

23.10.2017

Корпус русских публицистических текстов второй половины 19 века

(<http://smalt.karelia.ru/corpus/index.phtml>)

Для проектной работы мы выбрали Корпус русских публицистических текстов второй половины 19 века Петрозаводского государственного университета. Так как корпус включает в себя произведения сравнительно небольшого промежутка времени, то он небольшой - всего 49 текстов. Он имеет два слоя разметки: морфологическую и синтаксическую.

Дизайн:

Оформление сайта очень спокойное и сдержанное, нет ничего лишнего, избыточного. Цветовая гамма гармоничная, не режет глаз. Главной страницей является страница поиска. Она поделена на две колонки. Слева перечислены произведения с указанием автором, к сожалению, нет даты создания. Их можно открыть и просмотреть. Большим плюсом является возможность выбрать определенные тексты для работы; можно выбрать одно, несколько, все, причем для этого есть специальное поле, где можно поставить галочку (она стоит по умолчанию). Справа находится столбец «Параметры поиска». Есть поиск по точной форме, по морфологическим признакам, по синтаксическим признакам, которые можно выбрать из предлагаемого списка.

[illegible]

Меню находится в правом верхнем углу, в том же блоке, что и название корпуса. Это неудобно, непривычно, но оно особо не нужно. Кнопка «Принципы создания» переводит нас в раздел, где есть ссылки на статьи по теме, но эти ссылки ведут на несуществующие страницы. Раздел «Помощь» номинально есть, но он пустой. Зато можно найти информацию о создателях корпуса, причем не только их имя-фамилию, но и род деятельности и адреса электронной почты. Также можно зайти в раздел с полезными статьями про корпусы и корпусную лингвистику, которые легко оттуда скачиваются в формате PDF. Кроме того, в меню есть кнопка «Авторизоваться», но непонятно, как это сделать, ведь нет возможности зарегистрироваться для начала.

Еще одним недостатком сайта является отсутствие привычного «подвала» внизу страницы. Из-за этого появляется ощущение незаконченности в разработке сайта.

Нельзя не отметить как достоинство адаптацию сайта при изменении масштаба или под мобильную версию. Это удобная и поэтому выигрышная деталь при взгляде на корпус со стороны.

В целом, сложно сказать, нравится ли «внешний вид» и устройство корпуса. С одной стороны он покоряет лаконичностью и простотой, с другой - нерабочие ссылки и неудобная разметка сайта сильно разочаровывают.

Onboarding:

Данный ресурс найти легко - это первая ссылка при поиске по названию корпуса. Более того можно перейти в этот корпус через ссылку в НКРЯ. Форму поиска искать не надо, так как, мы уже отмечали, она открывается как главная страница. Все поля поиска подписаны, поэтому ориентироваться, как и по каким параметрам искать, крайне просто. Можно с уверенностью сказать, что сайт интуитивно очевиден и понятен, даже несмотря на отсутствие примеров и образцов запросов.

Помощь пользователю:

Как уже было сказано выше, на сайте отсутствуют любые подсказки, life hacks и инструкции для начинающих пользователей. Несомненно, по этому критерию данный корпус проигрывает многим другим, но так как он понятен интуитивно, это небольшая проблема. Более серьезным минусом является отсутствие описания корпуса и его структуры, поэтому пользователю сложно получить целостное представление о данном ресурсе.

Примеры запросов:

Теперь мы переходим непосредственно к практическому использованию данного корпуса для лингвистических исследований, его функционалу и

возможностям. Надо сразу отметить, что в текстах используется старая орфография; в поисковое поле же надо вводить слово (это указано, что является достоинством) с современным написанием(корпус учитывает данный факт).

1) Для начала проверим корпус на поиск слов по маске. Введем маску «цы*» в поле для поиска по точной форме. Результат выглядит как список текстов, где встречается слово, начинающееся с «цы-». Указаны авторы и названия текстов. Нужное слово выделено цветом, при нажатии на него открывается окно с информацией о слове, то есть его категории, в том числе указано его устаревшее написание. Также можно просмотреть контекст искомого слова или же прочитать произведение целиком, что очень удобно.

Корпус русских публицистических текстов второй половины 19 века

Поиск | Принципы создания | Помощь | Авторы | Статьи | Авторизоваться

Список произведений

☒ Все тексты

☒ Акула // В. И. Даль, Время

☒ Безцветныя явленія // Федор Достоевский, Время

☒ В редакцію "Сына Отечества" // Федор Достоевский, Время

☒ Водопой // В. И. Даль, В. И. Даль

☒ Вопросъ объ университетахъ //

Параметры поиска

Поиск по точной форме (совр. написание)

цы*

НайтиОчистить

Поиск по морфологическим признакам

Корпус русских публицистических текстов второй половины 19 века

Поиск | Принципы создания | Помощь | Авторы | Статьи | Авторизоваться

Слово: "цы"

Результаты поиска

1. .: Гуляя въ праздники, матрость угостили себя и товарищей музыкой: они встрѣтили **цыгана** со скрипкой. [\[контекст\]\[статья\]](#) // В. И. Даль, В. И. Даль, Играй назадъ

2. Натѣшившись еволю, хозяинъ пирушки этой далъ **цыгану** гривну и пошелъ было съ товарищами своимъ путемъ, но скрипачъ нашелъ, что ему за труды мало гривны и сталъ неотступно просить еще хоть лятака. [\[контекст\]\[статья\]](#) // В. И. Даль, В. И. Даль, Играй назадъ

3. Одинъ камаринскій больше стоитъ, говорилъ **цыганъ** жалбно и настоятельно, а я камаринскаго сыграю вамъ разъ десятокъ. [\[контекст\]\[статья\]](#) // В. И. Даль, В. И. Даль, Играй назадъ

4. Участіе **цыгана** въ сватовствѣ Голопуленка. вопросъ, сдѣланный Хиверю попovichу: «не сломали ли вы шеи», и наконецъ самая свадьба, свернутая наскоро, на ярманѣ, — все это черты, «профанирующія» украинскій народъ съ его нравами и обычаями; отъ такой свадьбы, — свадьбы безъ честного корова, безъ дружокъ и свѣтилокъ, съ ужасомъ долженъ быть отвернуться добрый, учтивый и во всѣхъ отношеніяхъ прекраснй украинскій народъ. [\[контекст\]\[статья\]](#) // Время, Dubia, Критики-этнографы

5. Пушкинъ проводилъ ту же мысль въ своей поэмѣ «**Цыгане**», но Пушкинъ, какъ великій художникъ, выбралъ изъ среды кочующаго племени такіе идеальныя личности, что сравнительно съ образованнымъ Алеко они кажутся и человѣчнѣе, и даже глубже его въ пониманіи человѣческаго сердца. [\[контекст\]\[статья\]](#) // Время, Dubia, Письмо к редактору

6. Еслибы Алеко ужился между идеальными пушкинскими **цыганами**, онъ могъ бы еще быть счастливъ; онъ самъ нарушилъ это счастье, самъ убилъ свою свободу нарушая свободу другихъ. [\[контекст\]\[статья\]](#) // Время, Dubia, Письмо к редактору

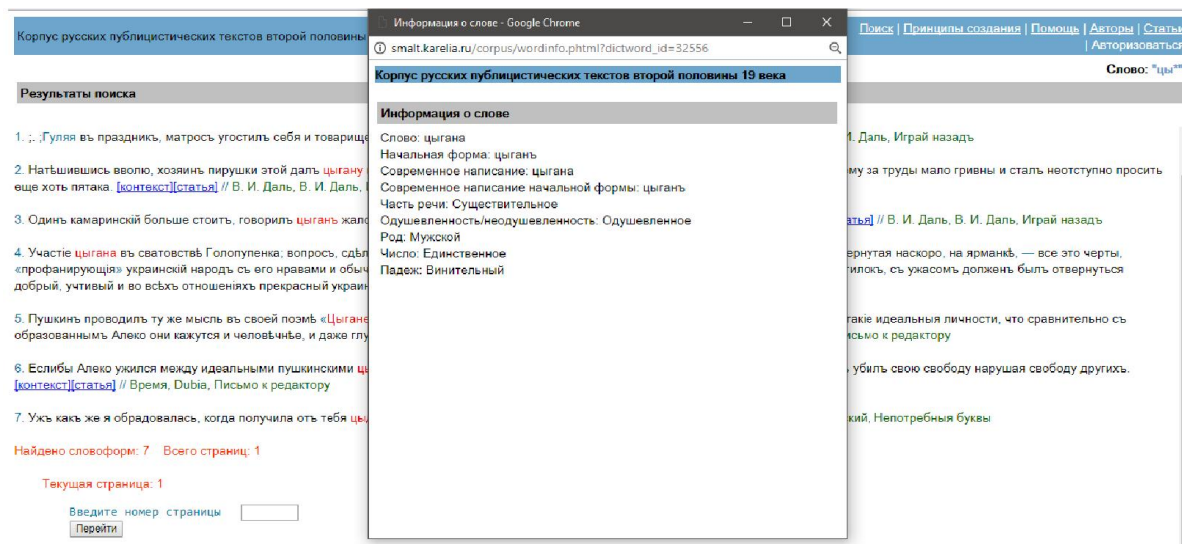
7. Ужъ какъ же я обрадовалась, когда получила отъ тебя **цыдулочку** - просто и расказать не умю. [\[контекст\]\[статья\]](#) // Время, Федор Достоевский, Непотребныя бувы

Найдено словоформ: 7 Всего страниц: 1

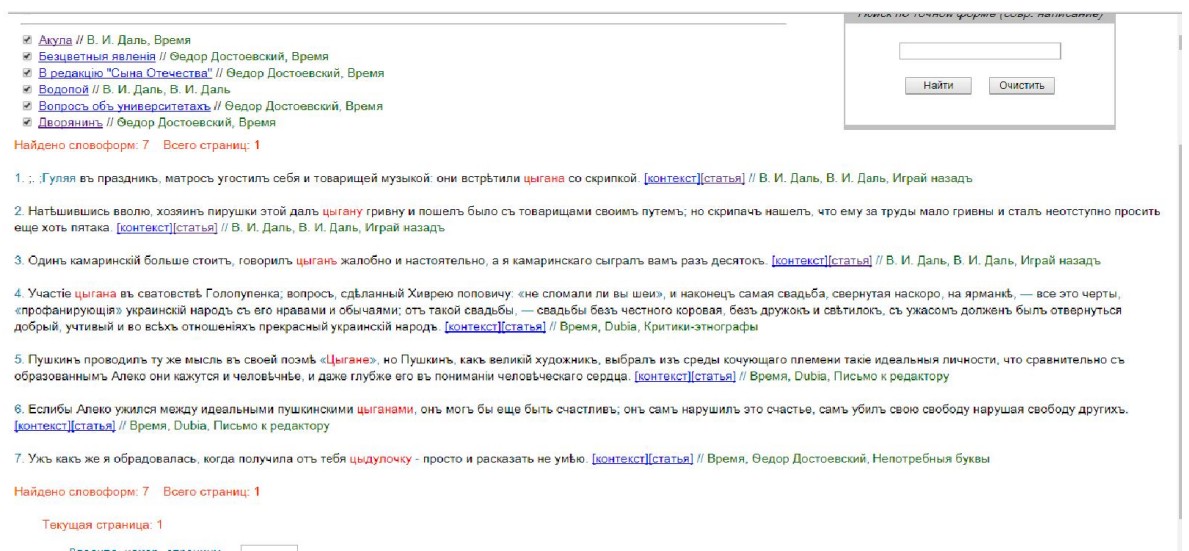
Текущая страница: 1

Введите номер страницы

Перейти



При нажатии на кнопку «Поиск» с намерением возвратиться на главную страницу видишь “интересную” картину: ты остаешься на той же странице с результатами выдачи, но поиск открывается на половину окна, то есть экран становится как бы поделенным горизонтальной линией- сверху поиск, снизу результат.



2) Теперь введем запрос сложнее. Посмотрим, сколько раз в корпусе встречается слово «человек» в форме предложного падежа единственного числа. Получаем, что данная словоформа встречается в корпусе 140 раз.

Федор Достоевский, Время

☒ [Водопой](#) // В. И. Даль, В. И. Даль

☒ [Вопросъ объ университетахъ](#) //

Федор Достоевский, Время

☒ [Дворянинъ](#) // Федор Достоевский,

Время

☒ [Девятнадцатый номер "Дня"](#) //

Владиславлев М.И., Время

☒ [Жуковский и романтизмъ](#) // Федор

Достоевский, Время

☒ [Журнальная заметка](#) // Федор

Достоевский, Время

☒ [Журнальные интересы](#) // Dubia,

Время

☒ [Закладъ](#) // В. И. Даль. В. И. Даль

Поиск по морфологическим признакам

Слово

человек

Морфологические признаки

S,Loc,Sg

?

Найти

Очистить

Поиск по синтаксическим признакам

Корпус русских публицистических текстов второй половины 19 века

[Поиск](#) | [Принципы создания](#) | [Помощь](#) | [Авторы](#) | [Статьи](#)

| [Авторизоваться](#)

Искомые параметры: "S,Loc,Sg" Слово: "человек"

Результаты поиска

Найдено словоформ: 140 Всего страниц: 15

1. Рассказываютъ, что мысль объ устройствѣ такой лечебницы (можетъ быть вызванная примѣромъ существующей въ Петербургѣ Максимилиановской лечебницы, — учрежденія крайне благотѣльнаго для бѣдныхъ больныхъ, приносящаго многимъ изъ нихъ, какъ намъ достовѣрно извѣстно, существенную пользу) давно ходила въ ярославскомъ обществѣ, но выражалась въ видѣ отдѣльныхъ желаній; теперь же осуществлена обществомъ ярославскихъ врачей, соединившихся, въ числѣ восемнадцати **человѣкъ**, на такое доброе дѣло, и совершившихся съ его помощію нѣкоторыхъ благотворителей. [\[контекст\]\[статья\]](#) // Время, Федор Достоевский, Мелочи
2. Въ извѣстіи упомянуто, что въ Ярославскую лечебницу съ 18 марта по 1 апрѣля приходящихъ больныхъ было сто тридцать **человѣкъ**. [\[контекст\]\[статья\]](#) // Время, Федор Достоевский, Мелочи
3. Кому не приходилъ въ голову мгновенный вопросъ: что дѣлать — пройти или остановится? кто это и какой это **человѣкъ**, — откуда и какъ пришолъ онъ на эту дорогу, темную, неровную, лежащую подъ удушливой и гнилой атмосферой? [\[контекст\]\[статья\]](#) // Время, Федор Достоевский, Мелочи
4. разумѣется, что **человѣчность**, выразившаяся въ такомъ кругу, должна была обратиться на животныхъ, хотя впрочемъ если бы она обратилась и на растенія, мы не удивились бы, потомучто и въ обращеніи съ растеніями можетъ **человѣкъ** являться то **человѣчнымъ**, то

Но если просматривать результаты, то мы видим, что выданные словоформы не соответствуют запросу. Самое интересное, что морфологическая разметка неподходящих нам слов в выдаче полностью правильная: например, в первом тексте словоформа «человек» размечена как форма генитива множественного числа. Потом мы меняем морфологические признаки того же самого слова, но результат получается таким же.

Искомые параметры: "S,Loc,Ssg" Слово: "человек"

Результаты поиска

Найдено словоформ: 140 Всего страниц: 15

1. Рассказываютъ, что мысль объ устройствѣ такой лечебницы (можетъ быть вызванная примѣромъ существующей въ Петербургѣ Максимилиановской лечебницы, — учрежденія крайне благотѣльнаго для бѣдныхъ больныхъ, приносящаго многимъ изъ нихъ, какъ намъ достовѣрно извѣстно, существенную пользу) давно ходила въ ярославскомъ обществѣ, но выражалась въ видѣ отдѣльныхъ желаній; теперь же осуществлена обществомъ ярославскихъ врачей, соединившихся, въ числѣ восемнадцати **человѣкъ**, на такое доброе дѣло, и совершившихся съ его помощію нѣкоторыхъ благотворителей. [\[контекст\]\[статья\]](#) // Время, Федор Достоевский, Мелочи
2. Въ извѣстіи упомянуто, что въ Ярославскую лечебницу съ 18 марта по 1 апрѣля приходящихъ больныхъ было сто тридцать **человѣкъ**. [\[контекст\]\[статья\]](#) // Время, Федор Достоевский, Мелочи
3. Кому не приходилъ въ голову мгновенный вопросъ: что дѣлать — пройти или остановится? кто это и какой это **человѣкъ**, — откуда и какъ пришолъ онъ на эту дорогу, темную, неровную, лежащую подъ удушливой и гнилой атмосферой? [\[контекст\]\[статья\]](#) // Время, Федор Достоевский, Мелочи

Из этого мы делаем вывод (проверив догадку еще на паре других слов), что поиск по морфологическим параметрам не работает, несмотря на морфологически правильную разметку слов.

3) Наконец, проверим, как работает поиск по синтаксическим признакам. Надо сказать, что под синтаксическими признаками подразумевается вхождение слова в некоторые конструкции, а не его роль в предложении. Хочется сразу заметить, что свод данных параметров кажется удобным, так как все признаки разделены на три категории: «Свободные структурные двухкомпонентные схемы», «Свободные структурные однокомпонентные схемы», «Фразеологизированные структурные схемы», - которые предлагают свои варианты синтаксических конструкций. При наведении на определенную синтаксическую схему появляется пример предложения с данной конструкцией, это very useful.

[Свободные структурные двухкомпонентные схемы](#) | [Свободные структурные однокомпонентные схемы](#) | [Фразеологизированные структурные схемы](#)

Свободные структурные двухкомпонентные схемы

Раздельнопредикативные схемы

- ☐ Подлежащно-сказуемые схемы
 - ☐ С координируемыми главными членами
 - ☐ N1 + Vf
 - ☐ N1 + (cop) + N1
 - ☐ N1 + Adj
 - ☐ N1 + part
 - ☐ С неkoordinируемыми главными членами
 - ☐ N1 + N2... (Adv)
 - ☐ N1 + Inf
 - ☐ N1 + (cop) + Praed
 - ☐ Inf + Praed(part)
 - ☐ Inf + (cop) + N1
 - ☐ Inf + Vf3s
 - ☐ Inf + Pron neg
 - ☐ Inf + Inf
- ☐ Не подлежащно-сказуемые схемы
 - ☐ N2 + (he) Vf3s
 - ☐ N2/N4 + (he) Praed(part)
 - ☐ N4 + Vf3s
 - ☐ N2 + N1quant(Adv quant)
 - ☐ N2 + Het
 - ☐ N3 + Vf3s
 - ☐ N3 + Praed
 - ☐ N2 + никого/ничего

Слитнопредикативные схемы

- ☐ Praed(part) N4/N2
- ☐ N1quant(Adv quant) N2
- ☐ Praed(part) Inf
- ☐ Pron neg Inf
- ☐ Het N2
- ☐ Vf3s N2
- ☐ Vf3s Inf
- ☐ Никого/Ничего N2
- ☐ Ни N2

Посмотрим, сколько результатов будет при выборе раздельнопредикативной схемы с неkoordinируемыми главными членами «Инфинитив + отрицательное местоимение» без ввода определенного слова.

Поиск по синтаксическим признакам

Слово

Синтаксические признаки

InfPronneg

?

Найти

Очистить

Корпус русских публицистических текстов второй половины 19 века

[Поиск](#) | [Принципы создания](#) | [Помощь](#) | [Авторы](#) | [Статьи](#) | [Авторизоваться](#)

Искомые параметры: "InfPronneg" Слово: ""

Результаты поиска

First: 1754, Last: 1779
Array ([0] => 1754 [1] => 1755 [2] => 1756 [3] => 1757 [4] => 1758 [5] => 1759 [6] => 1760)
1. В ней даже зацѣпиться не за что, потомучто все равно шатко, равно безосновательно, и не торчить ни одно мѣстечко, имѣющее какой-нибудь положительный смысл, положительное достоинство. [\[контекст\]\[статья\]](#) // Время, Федор Достоевский, Безцветные явления

First: 3507, Last: 3513
Array ([0] => 3507 [1] => 3508 [2] => 3509 [3] => 3510 [4] => 3511 [5] => 3512 [6] => 3513)
2. Говорить объ этихъ будущихъ журналахъ пока нечего. [\[контекст\]\[статья\]](#) // Время, Dubia, Журнальные интересы

First: 928, Last: 933
Array ([0] => 928 [1] => 929 [2] => 930 [3] => 931 [4] => 932 [5] => 933)
3. А ждательство уже нечего — все! [\[контекст\]\[статья\]](#) // Время, Федор Достоевский, Литературные антикварии

First: 119, Last: 163
Array ([0] => 160 [1] => 161 [2] => 162 [3] => 163)
4. Поэтому можно было полагать, что какъ ни много остается еще недосказаннаго о Гоголь относительно его человѣческой личности и поэтической дѣятельности, съ ея роковымъ концомъ, но это недосказанное никакъ не можетъ касаться его рассказовъ изъ малороссійскаго быта, о которыхъ, казалось, все рѣшено и перерѣшано уже нечего. [\[контекст\]\[статья\]](#) // Время, Dubia, Критики-этнографы

Всего 14 примеров, причем все они удовлетворяют параметрам поиска. Но сразу видно, что помимо самих примеров есть лишние строки «First» и «Array», которые являются доказательством недоработки сайта.

К сожалению, не на все синтаксические схемы есть примеры из корпуса. Так, там нет ни одной конструкции, построенной по фразеологизированной схеме.

Достоинства и недостатки:

Все достоинства и недостатки данного корпуса можно представить в виде таблицы:

Benefits	Drawbacks
Приятный дизайн	Сайт иногда не работает, «вылетает»

Страница поиска открывается по умолчанию	Незаконченная разметка сайта и неудобное расположение меню
Есть возможность выбрать произведения для работы	Много нерабочих ссылок и пустых разделов
Работает поиск по маске слова	Нет инструкций, примеров запросов
Полные списки морфологических и синтаксических признаков	Не работает поиск по морфологическим признакам
	Нет примеров фразеологизированных конструкций
	Невозможно искать сразу по морфологическим и синтаксическим параметрам
	В списке морфологических параметров можно выбрать несочетаемые признаки (например, возвратность-невозвратность глагола)
	Нет статистики употребления того или иного слова
	Нет метаразметки, нет возможности задавать подкорпус

Вывод:

В целом, корпус кажется незаконченным с точки зрения как разработки сайта, так и его состава. Кажущиеся на первый взгляд удобные методы поиска оказываются неэффективными, вдобавок сам поиск не всегда работает. Довольно внушительное количество недостатков корпуса и маленький объем не позволяют использовать его для лингвистических исследований.