

# Korpusová lingvistika – nástroje a jejich využití



Počítačové nástroje pro češtinu  
Jaro 2020

Jakub Machura

415795@mail.muni.cz

# Co nás čeká?

## 1. Nástroje běžící na pražské platformě [korpus.cz](https://korpus.cz)

- WaG
- SyD
- Morfio
- Kwords
- Treq

# Co nás čeká?

## 2. Brněnský korpusový manažer **Sketch Engine**

- Word Sketch
- Tezaurus
- N-gramy
- Skell

# WaG – Slovo v kostce

- aplikace poskytuje rychlý a základní přehled o tom, jak se v korpusu zadané slovo používá
- data získaná z psaného, mluveného i paralelního korpusu
- funkce: *Vyhledat slovo*  
*Hledat ve dvou jazycích / přeložit*
- vyzkoušejte sami na <https://www.korpus.cz/slovo-v-kostce/>

# SyD

- nástroj pro korpusový průzkum variant
- viz <https://syd.korpus.cz/>
- Sy – varianty v současném jazyce (synchronní část)
  - vyzkoušejte např. *současně* × *najednou* × *naráz*
- D – varianty v průběhu historie (diachronní část)
  - vyzkoušejte *beruška* × *slunéčko*

# Morfio

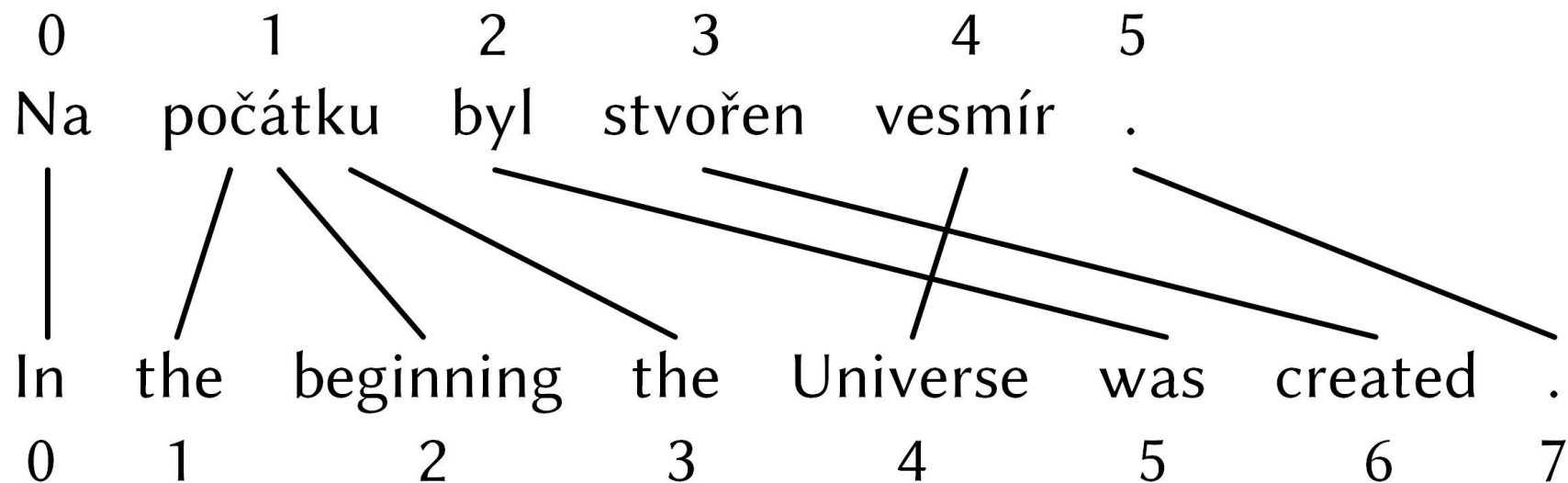
- nástroj pro slovotvornou analýzu
- identifikuje dvojice slov spjaté formálními slovotvornými vztahy
- více informací v hodině věnované slovotvorbě
- <https://morfio.korpus.cz/>

# KWords

- nástroj pro identifikaci klíčových slov
- součást projektu *A Needle in a Haystack*
- zkoumaný text je porovnáván s korpusem (referenční text)
- viz popis aplikace na <https://kwords.korpus.cz/>

# Treq

- databáze překladových ekvivalentů
- vytvořeno automaticky na základě dat z paralelního korpusu *InterCorp* (viz <https://intercorp.korpus.cz/>)
- snaha zarovnat slovo na slovo:





# Treq

- zadávání víceslovných výrazů
- možnost využití regulárních výrazů
- nástroj využívaný nejen překladateli
- až 39 jazyků
- vyzkoušejte nástroj na <https://treq.korpus.cz/>

# SKETCH ENGINE

- komerční korpusový manažer
- firma [Lexical Computing](#), zakladatel [Adam Kilgarriff](#)
- 500 korpusů (převážně webových) ve více než 90 jazycích
- <https://www.sketchengine.eu/>



# SKETCH ENGINE

## Přihlášení

- <https://auth.sketchengine.eu/#login>

# Log in

LOG IN

[Forgot password?](#)

[Need help logging in?](#)

or



Institutional login

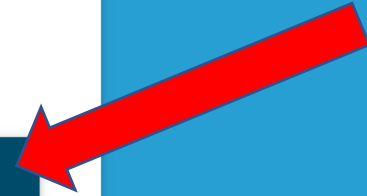


Sign in with Google

Don't have an account? [Sign up](#)

Or try [open corpora](#).

CLICK ON THAT



The Sketch Engine is a Corpus Query System allowing you to research how words behave.

## Which organisation would you like to sign in with?

Start typing the name of your organisation (e.g. Anywhere College) in the search box, and options will appear below:



**Continue**

or [Let me choose from a list](#)

 **UČO / GUEST ID** **SEKUNDÁRNÍ HESLO****PŘIHLÁSIT**

Tato služba vyžaduje ověření Vaší identity  
(**UČO / Guest ID**)

Pokud neznáte **sekundární heslo**,  
můžete si je nastavit prostřednictvím IS  
MU na stránce [změna hesla](#).

V případě problémů či dotazů kontaktujte prosím [helpdesk@ics.muni.cz](mailto:helpdesk@ics.muni.cz).



## CZECH WEB 2017 (CSTENTEN17)

INFO O KORPUSU

SPRAVOVAT KORPUS



Word Sketch

Kolokace a kombinace slov



Word sketch rozdíly

Porovnat kolokace dvou slov



Tezaurus

Synonyma a podobná slova



Konkordance

Příklady použití v kontextu



Paralelní konkordance

Hledání překladu



Seznam slov

Frekvenční seznam



N-gramy

Víceslovné výrazy (MWEs)



Klíčová slova

Extrakce terminologie



Trendy

Diachronická analýza, neologismy



Slovník jedním kliknutím

Návrh automatického slovníku

## NEDÁVNÉ KORPUSY

NOVÝ KORPUS

Czech Web 2017  
(csTenTen17)

Czech

10 502 222 474



SOUČASNÉ VÝSLEDKY

OBLÍBENÉ VÝSLEDKY

ZNAČKOVÁNÍ

hledat psaním





ZÁKLADNÍ

POKROČILÉ

MÉ KORPUSY

SDÍLENÉ SE MNOU

## JAZYKY

Vyberte jazyk a my pro vás zvolíme nejlepší korpus.

ARABIC

DANISH

DUTCH

ENGLISH

FRENCH

GERMAN

CHINESE

ITALIAN

JAPANESE

KOREAN

POLISH

PORTUGUESE

RUSSIAN

SPANISH

Více jazyků

hledat psaním



## RYCHLÝ ÚVODNÍ TUTORIÁL



How to start with Sketch Engine



Přehrát po...



Sdílet

# How to start



## in 2 minutes



Watch the tutorial



# Více info na <https://www.sketchengine.eu/guide/>

LOG IN Sign up **FREE Trial** subscribe to news



Home News & Events Pricing **Guide** About us Contact

## Sketch Engine User Guide

video lessons

face-to-face training

Glossary

### general screens

dashboard

select corpus

corpus info page

manage corpus

settings

### corpus tools

**nástroje**

word sketch

sketch difference

thesaurus

concordance

parallel concordance

wordlist

n-grams

keywords & terms

bilingual term extraction

trends

OneClick Dictionary

### dashboard

Click a place in the dashboard to get help.

The screenshot shows the Sketch Engine dashboard for the 'English Web 2015 (enTenTen15)' corpus. It features a grid of tool tiles, each with a numbered callout: 1. Word Sketch, 2. Word Sketch Differences, 3. Thesaurus, 4. Concordance, 5. Parallel Concordance, 6. Wordlist, 7. N-grams, 8. Keywords, 9. Trends, 10. One-Click Dictionary, 11. Manage Corpus, 12. Corpus Info, 13-17. Recently Used Corpora, 18. Recently Used Results, 19. Favourite Results, 20. Search results. A sidebar on the left contains navigation icons. A banner at the bottom right says 'Master the interface in 2 days!' with the location 'Prague Airport'.

### Corpus Building & download

build a corpus

### Corpora, languages, tagsets

[list of languages](#)

### Account, subscription, payment

[My account and password](#)

# Konkordance

- použití daného slova v kontextu
- slovní tvar, lemma, fráze, CQL



ZÁKLADNÍ

**POKROČILÉ**

JAK TO FUNGUJE

Typ dotazu ?

základní

**lemma**

fráze

slovo

znak

CQL

Slovní druh

jakýkoli

**podstatné jméno**

přídavné jméno

zájmeno

numeral

sloveso

přísluvce

nředložka

lemma

koronavirus

✓ A = a ?

Subkorpus ?

žádný (celý korpus)



Filtrovat kontext ? ▾

Typy textů ? ▾

OK



lemma **koronavirus** 258 (0,02 v mil.)



KWIC



Detaily

Levý kontext

KWIC

Pravý kontext

1	<input type="checkbox"/>	wikipedia.org	í původci nemoci napadají jen jisté zvířecí druhy. </s></s>	<b>Koronavirus</b> k1gInSc1/Koronavirus	přenosné gastroenteritidy (TGEV) se objevuje u prasat a z
2	<input type="checkbox"/>	wikipedia.org	nes, holubů, a japonských křepelek. </doc></s></s></doc>	<b>Koronavirus</b> k1gInSc1/Koronavirus	(Coronavirus) je společné označení pro čtyři rody virů obs
3	<input type="checkbox"/>	wikipedia.org	ěsnila, že původní údaje o počtu osob zemřelých na tento	<b>koronavirus</b> k1gInSc4/koronavirus	byly podhodnocené, a že ve skutečnosti zemřelo 282 lidí (
4	<input type="checkbox"/>	wikipedia.org	nc Felipe.[6] </s></s> Saúdskoarabské úřady oznámily, že	<b>koronavirus</b> k1gInSc1/koronavirus	MERS si v zemi vyžádal již 282 obětí. [7] </s></s> Za bojko
5	<input type="checkbox"/>	tribune.cz	oho u 22 pikornavirus, u sedmi metapneumovirus, u osmi	<b>koronavirus</b> k1gInSc1/koronavirus	, u dvou influenza virus typu A/B a parainfluenza virus a u
6	<input type="checkbox"/>	zet.cz	vyžádala 46 životů, nakaženo bylo celkem 94 lidí. </s></s>	<b>Koronavirus</b> k1gInSc1/Koronavirus	MERS může vyvolat kašel, vážné dýchací potíže, horečku
7	<input type="checkbox"/>	zet.cz	otilátky, o kterých si myslíme, že jsou specifické právě pro	<b>koronavirus</b> k1gInSc4/koronavirus	MERS nebo jemu podobný virus, v krvi jednohrbých velblc
8	<input type="checkbox"/>	tribune.cz	ně </s></s> V Česku se objevil druhý případ podezření na	<b>koronavirus</b> k1gInSc4/koronavirus	MERS. </s></s> Muž z Ostravy byl hospitalizován na infek
9	<input type="checkbox"/>	tribune.cz	> V České republice se objevil druhý případ podezření na	<b>koronavirus</b> k1gInSc4/koronavirus	MERS. </s></s> Muž z Ostravy byl hospitalizován na infek
10	<input type="checkbox"/>	tribune.cz	zdravotnický systém ji zvládne. </s></s> Infekci způsobuje	<b>koronavirus</b> k1gInSc1/koronavirus	z téhož kmene, jako je virus zodpovědný za onemocnění s
11	<input type="checkbox"/>	euro.cz	íkům otvírá už v deset hodin. </doc></s></s></doc> WHO:	<b>Koronavirus</b> k1gMnSc1/Koronavirus	MERS zatím není celosvětovou hrozbou </s></s> Přestože
12	<input type="checkbox"/>	euro.cz	dýchacích obtíží i selhání ledvin. </s></s> Celosvětově na	<b>koronavirus</b> k1gInSc4/koronavirus	MERS zemřelo už 152 lidí, naprostá většina v Saúdské Ar
13	<input type="checkbox"/>	aktualne.cz	sobuje symptomy podobné nákaze virem SARS. </s></s>	<b>Koronavirus</b> k1gMnSc1/Koronavirus	je označení pro čtyři rody virů, které způsobují méně záva
14	<input type="checkbox"/>	brnenskadrba.cz	áb </s></s> V Česku se objevil druhý případ podezření na	<b>koronavirus</b> k1gInSc4/koronavirus	MERS. </s></s> Třiatřicetiletý pracovník brněnské cestovn
15	<input type="checkbox"/>	euro.cz	zhlédnete. kontaktujte policii. </doc></s></s></doc> Nový	<b>koronavirus</b>	je dle WHO přenosný kontaktem s nakaženou osobou </s>



lemma **koronavirus** 258 (0,02 v mil.)



KWIC



Detaily

Levý kontext

KWIC

Pravý kontext

	Levý kontext	KWIC	Pravý kontext
1	wikipedia.org í původci nemoci napadají jen jisté zvířecí druhy.	<b>Koronavirus</b> k1glnSc1/Koronavirus	přenosné gastroenteritidy (TGEV) se objevuje u prasat a z
2	wikipedia.org nes, holubů, a japonských křepelek.	<b>Koronavirus</b> k1glnSc1/Koronavirus	(Coronavirus) je společné označení pro čtyři rody virů obs
3	wikipedia.org ěsnila, že původní údaje o počtu osob zemřelých na tento	<b>koronavirus</b> k1glnSc4/koronavirus	byly podhodnocené, a že ve skutečnosti zemřelo 282 lidí (
4	wikipedia.org nc Felipe.[6] Saúdskoarabské úřady oznámily, že	<b>koronavirus</b> k1glnSc1/koronavirus	MERS si v zemi vyžádal již 282 obětí. [7] Za bojko
5	tribune.cz oho u 22 pikornavirus, u sedmi metapneumovirus, u osmi	<b>koronavirus</b> k1glnSc1/koronavirus	, u dvou influenza virus typu A/B a parainfluenza virus a u
6	zet.cz		otíže, horečku
7	zet.cz		ohrbých velblc
8	tribun		zován na infek
9	tribun		zován na infek
10	tribun		onemocnění :
11	euro.		>> Přestože
12	euro.		v Saúdské Ar
13	aktua		ují méně záva:
14	brnen		ěnské cestovn
15	euro.		ou osobou </s:

virus SARS, se začala šířit z Blízkého východu a od loňska si vyžádala 46 životů, nakaženo bylo celkem 94 lidí. </s></p><p><s> Choroba, která připomíná virus SARS, se začala šířit z Blízkého východu a od loňska si vyžádala 46 životů, nakaženo bylo celkem 94 lidí. </s> <s> **Koronavirus** MERS může vyvolat kašel, vážné dýchací potíže, horečku, zápal plic a údajně i selhání ledvin. </s><s> Úmrtnost nakažených pacientů přesahuje 50 procent, což je mnohem víc než v případě příbuzného koronaviru SARS. </s></p><p><s> Důkazy pro přenos nemoci mezi lidmi existovaly, odborníci se ale vždy domnívali

# Word Sketch

- tzv. slovní profily
- sdružují kolokace slov na základě gramatických relací (podmět, přísudek, atribut, ...)
- hodnoty frekvence a skóre



ZÁKLADNÍ

**POKROČILÉ**

JAKO SEZNAM

JAK TO FUNGUJE

Hledání ?

krásný

Subkorpus ?

žádný (celý korpus)



Slovní druh ?

**auto**

adjective

adverb

noun

pronoun

verb

numeral

preposition

Minimální frekvence ?

auto

Minimální skóre ?

0

Přeložit

Typy textů ?

OK



krásný as přídavné jméno 3 354 950x



modifiers of "krásný"			
<b>moc</b>	39 843	9,1	...
moc krásné			
<b>opravdu</b>	26 054	8,27	...
opravdu krásné			
<b>nejen</b>	7 957	7,89	...
nejen krásné			
<b>uvěřitelně</b>	2 841	7,59	...
neuvěřitelně krásné			
<b>tak</b>	36 884	7,34	...
tak krásné			
<b>opět</b>	4 082	7,18	...
opět krásné			
<b>přirozeně</b>	1 555	7,07	...
přirozeně krásná			
<b>skutečně</b>	5 535	7,05	...
skutečně krásné			
<b>strašně</b>	2 501	7,03	...
strašně krásný			
<b>plno</b>	1 331	6,97	...
plno krásných			
<b>tu</b>	5 362	6,93	...
tu krásné			
<b>spoustu</b>	1 654	6,93	...
spoustu krásných			

nouns modified by "krásný"			
<b>výhled</b>	64 223	9,53	...
s krásným výhledem			
<b>příroda</b>	42 731	8,96	...
krásné přírodě			
<b>počasí</b>	43 816	8,85	...
krásné počasí			
<b>den</b>	108 899	8,77	...
krásný den			
<b>pláž</b>	33 314	8,56	...
krásné pláže			
<b>žena</b>	33 172	8,27	...
krásné ženy			
<b>prostředí</b>	36 253	7,94	...
v krásném prostředí			
<b>fotka</b>	20 955	7,89	...
krásné fotky			
<b>zážitek</b>	24 066	7,86	...
krásný zážitek			
<b>místo</b>	38 407	7,63	...
krásné místo			
<b>dívka</b>	17 267	7,63	...
krásné dívky			
<b>dárek</b>	19 731	7,6	...
krásný dárek			

"krásný" and/or ...			
<b>zdravý</b>	8 471	9,46	...
krásné a zdravé			
<b>klidný</b>	6 465	9,22	...
krásné a klidné			
<b>pohodový</b>	4 040	8,89	...
krásné a pohodové			
<b>zajímavý</b>	5 516	8,64	...
krásné a zajímavé			
<b>mladý</b>	4 399	8,51	...
mladá a krásná			
<b>čistý</b>	3 742	8,3	...
krásné a čisté			
<b>milý</b>	2 486	7,89	...
krásné a milé			
<b>originální</b>	2 238	7,82	...
krásné a originální			
<b>bohatý</b>	2 136	7,68	...
krásné a bohaté			
<b>elegantní</b>	1 924	7,65	...
krásné a elegantní			
<b>kvalitní</b>	2 999	7,59	...
krásné a kvalitní			
<b>veliký</b>	2 862	7,52	...
největší a nejkrásnější			



... is "krásný"			
<b>šťěstí</b>	2 323	9,53	...
šťěstí je krásná věc			
<b>život</b>	3 072	9,46	...
Život je krásný			
<b>fotka</b>	1 048	8,2	...
fotky jsou krásné			
<b>počasí</b>	1 020	8,16	...
Počasí bylo krásné			
<b>žena</b>	1 038	8,08	...
žena je krásná			
<b>svět</b>	1 107	7,92	...
svět je krásný			
<b>láska</b>	723	7,74	...
Láska je krásná			
<b>okolí</b>	727	7,73	...
V okolí je krásná			
<b>den</b>	1 014	7,62	...
den je krásný			
<b>pláž</b>	572	7,39	...
Pláž je krásná			
<b>hotel</b>	619	7,27	...
Hotel je krásný			
<b>Praha</b>	537	7,02	...
Praha je krásná			

words before "krásný"			
<b>fakt</b>	2 034	7,96	...
je fakt krásná			
<b>vskutku</b>	714	6,94	...
vskutku krásné			
<b>bezesporu</b>	587	6,79	...
bezesporu nejkrásnější			
<b>prostě</b>	1 901	5,59	...
prostě krásná			
<b>opravdu</b>	205	5,51	...
opravdu krásný			
<b>přece</b>	1 149	5,46	...
přece krásné			
<b>mimochodem</b>	236	5,41	...
mimochodem krásné			
<b>prý</b>	1 029	5,34	...
prý nejkrásnější			
<b>zase</b>	1 930	5,15	...
zase krásný			
<b>také</b>	9 062	5,1	...
také krásné			
<b>jistě</b>	1 022	5,09	...
jistě krásné			
<b>ra</b>	440	5,08	...
book of ra nejkrásnější			

words after "krásný"			
<b>ba</b>	34	5,01	...
krásný ba			
<b>jednou</b>	48	4,79	...
krásné jednou			
<b>ch</b>	18	4,17	...
z nejkrásnější ch			
<b>akorát</b>	14	4,17	...
krásné akorát			
<b>ach</b>	14	3,61	...
všechno zdá se krásnější Ach			
<b>pr</b>	45	3,59	...
krásné PR			
<b>jaký</b>	19	3,57	...
nejkrásnější jaký jsem			
<b>právě</b>	440	3,51	...
krásná právě			
<b>co</b>	843	3,49	...
to nejkrásnější co			
<b>kde</b>	15	2,93	...
Krásný Kde			
<b>ahoj</b>	23	2,85	...
krásné ahoj			
<b>bezesporu</b>	13	2,72	...
nejkrásnější bezesporu patří			



# Tezaurus

- The thesaurus in Sketch Engine is an automatically generated list of synonyms or words belonging to the same category (semantic field). The list is produced based on the context in which the words appear in the selected corpus. Only nouns, adjectives, verbs and adverbs are supported in most corpora



ZÁKLADNÍ

**POKROČILÉ**

JAK TO FUNGUJE

Hledání ?

krásný



Slovní druh ?

auto

**adjective**

adverb

noun

pronoun

verb

numeral

preposition

Shlukovat podobné položky ?

Maximum položek tezauru ?

1000

Minimální skóre tezauru ?

0

OK



krásný as přídavné jméno 3 354 950×



	Word	Frekvence ?	
1	nádherný	1 063 308	...
2	hezký	1 286 906	...
3	pěkný	1 677 821	...
4	úžasný	1 362 352	...
5	zajímavý	4 050 171	...
6	skvělý	3 022 707	...
7	příjemný	2 802 641	...
8	plný	3 222 664	...
9	dobrý	20 585 983	...
10	překrásný	158 798	...

Řádků na stránku

10



1-10 z 1 000



1

/ 100



# n-gramy

- sled po sobě jdoucích položek z dané posloupnosti
- unigram, bigram, trigram...
- kolokace vs. n-gram

The office building was demolished yesterday.

**5 bigramů:** *the office, office building, building was, was demolished, demolished yesterday*

**2 kolokace:** *office building, to demolish a building*



ZÁKLADNÍ

**POKROČILÉ**

JAK TO FUNGUJE

Délka n-gramu ?

**2** 3 4 5 6

Atribut ?

slovo

Frekvence min. ?

5

Frekvence max. ?

0

Zanořené n-gramy ?

Zahrnout neslovní výrazy

A = a ?

mimo tato slova: ?

Subkorpus ?

žádný (celý korpus)



Klíčové n-gramy

Additional criteria ?

vše

začínající písmeny

končící písmeny

**obsahující písmena**

začínající slovem

obsahující slovo

končící slovem

odpovídá regulárnímu výrazu

černý



**PŘIDAT TUTO PODMÍNKU**

Typy textů ?

OK



# ! 2–3-grams slovo (položky: 1 630, celková frekvence: 33 088)



	Word	↓ Počet ?
1	a černý	1 198 ...
2	černý čaj	904 ...
3	je černý	819 ...
4	černý a	644 ...
5	černý humor	535 ...
6	na černý	447 ...
7	černý trh	422 ...
8	černý rybíz	401 ...
9	černý pepř	309 ...
10	jako černý	297 ...
11	černý toner	280 ...
12	černý kašel	275 ...
13	má černý	255 ...

	Word	↓ Počet ?
14	nebo černý	232 ...
15	černý s	219 ...
16	černý nebo	218 ...
17	černý kůň	209 ...
18	černý je	191 ...
19	černý bez	183 ...
20	ten černý	175 ...
21	i černý	171 ...
22	černý pruh	164 ...
23	se černý	161 ...
24	velký černý	155 ...
25	čáp černý	154 ...
26	černý se	154 ...

	Word	↓ Počet ?
27	byl černý	149 ...
28	že černý	141 ...
29	černý dým	137 ...
30	černý pátek	135 ...
31	celý černý	134 ...
32	pro černý	128 ...
33	černý kouř	125 ...
34	černý na	122 ...
35	černý pasažér	119 ...
36	černý plech	116 ...
37	černý povrch	110 ...
38	černý pes	106 ...
39	černý lak	106 ...

	Word	↓ Počet ?
40	jen černý	105 ...
41	bez černý	103 ...
42	až černý	103 ...
43	to černý	98 ...
44	černý plast	97 ...
45	datel černý	97 ...
46	černý kontinent	93 ...
47	černý trh s	92 ...
48	černý jako	90 ...
49	o černý	88 ...
50	černý den	87 ...

Máte přístup pouze k 1 000 položkám. [Získat více](#)

Řádků na stránku

50

1–50 z 1 000



1 / 20



# SkELL

- webová aplikace pro učení jazyka na základě dat z korpusu
- česká verze: <https://csskell.sketchengine.co.uk/run.cgi/skell>
- verze pro angličtinu, ruštinu, italštinu, němčinu nebo estonštinu



Děkuji za pozornost

