

Национальный корпус польского языка

Narodowy Korpus Języka Polskiego

Маслова Мария (mashamaslova85@gmail.com)

Гребёнкина Мария (ya.grebenkina-maria-74@ya.ru)

Кошкарёва Диана (di.coshckarewa@yandex.ru)

Миронов Михаил (orderlyoftime@gmail.com)

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 пелкрой 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*,
- Астериск означает, что в выдачу будут включены слова с любым количеством символов, предшествующих астерису, 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. н.,
- Плюс по значению равен астерису, 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 interface. At the top, there is a search bar with the query 'herbatniki'. Below it, the corpus is specified as 'the full NKJP corpus (1400M segments)'. The search results are displayed in a table with 17 rows. Each row contains a number, a snippet of text, the word 'herbatniki', its grammatical information, and a source reference.

№	Snippet	herbatniki	Grammatical info	Source
1.	... ciasto mączne, czekoladowe	herbatniki	herbatnik sobot.pl nom m3	Bambo. Trzasko nazywał to
2.	w węglowodany. Trzeba wycofać	herbatniki	herbatnik sobot.pl nom m3	biażkopy, boleczki.
3.	Możemy do tego wykorzystać	herbatniki	herbatnik sobot.pl nom m3	placki czekoladowe, lizaki
4.	ciasta korzenne i ciasteczka	herbatniki	herbatnik sobot.pl nom m3	ciasteczka z migdałami i
5.	placki kukurydziane, a nawet	herbatniki	herbatnik sobot.pl nom m3	Zywność pochodzi z programu
6.	energii, podobnie jak czekoladowe	herbatniki	herbatnik sobot.pl nom m3	lub wafelki (niestety rozpuszczają
7.	prémium mogły za niego kupić	herbatniki	herbatnik sobot.pl nom m3	małame, kanapki albo owoce
8.	miadzio: pieczono jabłko	herbatniki	herbatnik sobot.pl nom m3	Ostatni: zupa pomidorowa z
9.	... pieczywo ciakierne twarde i	herbatniki	herbatnik sobot.pl nom m3	wafle), suszone
10.	... pieczywo ciakierne twarde i	herbatniki	herbatnik sobot.pl nom m3	wafle), suszone
11.	herbatki do filiżanki i poprzyjmu	herbatniki	herbatnik sobot.pl nom m3	z porcelanowego talerzyka. Tak
12.	szkoly chce tylko przetrwać i	herbatniki	herbatnik sobot.pl nom m3	– powiedziała nam jedna z
13.	w stanie wolnym mediewiedzi najdale	herbatniki	herbatnik sobot.pl nom m3	albo ciakierki? O: używacie
14.	magnezu Czekoladowa słodycz Kakao	herbatniki	herbatnik sobot.pl nom m3	JA) Petki Dzien dobry
15.	oraz napoje chłodzące a także	herbatniki	herbatnik sobot.pl nom m3	słodkimi i słodkimi
16.	... pieczywo ciakierne twarde i	herbatniki	herbatnik sobot.pl nom m3	wafle), suszone
17.	Wycożać z niego to nie	herbatniki	herbatnik sobot.pl nom m3	i wieszka chorow do reki

Пример запроса:

От одной до трёх букв *k* перед *t*

Poliqarp search engine for NKJP data

QUERY
SETTINGS
FILE A BUG
HELP

Query:

Corpus:

Results

Found 27 results
Displaying results 1—27

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

Query:

Corpus:

Results

Found 350 results
Displaying results 1—25

Next 25

1.	-moralistę, który nosił	<u>czernoną legitymację</u>	komunistycznej partii, o to
2.	, a jest to Polski	<u>Czerwony Krzyż</u>	, którego pracami kieruję od
3.	, a straż pożarna i	<u>Czerwony Krzyż</u>	wypełniają tylko lukę w ratownictwie
4.	, a także w krajowej	<u>czernonej księdze</u>	roslin (Zarzycki, Kaźmierczakowa
5.	, amia rozrywana jak postaw	<u>czernonego sukna</u>	między walczących między sobą polityków
6.	, Caritas Polska, Polski	<u>Czerwony Krzyż</u>	, Polski Komitet Pomocy Społecznej
7.	, czarne, białe lub	<u>czernone porzeczki</u>	i agrest:- -
8.	, Drygaly, Nowogród,	<u>Czerwony Dwór</u>	, Giżycko i Borki.
9.	, Drygaly, Nowogród,	<u>Czerwony Dwór</u>	, Giżycko i Borki.
10.	, jak kto woli,	<u>czernona księga</u>	opisująca praktyki rozszarpywania środków finansowych
11.	, jakkolwiek zaproponowała dłuższą	<u>czernoną plamkę</u>	" w formie nagranego tekstu
12.	, która dotyczy finansowania Polskiego	<u>Czerwonego Krzyża</u>	, pragnę poinformować Wysoką Izbę
13.	, który uwolnił nas od	<u>czernonej bandy</u>	, przepraszam bardzo (Okłaski
14.	, malin, czarnej i	<u>czernonej porzeczki</u>	, wiśni i jabłek.
15.	, maliny, czarnej i	<u>czernone porzeczki</u>	, wiśnie urozone bez cukru