

Национальный корпус польского языка

Narodowy Korpus Języka Polskiego

1. Данные о ресурсе

Национальный корпус польского языка (Narodowy Korpus Języka Polskiego) – это интернет-ресурс, основанный на собрании текстов на польском языке. Был создан в 2008 году. Делится на два механизма поиска: Поликарп и Пелкра. Первый работает на двух языках: польский и английский, однако, к сожалению, второй работает только на польском.

2. Дизайн и устройство корпуса

Шрифт, преимущественно используемый в корпусе, – Times New Roman, встречающийся повсеместно благодаря тому, что считается в некотором роде классическим.

Цветовая гамма сайта достаточно приятна на вид. Сочетание серого, белого и красного не раздражает взгляд, а белый и красный к тому же отсылают к цветам польского флага, что демонстрирует внимание разработчиков к деталям.

Интерфейс корпуса вполне удобен. Слева сгруппированы ссылки на материалы об истории корпуса, команде разработчиков, публикациях, благодарностях и тому подобных вещах, а также ссылка непосредственно для работы с корпусом, причем доступна она только в англоязычной версии – в польской по этому же адресу находятся сведения об использовании корпуса в различных проектах, что несколько сбивает с толку. Впрочем, вне зависимости от выбранного языка, ссылки для работы будут находиться справа.

Нарекания вызывает в первую очередь страница с подсказками, у которой отсутствует какое бы то ни было оформление – это просто текст. Впрочем, найти их достаточно легко, а это все-таки важнее. Ссылка на страницу находится слева от меню поиска, подсказки касаются в основном разметки, используемой в корпусе, и того, как следует оформлять запросы.

Версия сайта для мобильных устройств работает корректно, но вот десктопная подкачала: к сожалению, в ней не всегда работает вывод диакритических знаков.

3. Глазами новичка

Взгляд 1: Корпус довольно легко найти: достаточно выполнить запрос в любой поисковой системе. Можно заметить, что на сайте НКРЯ есть ссылка на корпус польского. Это делает поиск данного корпуса гораздо легче для русскоговорящих пользователей. Для англоговорящих пользователей поиск тоже достаточно прост. Сайт функционирует на двух языках: английском и польском. Форму поиска можно найти на всех страницах сайта, она находится в меню справа. На самой странице поиска под полем ввода есть специальные клавиши с польскими буквами, которые могут не присутствовать на клавиатуре пользователя. Слева находится небольшое меню с пунктами: поиск, настройки, сообщить об ошибке, помощь. В разделе “помощь” находится достаточно подробное описание работы и элементарные примеры запросов. Поиск осуществляется достаточно быстро, особенно для простых запросов.

Взгляд 2: Корпус симпатично выглядит и красиво оформлен. Возможен ввод диакритических символов при отсутствии таковых на клавиатуре. Можно работать и с мобильного телефона. Сам процесс несложен, однако работу с пелкррой сильно затрудняет то, что она только на польском. Результаты ищутся довольно быстро, но смена настроек иногда проходит затруднительно. Интересно наличие «слов дня» - самых частотных слов за последний день или неделю.

4. Продвинутый функционал (поиск по корпусу)

Ресурс позволяет осуществлять поиск словоформ, лексем, поиск по грамматическим категориям и метатекстовым данным. Семантическая, синтаксическая разметка отсутствует.

В строку ввода можно вбивать непосредственно словоформу (дополняя запрос спец. символами) или осуществлять поиск при помощи атрибутов:

- base задаёт поиск словоформ по лексеме
- orth задаёт поиск определённой словоформы или сегмента
- pos задаёт поиск словоформ по грамматическому значению
- Возможно использовать как атрибут название грамматической категории

attribute	possible values
number	sg pl
case	nom gen dat acc inst loc voc
gender	m1 m2 m3 f n
person	pri sec ter
degree	pos comp sup
aspect	imperf perf
negation	aff neg
accentability	akc nakc
post-prepositionality	npraep praep
accommodability	congr rec
agglutination	agl nagl
vocalicity	nwok wok
fullstoppedness	pun npun

Для создания сложных запросов можно использовать специальные символы и цифры:

?, *, +, ., ,, |, {, }, [,], (,)

- | - логическое ИЛИ. Результатом запроса "Ala|Ela" будет *Ala* или *Ela*,
- Внутри квадратных скобок заключены эквивалентные друг другу символы, то есть результат запроса [AE]la будет таким же, как и в предыдущем пункте
- Знак вопроса маркирует предшествующий символ как необязательный, то есть результатом запроса "beza?" будет и *bez*, и *beza*,
- Точка заменяет один символ (но не его отсутствие). Выдача по запросу "bez." включит *beza*, *bezy*, *bezq*, и т. д., но не *bez* или *bezami*,
- Астерикс означает, что в выдачу будут включены слова с любым количеством символов, предшествующих астериксу, в указанном месте. Например, по запросу "a*by" будут выданы варианты *by* (ноль *a*), *aby*, *aaaaby*, и т. д.,
- Сочетание .* заменяет любое количество любых символов. В выдачу по запросу "Ala.*" будут включены все слова, начинающиеся с *Ala*: *Ala*, *Alabama* и т. п.,
- Плюс по значению равен астериксу, но число заменяемых символов строго больше нуля.
- Выражение {n,m} означает от *n* до *m* повторений предшествующего символа. Результатом запроса "a{1,3}b.*" будут словоформы, содержащие от одной до трех *a* перед *b*: *aby*, *aaaby*, *absolutnie*, и т. д. Результатом запроса ".*(la){3,}.*" будут выражения, содержащие как минимум три *la* подряд: *tralalala*, *sialalala*,
- Сочетание /i в конце запроса определяет нечувствительность к регистру
- & - логическое И
- ! – отрицание

Poliqarp также предоставляет возможность поиска определённой фразы только *внутри* одного предложения или абзаца. Для этого используются определители within s и within p соответственно.

Помимо прочего, с помощью системы атрибутов можно задавать значения метатекстовых атрибутов, таких как имя автора, год создания, источник и др. Например, запрос “[pos=verb] meta channel=mowiony” выдаст все глаголы из текстов-записей устной речи.

Преимуществом поисковой системы является возможность дальнейшей офлайн работы с выдачей корпуса: результаты легко загружаются в excel или html.

Пример запроса:

Слово herbatniki (печенье)

The screenshot shows the PoliQarp search engine interface. The query is "herbatniki". The corpus is "the full NKJP corpus (1400M segments)". The search results are displayed as a list of 17 items, each showing a snippet of text and the corresponding word and its grammatical information.

Rank	Snippet	Word	Grammatical Information
1.	... ciasto maciorkę, czekoladowe	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
2.	... w tygodniadach. Trzeba wycofać	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
3.	... Możemy do tego wykorzystać	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
4.	... ciasta korzenne i ciasteczka	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
5.	... placki kukurydziane, a nawet	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
6.	... energię, podobnie jak czekoladowe	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
7.	... pomysłki mogły za niego kupić	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
8.	... miadzioł: pierszenie jabłko	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
9.	... , pieczywo cukierkowe trwałe	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
10.	... , pieczywo cukierkowe trwałe	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
11.	... herbatę do filiżanek i poprzyruci	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
12.	... szkody chce tylko przetrwać	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
13.	... w stanie wolnym niedźwiedź znajduje	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
14.	... magnezu. Czekoladowa słodycz Kakao	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
15.	... oraz napoje chłodzące a także	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
16.	... , pieczywo cukierkowe trwałe	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]
17.	... Wyciągnę z niego to co	herbatniki	[herbatnik: substantivus masculinus pluralis nominativus]

Пример запроса:

От одной до трёх букв k перед t

The screenshot shows the PoliQarp search engine interface. The query is "k{1,3}t". The corpus is "balanced NKJP subcorpus (300M segments)". The search results are displayed as a list of 13 items, each showing a snippet of text and the corresponding word and its grammatical information.

Rank	Snippet	Word	Grammatical Information
1.	... obietnic w rodzaju j a	k	[kolo: brevis-pun] t [tom: brevis-pun] y i k o i
2.	... Wyciągnąłem rękę, k	k	[kolo: brevis-pun] t [tom: brevis-pun] o i uderzał metalowym prętem
3.	... Jednak a	k	[kolo: brevis-pun] t [tom: brevis-pun] konwersji zawsze okazywał się procesem
4.	... p a n j a	k	[kolo: brevis-pun] t [tom: brevis-pun] o j e s t
5.	... No, dobrze, lecz k	k	[kolo: brevis-pun] t [tom: brevis-pun] o właściwie zawiał na krzyżu
6.	... nawet, że w ogóle k	k	[kolo: brevis-pun] t [tom: brevis-pun] o k o l w
7.	... I n s p e	k	[kolo: brevis-pun] t [tom: brevis-pun] o r k a?
8.	... żeby skopiować n i e	k	[kolo: brevis-pun] t [tom: brevis-pun] o r e, jak
9.	... żydowskim a zaraz po strzałach k	k	[kolo: brevis-pun] t [tom: brevis-pun] o i przeszukuje legowiska w
10.	... Literackie" 1939 D o	k	[kolo: brevis-pun] t [tom: brevis-pun] o r H a c
11.	... św. Zyty, A k	k	[kolo: brevis-pun] t [tom: brevis-pun] o r z y c
12.	... a ś c n a	k	[kolo: brevis-pun] t [tom: brevis-pun] o r y m g
13.	... - to musiał być naprawdę k	k	[kolo: brevis-pun] t [tom: brevis-pun] o i. "Dobrze

Пример запроса:

Все словоформы czerwony (красный), за которыми следует существительное единственного числа не в дательном падеже, из текстов, созданных позднее 1980 года
[base=czerwony][pos=subst & number=sg & case!=dat] meta created>1980

Query: [base=czerwony][pos=subst & number=sg & case!=dat] meta created>1980
a e s l u o z z A C E L N O S Z Z
Corpus: the full NKJP corpus (1800M segments)
Search

Results

Found 350 results
Displaying results 1—25

Next 25

1.	-moralistę, który nosił	czerwoną legitymację	komunistycznej partii, o to
2.	, a jest to Polska	Czerwony Krzyż	, którego pracami koczują od
3.	, a straż pożarna i	Czerwony Krzyż	wypełniają tylko lukę w ratownictwie
4.	, a także w krajowej	czerwonej księdze	roślin (Zarzycki, Kaźmierczakowa
5.	, amia rozrywana jak postać	czerwonego sukna	między walczących między sobą polityków
6.	, Caritas Polska, Polski	Czerwony Krzyż	, Polski Komitet Pomocy Społecznej
7.	, czarne, białe lub	czerwone porzeczki	i agrest:-
8.	, Drygały, Nowogród,	Czerwony Dwór	, Giżycko i Borki.
9.	, Drygały, Nowogród,	Czerwony Dwór	, Giżycko i Borki.
10.	, jak kto woli,	czerwona księga	opisująca praktyki rozszarpywania środków finansowych
11.	, jakkolwiek zaproponowała	czerwoną plamkę	" w formie nagrzanego tekstu
12.	, która dotyczy finansowania	Czerwonego Krzyża	, pragnę poinformować Wysoką Izbę
13.	, który uwolnił nas od	czerwonej bandy	, przepraszam bardzo (Okłaski
14.	, maślna, czarnej i	czerwonej porzeczki	, wiśni i jabłek.
15.	, maślane, czarne i	czerwone porzeczki	, wiśnie mrożone bez cukru