты леммы, используя простые шаблоны (см. Шаг 5).

## ★ Леммы

[Поиск лемм по словоформам](#) | [Список длинных лемм](#) | [Создать новую](#)

[Расширенный поиск ↓](#)

Найдено 3 записи.

№	Лемма	Язык	Часть речи	Толкование	Словоформы *	Примеры **
1	<a href="#">маа</a>	карельский: собственно карельское наречие	существительное	земля, суша	2	0
2	<a href="#">маа</a>	карельский: ливвиковское наречие	существительное	1) земля, суша 2) почва	2	0 / 2 / 2
3	<a href="#">маа</a>	вепсский	существительное	1) земля, суша 2) почва	2	0 / 2 / 2

\* - Количество словоформ с грамматическими признаками [+ кол-во словоформ без грамматических признаков]

\*\* - Количество проверенных примеров / Количество непроверенных примеров / Общее количество

Рис. 31. Поиск лемм в словаре ВепКар по начальной форме слова. Поле «лемма» с иконками «помощь» и «специальные символы». Ссылка «расширенный поиск» для открытия полной формы

### Шаг 24. Как указать язык/наречие при поиске леммы?

- Выполните шаг 22, щелкнув по ссылке «Леммы».
- Щелкните на поле «Выберите язык».
- Выберите нужный язык/наречие (Рис. 32).

Выберите язык

- Выберите язык
- вепсский (18 618)
- карельский: ливвиковское наречие (27 445)**
- карельский: людиковское наречие (6 508)
- карельский: собственно карельское наречие (10 693)

Рис. 32. Выбор языка из выпадающего списка

### Шаг 25. Как осуществлять расширенный поиск в словаре ВепКар?

- Выполните шаг 22, щелкнув по ссылке «Леммы».
- Щелкните на ссылке «расширенный поиск» справа вверху над поисковой формой (Рис. 31).

В форме откроются дополнительные поля (Рис. 33).

Рис. 33. Расширенный поиск лемм в словаре VenKar

## Шаг 26. Как найти вепсское (карельское) слово по русскому толкованию?

- Выполните шаг 25.
- Введите текст в поле «толкование» (Рис. 33).
- Нажмите кнопку **ПОКАЗАТЬ** (Рис. 34).

Найдено 14 записей.

№	Лемма	Толкование	Словоформы *	Примеры **
1	<a href="#">jumalanlehmäine</a>	божья коровка	37	0
2	<a href="#">lehmy</a>	корова	37	5 / 113 / 118
3	<a href="#">lehmykarju</a>	стадо коров	37	0 / 1 / 1
4	<a href="#">lehmänkello</a>	колокольчик на шее у коровы	37	0 / 1 / 1
5	<a href="#">lehmänmaido</a>	коровье молоко	37	0
6	<a href="#">lehmännahku</a>	шкура коровы	37	0
7	<a href="#">lehmänsarvi</a>	рог коровы	54	0
8	<a href="#">lehmö</a>	корова	2	0

Рис. 34. Поиск лемм в словаре VenKar по языку (ливвиковское наречие карельского языка), части речи (существительное) и толкованию (%коров%)

## Шаг 27. Как найти вепсское (карельское) слово по словоформе?

Можно искать слово не только по начальной, но и по любой его грамматической форме - словоформе.

- Выполните шаг 22, щелкнув по ссылке «Леммы».
- Щелкните на ссылку «Поиск лемм по словоформам»<sup>18</sup> (Рис. 31).

Откроется страница с новой поисковой формой (Рис. 35). В этом поиске Вы можете указать словоформу (слово или специализированный шаблон, см. шаг 6).

- Если известна форма слова, заполните поле «словоформа 1».
- Если известны грамматические признаки слова, выберите нужное значение в списке «грамматические признаки для словоформы».
- Нажмите кнопку **ПОКАЗАТЬ**.

<sup>18</sup> Поиск лемм по словоформам: [http://dictorpus.krc.karelia.ru/ru/dict/lemma/by\\_wordforms](http://dictorpus.krc.karelia.ru/ru/dict/lemma/by_wordforms)

## ✳ Поиск лемм по словоформам

Расширенный поиск | Создать новую

Выберите язык  диалект  Выберите часть речи

словоформа 1  Грамматические признаки для словоформы

**ПОКАЗАТЬ** по  записей [добавить словоформу](#)

Рис. 35. Поиск лемм в словаре VenKar по словоформам

**Шаг 28. Как найти все глаголы новописьменного севернокарельского языка, начинающиеся с буквы ‘т’, у которых начальная форма (инфинитив) оканчивается на ‘uo’ или ‘yö’, а форма 1 лица ед.ч. настоящего времени оканчивается на ‘un’ или ‘yn’?**

- Выполните шаг 27.
- В поле «язык» выберите «карельский: собственно карельское наречие».
- В поле «диалект» выберите «Новописьменный севернокарельский».
- В поле «часть речи» выберите «глагол».
- В поле «словоформа 1» введите шаблон «<sup>^</sup>t.+[uy][oö]\$»
- В поле справа выберите «1 инфинитив»
- Нажмите иконку «добавить словоформу» (Рис. 35).
- В новом поле «словоформа 2» введите шаблон «[uy]n\$»
- В поле справа выберите «индикатив, презенс, 1 л., ед. ч., полож. ф.»
- Нажмите кнопку **ПОКАЗАТЬ** (Рис. 36).

карельский: собственно карельско  \* Новописьменный севернокарельский  глагол (13 318)

<sup>^</sup>t.+[uy][oö]\$  131. I инфинитив

[uy]n\$  1. индикатив, презенс, 1 л., ед. ч., полож. ф.

**ПОКАЗАТЬ** по  записей

Найдено 25 записей.

№	Лемма	Толкование	Словоформы *	Примеры **
1	<a href="#">mahtuo</a>	1) вмещаться, входить 2) ладить, жить в ладу	242 (mahtuo, mahun)	0 / 23 / 23
2	<a href="#">maistuo</a>	иметь какой-либо вкус	125 (maistuo, maissun)	0 / 4 / 4
3	<a href="#">mataltuu</a>	становиться ниже, мелеть	125 (mataltuu, matallun)	0 / 1 / 1
4	<a href="#">matkeutuo</a>	1) собираться в дорогу 2) прийти (о новости)	125 (matkeutuo, matkeuvun)	0 / 3 / 3
5	<a href="#">matoutuo</a>	зачервиветь	125 (matoutuo, matouvun)	0
6	<a href="#">mehuo</a>	течь, просачиваться	125 (mehuo, mehun)	0 / 1 / 1
7	<a href="#">meruutuo</a>	плавиться, топиться (о масле, жире)	125 (meruutuo, meruuvun)	0
8	<a href="#">mečittyö</a>	зарастать лесом	125 (mečittyö, mečityn)	0 / 2 / 2

Рис. 36. Пример поиска всех глаголов новописьменного севернокарельского языка, начинающихся с буквы ‘т’, у которых начальная форма (инфинитив) оканчивается на ‘uo’ или ‘yö’, а форма 1 лица ед.ч. настоящего времени оканчивается на ‘un’ или ‘yn’

**Шаг 29. Как найти обратный словарь?**

## ✳ Обратный словарь лемм

Выберите язык  Выберите часть речи  по 10 записей **ПОКАЗАТЬ**

Рис. 37. Поиск по обратному словарю лемм.

- Выполните шаги 1 и 2.
- В верхнем меню найдите пункт «СЛОВАРЬ» и щелкните на него мышкой.
- В выпадающем меню щелкните по ссылке «Обратный словарь лемм».
- В загрузившейся форме (Рис. 37) щелкните на поле «Выберите язык» и выберите нужный язык/наречие.
- В поле «Выберите часть речи» выберите нужную часть речи (если нужно).
- Нажмите кнопку **ПОКАЗАТЬ**.

В загрузившейся таблице (Рис. 38) будут выведены слова, отсортированные по алфавиту по последней букве.

карельский: людиковское нареч  Выберите диалекты  существительное (37 202)

лемма  по 10 записей **ПОКАЗАТЬ**

Найдена 1 561 запись.

№	Лемма	Часть речи
101	šleja	NOUN
102	agja	NOUN
103	sija	NOUN
104	magadussija	NOUN
105	oja	NOUN
106	koja	NOUN
107	kirja	NOUN
108	nahka	NOUN
109	pahka	NOUN
110	tähkä	NOUN

Рис. 38. Поиск по обратному словарю лемм

## Заключение

В этом руководстве для пользователей приведены основные понятия корпусной лингвистики, описана разметка, типология и обзор современных лингвистических корпусов. Показаны основные виды работы с лингвистическим корпусом ВеКар: регистрация, переключение языка интерфейса, поиск в словаре и текстах с помощью простых и специализированных шаблонов. Приведены примеры сложных запросов, в которых указывается сразу несколько (много) параметров поиска.

Надеемся, что предложенное пособие будет востребовано в сферах среднего и высшего образования в процессе преподавания вепского и карельского языков, а также корпусной лингвистики. Замечания и предложения по улучшению руководства можно отправлять разработчикам корпуса ВеКар по адресу: [nataly@krc.karelia.ru](mailto:nataly@krc.karelia.ru).

## Тезаурус

**Корпус (лингвистический)** — коллекция текстов, специально отобранных, размеченных по различным лингвистическим параметрам и обеспеченных системой поиска.

**Корпус аннотированный/размеченный** — корпус текстов, в котором содержатся специальные метки, позволяющие получать из корпуса данные (статистику, языковые примеры и др.) по каким-либо лингвистическим параметрам (части речи, грамматической форме, синтаксической функции и т. п.) [Захаров, Богданова, 2020].

**Корпусная лингвистика** — раздел компьютерной лингвистики, занимающийся разработкой общих принципов построения и использования лингвистических корпусов (корпусов текстов) с применением компьютерных технологий [Захаров, Богданова, 2020].

**Корпусный менеджер (корпус-менеджер)** — специальная информационно-поисковая система, включающая программные средства для поиска данных в корпусе, получения статистической информации.

**Лемма** - 1) начальная форма слова; 2) словарная статья в словаре ВепКар.

**Лексико-грамматический поиск** — поиск заданной последовательности лемм и/или словоформ, обладающих определенными грамматическими характеристиками.

**Миноритарный язык** — язык национального (этнического) меньшинства [ССТ: 129].

**Морфологический анализ** — определение леммы и ее грамматических характеристик.

**Объем корпуса** — информация об общем объеме корпуса и о количестве извлеченных из текста примеров.

**Разметка (аннотация) текста** — введенная вручную или автоматически лингвистическая или метатекстовая информация обо всех выбранных единицах текста (тексте, предложении, словосочетании, словоформе, морфеме).

**Репрезентативность** — свойство корпуса, заключающееся в статистически достоверном представлении языка или его части и достигаемое за счет необходимого объема и жанрового многообразия текстов [Копотев, 2014].

**Сбалансированность** — свойство корпуса, определяющее насколько равномерно представлены тексты разных типов.

## Список использованной литературы

[Albert, 2012] Albert, Saul. The Audio BNC. 2012. URL: <https://saulalbert.net/blog/the-audio-bnc/> (дата обращения: 04.11.2021).

[Arkhangelskiy, 2020] Arkhangelskiy, T., 2020. Web Corpora of Volga-Kama Uralic Languages // Finno-Ugric Languages and Linguistics. Vol. 9. No. 1-2. (2020), P. 58–66.

- [Klyachko, 2019] Klyachko E. L., Sorokin A. A., Krizhanovskaya N. B., Krizhanovsky A. A., Ryazanskaya G. M. LowResourceEval-2019: a shared task on morphological analysis for low-resource languages // Computational Linguistics and Intellectual Technologies : papers from the Annual conference “Dialogue”. — М.: РГГУ, 2019. — Вып. 18 (25). — С. 45—62.
- [Krizhanovskaya et al., 2022] Krizhanovskaya, N., Novak, I., Krizhanovsky, A., & Pellinen, N. Morphological inflectional rules for Karelian Proper verbs. Eesti Ja Soome-Ugri Keeleteaduse Ajakiri. Journal of Estonian and Finno-Ugric Linguistics, 13(2). 2022. Pp. 47–78. <https://doi.org/10.12697/jeful.2022.13.2.02>
- [Meyer, 2002] Meyer, Charles. English Corpus Linguistics: An Introduction. Cambridge: Cambridge University Press, 2002. URL: [http://assets.cambridge.org/97805218/08798/frontmatter/9780521808798\\_frontmatter.pdf](http://assets.cambridge.org/97805218/08798/frontmatter/9780521808798_frontmatter.pdf) (дата обращения: 07.06.2021).
- [Баранов, 2001] Баранов А. Н. Введение в прикладную лингвистику. – М, 2001. – Т. 2.
- [Бойко и др., 2021] Бойко Т. П., Зайцева Н. Г., Крижановская Н. Б., Крижановский А. А., Новак И. П., Пеллинен Н. А., Родионова А. П., Трубина Е. Д. Лингвистический корпус ВепКар — «заповедник» прибалтийско-финских языков Карелии // Труды КарНЦ РАН. No 7. Комплексные научные исследования КарНЦ РАН. 2021. С. 100-115 DOI: 10.17076/them1415
- [Герд, 2008] Герд А. С. Национальный корпус русского языка – Словарная картотека – Академический словарь // Тр. Междунар. конф. «Корпусная лингвистика–2008» – 2008. – с. 143–148. URL: [http://www.project.phil.spbu.ru/corpora2011/Works2008/Gerd\\_143\\_148.pdf](http://www.project.phil.spbu.ru/corpora2011/Works2008/Gerd_143_148.pdf) (дата обращения: 21.06.2021).
- [Зайцева, 2012] Зайцева Н.Г. Вепские причитания в фокусе корпусной лингвистики и лингвофольклористики // Материалы ХLI Международной филологической конференции. 26-31 марта 2012 г. Секция "Уралистика". СПб.: Филологический факультет СПбГУ. 2012. С. 16-26.
- [Захаров, Богданова, 2020] Захаров В. П., Богданова С. Ю. Корпусная лингвистика: учебник. 3-е изд., перераб. — СПб.: Изд-во С.-Петербур. ун-та, 2020. — 234 с.
- [Копотев, Янда, 2006] Копотев М. В., Янда Л. Национальный корпус русского языка // Вопросы языкознания. № 5. 2006. С. 149–155.
- [Копотев, 2014] Копотев, Михаил. Введение в корпусную лингвистику: Учебное пособие для студентов филологических и лингвистических специальностей университетов. Прага, 2014
- [Крижановская, Крижановский, 2020] Крижановская Н.Б., Крижановский А.А. От корпуса ВепКар к лингвистической платформе // Материалы Всероссийской научной конференции с международным участием «Бубриховские чтения: задокументированное народное слово» (Петрозаводск, 27–28 октября 2020 г.) [Электронный ресурс]. 2020. С. 157–159. URL: <http://mathem.krc.karelia.ru/publ.php?id=19787> (дата обращения: 07.06.2021).
- [ССТ] Словарь социолингвистических терминов. М., 2006.

[Плунгян, Сичинава, 2004] Плунгян В. А., Сичинава Д. В. Национальный корпус русского языка: опыт создания корпуса текстов современного русского языка // Труды международной конференции «Корпусная лингвистика-2004». — СПб.: Издательство Санкт-Петербургского университета, 2004. С. 216–238. URL: [https://events.spbu.ru/eventsContent/files/corpling/corpora2004/Sitchinava\\_art.pdf](https://events.spbu.ru/eventsContent/files/corpling/corpora2004/Sitchinava_art.pdf)