

# KonText a Sketch Engine – shody a rozdíly korpusových vyhledávačů v českém prostředí<sup>1</sup>

*Klára Osolsobě*

osolsobe@phil.muni.cz

## Anotace

Přednáška se zaměří na praktický úvod do práce s jazykovými korpusy češtiny. Ukážeme nejdříve základní funkce nástroje pro práci s korpusy KonText a porovnáme je s analogickými funkcemi nástroje Sketch Engine. Na závěr vysvětlíme, kterými funkcemi disponuje nástroj Sketch Engine navíc oproti nástroji KonText a budeme demonstrovat, jak lze tyto funkce využít pro studium češtiny.

## Úvod

Jednou z podmínek, kterou musí jazykový korpus jako zdroj dat pro studium jazyka splňovat, je v moderní době elektronické uložení a elektronická přístupnost jazykových dat, která korpus tvoří. Pro češtinu existuje v současnosti velké množství jazykových korpusů, které tuto podmínku splňují. Přestože roku 1994 vzniklo centrální pracoviště, jímž je Ústav Českého národního korpusu na filozofické fakultě Univerzity Karlovy v Praze (<http://ucnk.ff.cuni.cz/cs/>), dosud přetrvává v českém prostředí jistá variabilita nástrojů počítačového zpracování přirozeného jazyka (dále NLP). Nejobvyklejším základním nástrojem pro práci s jazykovým korpusem je korpusový vyhledávač (manažer), který v zásadě umožňuje klást dotazy na jazykové jednotky a zobrazovat jejich výskyt ve formátu KWIC (**key word in context**).

Standardní korpusový vyhledávač (manažer) disponuje dalšími funkcemi, které umožňují **ukládat** vyhledaná data, **třídít** je a **analyzovat** je, přičemž nejobvyklejším typem analýzy je **frekvenční analýza** a dále analýzy založené na složitějších **statisticky ověřených metodách** pracujících s jazykovými daty. Výsledkem jsou potom frekvenční seznamy a seznamy kandidátů na relevantní kolokace pro klíčovou jednotku (klíčové slovo).

## Situace v českém prostředí

V českém prostředí (České republice) lze z uživatelského hlediska pracovat s několika podobami korpusových vyhledávačů, které se váží na akademická pracoviště. Na tomto místě pomineme vyhledávač B. Lehečky vytvořený pro Staročeskou textovou banku v rámci projektů pracoviště Ústavu pro jazyk český AV ČR (<http://vokabular.ujc.cas.cz/banka.aspx?idz=STB>) a budeme se věnovat patrně uživatelsky nejrozšířenějšímu vyhledávači, jímž je nástroj **KonText**. ([https://kontext.korpus.cz/first\\_form](https://kontext.korpus.cz/first_form)) a dále komerčnímu produktu firmy Lexical Computing

---

<sup>1</sup> Text vznikl za podpory projektu **MUNI/A/0862/2017** Čeština v jednotě synchronie a diachronie – 2018.

**Sketch Engine** (<https://www.sketchengine.co.uk/>), který vybíráme zejména proto, že nástroj – webové rozhraní **KonText** s ním historicky souvisí.

### **KonText a Sketch Engine historicky**

**KonText** je rozšířenou a graficky upravenou verzí původní aplikace **NoSketch Engine**, vyvíjenou v ÚČNK FF UK a na ÚFAL MFF UK pod licencí **GNU GPL 2**. Stejně jako **NoSketch Engine**<sup>2</sup> používá i KonText jako výpočetní server program Manatee.

Vyhledávač (grafické rozhraní) **KonText** slouží všem uživatelům, kteří pracují s korpusy zpřístupňovanými Ústavem Českého národního korpusu. Projekt **Český národní korpus** umožňuje pracovat s obecnými vyváženými korpusy synchronními i diachronními, psaného i mluveného jazyka. Kromě toho jsou k dispozici i různé specializované korpusy, mezi nimi i korpusy paralelní (více informací viz <http://wiki.korpus.cz/doku.php/cnk:uvod>). Velké synchronní psané korpusy jsou lematizovány a morfologicky (některé i syntakticky) označovány.

**Sketch Engine** je nástroj pro výzkum fungování přirozeného jazyka určený lingvistům, lexikografům, překladatelům a učitelům i studentům. Umožňuje přístup ke čtyřem stům korpusů, které představují bohatství devadesátky jazyků (přehled a seznam viz <https://www.sketchengine.co.uk/user-guide/user-manual/corpora/corpora-list/>).

Pro akademickou obec (studenty, pedagogy) Masarykovy Univerzity je dostupný již delší dobu. Nově od dubna 2018 do dubna 2022 je volně dostupný pro nekomerční využití přes infrastrukturu EU ELEXIS (<https://www.sketchengine.co.uk/elexis/>).

### **Přehled základních funkcí obou nástrojů**

Pokud chceme užívat jakýkoliv korpusový vyhledávač, volíme na začátku korpus, s nímž chceme pracovat. Odpověď na otázku, který z obou nástrojů, o nichž je řeč výše, je lepší, je podmíněna badatelským záměrem. Ten totiž musí vzít v úvahu i to, který korpus je vhodný k jeho uskutečnění.

Výhodou korpusů nabízených přes ÚČNK je dlouhodobá orientace pracoviště na tvorbu vyvážených korpusů češtiny a na dobrou úroveň lingvistických anotací. Výhodou nástroje Sketch Engine je pestrá nabídka korpusů různých jazyků. Se zřetelem k češtině spatřujeme jistě omezení v tom, že nabízí menší výběr korpusů než je k dispozici v ÚČNK. Nespornou výhodou nástroje SketchEngine je řada funkcí, které lze využít zejména v moderní korpusově orientované lexikografii. Některé z nich chceme ukázat, neboť je mohou využít i pokročilí studenti češtiny, a sice jako svého druhu on-line slovník.

---

<sup>2</sup> Manažer užívaný dříve pro práci s korpusy ÚČNK, který je nekomerční volně dostupnou (open-source) verzí softwarového nástroje Sketch Engine. Na rozdíl od něho nedisponuje některými funkcemi. Více: <https://nlp.fi.muni.cz/trac/noske>.

## Shodné funkce obou nástrojů

Nadále se tedy zmíníme o těch funkcích, které jsou společné, a ukážeme na praktických příkladech, jak je lze použít.

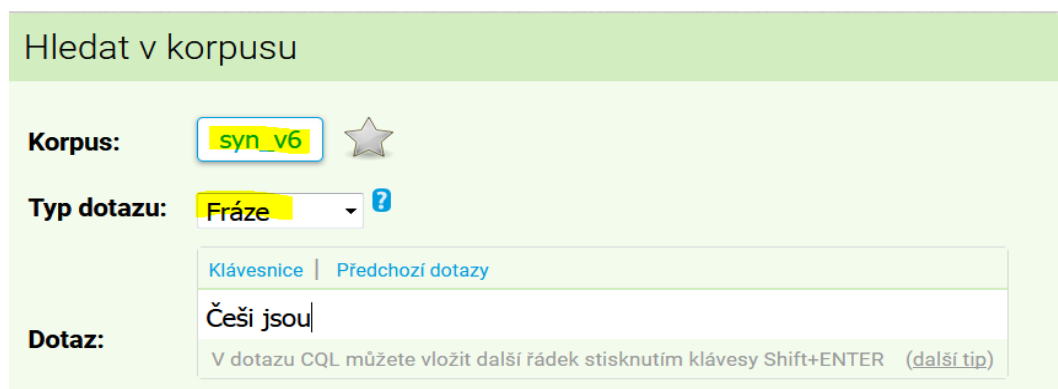
Základní funkcí každého korpusového vyhledávače je okno, do kterého lze zapsat dotaz na jazykovou jednotku, kterou chceme v korpusu najít. Můžeme se tázat na konkrétní slovní tvar (**word**), na slovní spojení (**frázi**), popřípadě můžeme konstruovat složitější dotaz za pomoci regulárních výrazů v jazyce **cql** (**corpus query language**). Pokud je korpus lemmatizovaný a značkový, pak se lze dotazovat také na všechny tvaroslovné varianty patřící pod jeden základní tvar, tedy na systémové slovo neboli **lemma**. **Typ dotazu** je v obou vyhledávačích třeba zvolit a dodržet formální náležitosti. V obou nástrojích fungují přístupy k manuálům – uživatelským příručkám<sup>3</sup>, v nichž jsou popsány všechny podrobnosti včetně například základních rysů jazyka **cql**.

Odpovědí na dotaz je vyhledaná jazyková jednotka ve všech výskytech ve zvoleném korpusu ve formátu KWIC (**key word in context**).

Mějme následující úkol. Chtěli bychom z korpusů zjistit, co říkají Češi sami o sobě. Budeme postupovat co nejjednodušeji tak, abychom přitom demonstrovali funkce obou nástrojů a zároveň porovnali výsledky na různých korpusech, které jsou přes srovnávané nástroje přístupné.

V rozhraní KonText zvolíme největší v současnosti dostupný korpus synchronní psané češtiny **syn\_v6** a přes **Typ dotazu** Fráze vyhledáme všechny výskyty spojení *Češi jsou* ve zvoleném korpusu (Obrázek 1).

Obrázek 1



Tutéž akci zopakujeme pomocí nástroje Sketch Engine, přičemž budeme pracovat s velkým webovým korpusem **czTenTen12** (Obrázek 2).

Obrázek 2

<sup>3</sup> V rozhraní KonText se přes menu **Nápověda/Help** → **Uživatelská příručka/User Manual** se podle volby jazyka dostaneme k podrobné internetové příručce v češtině nebo angličtině. V nástroji Sketch Engine je v liště přímo odkaz **Uživatelská příručka/Help** oba odkazy vedou k manuálu v angličtině.

Jednoduchý dotaz:

[Typy dotazů](#) [Kontext](#) [Typy textů](#) [?](#)

Typ dotazu  jednoduchý  lemma  fráze  slovo  znak  CQL

Lemma:  Slovní druh (PoS): nspecifikováno ▾

**Fráze:**  Slovní druh (PoS): nspecifikováno ▾

Slovní tvar:   citlivé na velikost znaků

Znak:  Implicitní atribut: word

CQL:

[Přehled značek](#) [CQL editor](#)

Výsledkem bude v obou případech konkordanční seznam (Obrázek 3 a 4).

Obrázek 3

víru v lidskost , ale i v to , že	<b>Češi jsou</b>	přece jen národem muzikantů a ne zlodějů , " dodaly
na dnešním losování základních skupin v Kyjevě zvláštní premiéru .	<b>Češi jsou</b>	totiž poprvé v samostatné historii žebříčkově nejhorší zemí , která
řekl . Třídíme nejvíce odpadu ZLATÁ MEDAILE V RECYKLACI .	<b>Češi jsou</b>	v třídění odpadu první v Evropě , lidé z Pardubického
ještě lépe postaveného Šmicra , přitukává mu balon a všichni	<b>Češi jsou</b>	v extázi ! Tedy ve fotbalovém nebi pro puristy .
3 . Co máte a nemáte rád na Česích ?	<b>Češi jsou</b>	obecně velmi milí , laskaví a skromní lidé . Někdy

Obrázek 4

a neúspěchem ! A at' se to někomu líbí nebo ne ,	<b>Češi jsou</b>	favorité . Blesk bodoval deset kategorií ( 5
musí Seveřany jednoznačně smáznout ! Forma	<b>Češi jsou</b>	zkrátka mazáci ! Před EURO " simulovali " že jim
tři výhry a při každé z nich nejprve prohrávali .	<b>Češi jsou</b>	mistři v otáčení zápasů ! Češi : 5 Dáni : 4
. Zraněný je obránce Jensen a útočník Sand .	<b>Češi jsou</b>	naopak zdraví jako rybičky . Češi : 5 Dáni : 3
je proti . Vyplývá to z průzkumu agentury STEM .	<b>Češi jsou</b>	podle průzkumu poměrně dobře informováni o

Dalšími funkcemi, jimiž běžně korpusové vyhledávače disponují, jsou **ukládání** vyhledaných dat, jejich další zpracování pomocí **filtrů**, které slouží k výběru/odstranění části dat, **frekvenční analýza** vyhledaných dat a **zobrazování** dat i metadat<sup>4</sup>.

K čemu slouží funkce **Uložit/Save**? Pomocí této funkce můžeme vyhledaná data uložit v textovém formátu a dále s nimi pracovat off-line. Pokud hledáme v korpusech ilustraci nějakého jevu v jazyce, pak nám tato funkce postačí. Příkladem je třeba využití jazykových korpusů v pedagogické praxi.<sup>5</sup>

K čemu slouží funkce **Filtr/Filter**? Pomocí této funkce můžeme například vybrat část konkordance. Představme si, že se chceme (v konkordančním seznamu uvedeném výše) podívat pouze na ty doklady, kdy za klíčovým slovesným tvarem *jsou* následuje adjektivum. V rozsahu hledání definovaném jako první pozice vpravo (interval <1,1>) bezprostředně po

<sup>4</sup> Metadata jsou jednak informace, které identifikují zdroje vyhledaných dokladů (typ textu, autora, žánr, ...), jednak lingvistické informace vložené do korpusu, přes něž lze vyhledávat na různých rovinách lingvistické abstrakce (například lemmata, slovní druhy a další gramatické významy reprezentované morfologickými značkami, tzv. tagy).

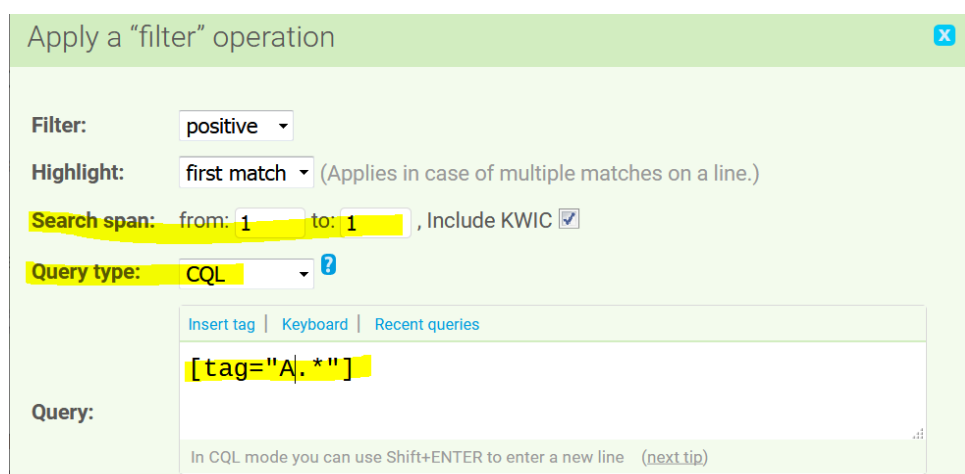
<sup>5</sup> Takto například pracují učitelé cizích jazyků s korpusy, když v nich hledají autentické věty jako ilustrační materiál nebo při tvorbě různých cvičení.

klíčovém slovu vyhledáme všechna adjektiva. Takto lze filtr použít, pokud je text morfologicky označován. V jazyce cql definujeme značku podle použitého tagsetu<sup>6</sup>.

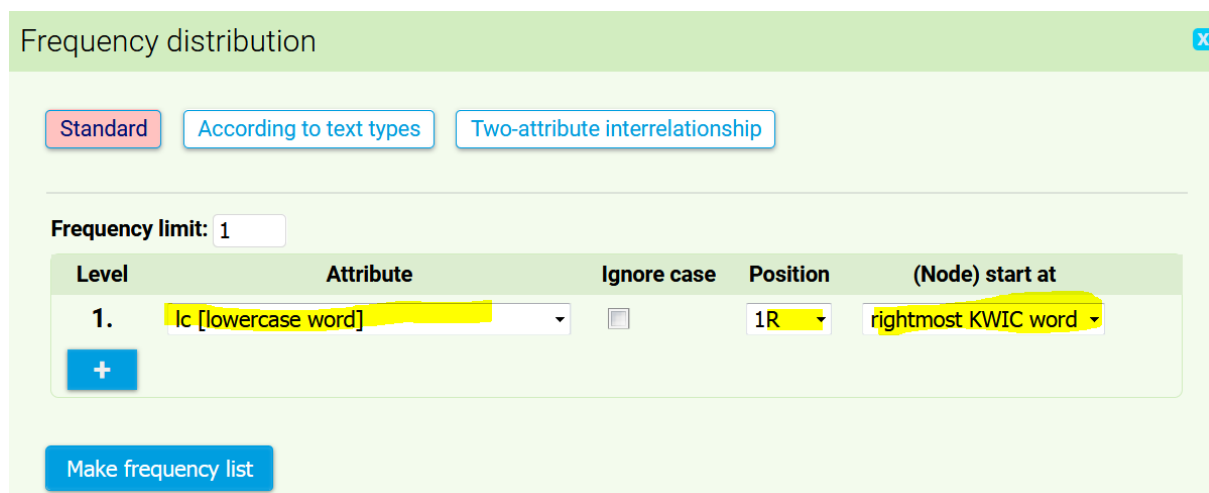
K čemu slouží funkce **Frekvence/Frequency**? Její pomocí můžeme vyhledaná data analyzovat z hlediska frekvence. Představme si, že chceme ve sledovaných korpusech zjistit frekvenci adjektiv vyskytujících se ve spojení *Češi jsou <u>stakoví</u>*.

Na screenshotech níže vidíme, jak se pracuje v obou systémech s popsányi funkcemi a jaké výsledky získáme ze zvolených korpusů. Na první trojici screenshotů (Obrázek 5, 6, 7) vidíme zadání v rozhraní KonText a výsledek frekvenční analýzy dat z korpusu syn\_v6, druhá trojice screenshotů (Obrázek 8, 9, 10) ukazuje zadání v nástroji Sketch Engine a výsledky frekvenční analýzy dat korpusu czTenTen12.

Obrázek 5



Obrázek 6



Obrázek 7

<sup>6</sup> V českém prostředí existuje vícero tagsetů. Korpusy přístupné přes grafické rozhraní **KonText** jsou značkovány tagsetem, který v lingvistické literatuře bývá označován jako „pražský“ /viz více zde: <https://wiki.korpus.cz/doku.php/seznamy:tagy>). Korpusy budované na FI MU užívají tagset označovaný jako „brněnský“ (viz <https://www.sketchengine.co.uk/tagset-reference-for-czech/>).

# Frequency list

1

Minimum frequency: 1

Apply

Total: 986 items (20 pages)

	Filter	lc [lowercase word]	Freq	
1	p / n	konzervativní	127	<div style="width: 127px; height: 10px; background-color: #007bff;"></div>
2	p / n	zvyklí	120	<div style="width: 120px; height: 10px; background-color: #007bff;"></div>
3	p / n	ochotni	81	<div style="width: 81px; height: 10px; background-color: #007bff;"></div>
4	p / n	dobří	77	<div style="width: 77px; height: 10px; background-color: #007bff;"></div>
5	p / n	schopni	56	<div style="width: 56px; height: 10px; background-color: #007bff;"></div>
6	p / n	nejlepší	50	<div style="width: 50px; height: 10px; background-color: #007bff;"></div>
7	p / n	skvělí	45	<div style="width: 45px; height: 10px; background-color: #007bff;"></div>
8	p / n	největší	43	<div style="width: 43px; height: 10px; background-color: #007bff;"></div>
9	p / n	velcí	40	<div style="width: 40px; height: 10px; background-color: #007bff;"></div>
10	p / n	lepší	40	<div style="width: 40px; height: 10px; background-color: #007bff;"></div>

Obrázek 8

## Filter konkordance

Filtr:  pozitivní  negativní

Vybraný token:  první  poslední

Rozsah hledání: od 1 do 1  zahrnout KWIC

Jednoduchý dotaz:  [Typy dotazů](#) [Typy textů](#)

Typ dotazu  jednoduchý  lemma  fráze  slovo  znak  CQL

Lemma:  Slovní druh (PoS): nespecifikováno

Fráze:

Slovní tvar:  Slovní druh (PoS): nespecifikováno  citlivé na velikost znaků

Znak:

CQL:  Implicitní atribut: word

[Přehled značek](#) [CQL editor](#)

Obrázek 9

## Víceúrovňové frekvenční rozložení ?

Frekvenční limit: 0

první úroveň

Atribut: **word (lowercase)**

Ignorovat velikost písmen

- 6L
- 5L
- 4L
- 3L
- 2L
- 1L
- Node
- 1R

Pozice: 2R

druhá úroveň

Atribut: word

Ignorovat velikost písmen

- 6L
- 5L
- 4L
- 3L
- 2L
- 1L
- Node
- 1R
- 2R

Pozice: 2R

Vytvořit seznam frekvencí

Obrázek 10

<a href="#">word (lowercase)</a>	<a href="#">Frekvence</a>	Items: 749    Total frequency: 1,775
<a href="#">P   N</a> známí	60	<div style="width: 60%;"></div>
<a href="#">P   N</a> zvyklí	57	<div style="width: 57%;"></div>
<a href="#">P   N</a> konzervativní	45	<div style="width: 45%;"></div>
<a href="#">P   N</a> smějící	40	<div style="width: 40%;"></div>
<a href="#">P   N</a> schopní	36	<div style="width: 36%;"></div>
<a href="#">P   N</a> ochotní	34	<div style="width: 34%;"></div>
<a href="#">P   N</a> dobří	33	<div style="width: 33%;"></div>
<a href="#">P   N</a> nejlepší	26	<div style="width: 26%;"></div>
<a href="#">P   N</a> líní	25	<div style="width: 25%;"></div>
<a href="#">P   N</a> hrdí	23	<div style="width: 23%;"></div>

K čemu slouží funkce **Zobrazení/View**? Její pomocí měníme zobrazení konkordančního seznamu a zobrazení metainformací. Na ilustračním screenshotu (Obrázek 11) vidíme, že výskyt negativně hodnotícího adjektivum *líní*, které se jako jediné negativní adjektivum objevuje mezi prvními deseti nejfrekventovanějšími adjektivy ve sledovaném kontextu, je pozorovatelně svázán s žánrem blogu, který má z definice subjektivní charakter.

Obrázek 11

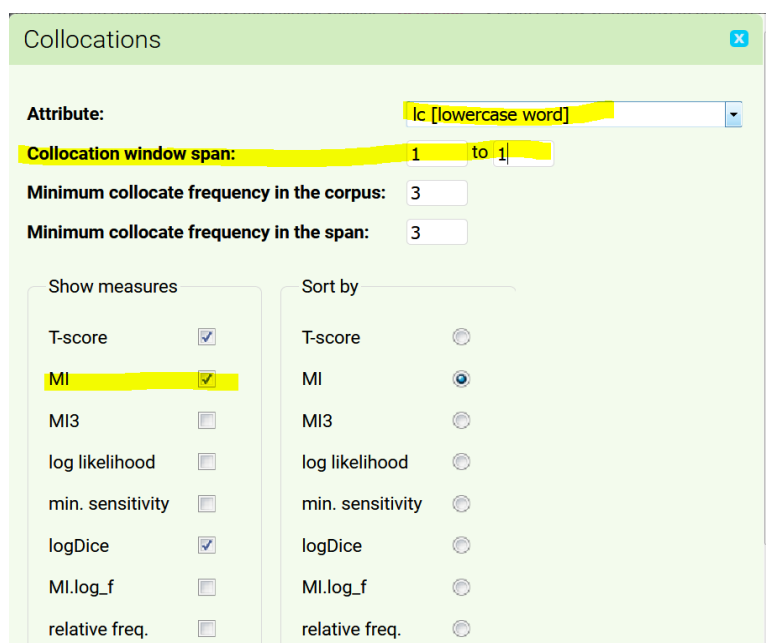
<a href="#">vyhledavani.ihned.cz</a>	cizinci ) jsou jediná , kteří opravdu pracují . <b>Češi jsou líní</b> jak vši a jen hledají výmluvy . </p><p> : I Mluv za
<a href="#">jansky.blog.idnes.cz</a>	uplně nejhloupější mi připadají argumenty , že <b>Češi jsou líní</b> a emigranti pracovití ( já a priori neříkám
<a href="#">blog2.idnes.cz</a>	až jako senilní dědek , tj . tak za 30 - 40 let . <b>Češi jsou líní</b> kverulanti , když něco dělají špatně , tak
<a href="#">zpravy.idnes.cz</a>	KČ . </p><p> A pak nějaký pisálek nebo politik napíše , že <b>Češi jsou líní</b> dojíždět za prací . Ale , že nejsou
<a href="#">zdn.cz</a>	, v němž se během čtyřiceti let nic nezmění . A <b>Češi jsou líní</b> myslet dopředu . Vzpomeňme , jak dlouho
<a href="#">noveclanky.cz</a>	letech využívají nabídky levnějších ... </p><p> <b>Češi jsou líní</b> . I když jsou farmářské trhy stále na
<a href="#">zpravy.zona.cz</a>	na ITC znalosti a schopnosti zaměstnanců . </p><p> <b>Češi jsou líní</b> recyklovat drobné elektrospotřebiče .
<a href="#">blog.aktualne.centrum.cz</a>	, jak nízká je produktivita zaměstnanců , ba že <b>Češi jsou líní</b> . Kdeko z majitelů firem má krizi . Pro
<a href="#">blog.aktualne.centrum.cz</a>	. A většina debatérů kupodivu taky ... že <b>Češi jsou líní</b> , vlastní osud je jim fuk apod . Nuže , vy
<a href="#">janasimonova.blog.idnes.cz</a>	přimluvíla u redaktorů Idnes , at' nepíše , že <b>Češi jsou líní</b> třídit odpad , ale at' už v názvu článku

Frekvenční analýzy lze ovšem dělat také sofistikovaněji a využít k nim metody corpus driven výzkumu kolokací. Oba nástroje umožňují získat kandidáty na relevantní kolokace klíčových

slov vyhledaných v korpusech. Využívají k tomuto účelu statistické výpočty, které dávají dobré výsledky při práci s jazykovými daty<sup>7</sup>. Oba nástroje disponují s funkcí **Kolokace/Collocations**.

Na následujících screnshotech ukážeme přehled statisticky významných kolokací na prvním místě vpravo od klíčového spojení (fráze) *Češi jsou*. Demonstrujeme tak jak funkci obou nástrojů (Obrázek 12 a 14), tak i rozdíly ve výsledcích prosté frekvenční analýzy v porovnání se sofistikovanější analýzou relevantních kolokací (Obrázek 13 a 15). Tentokrát jsme nesledovali jen adjektiva, ale všechna slova napravo od spojení (fráze) *Češi jsou*. Výsledky jsou pouze ilustrací k fungování obou nástrojů.

Obrázek 12



Obrázek 13

<sup>7</sup> Na obrázcích níže jsou kolokace seřazeny podle asociační míry MI-score (Obrázek 14) a LogDice (Obrázek 15). Viz více [https://wiki.korpus.cz/doku.php/pojmy:asociacni\\_miry?redirect=1#dice\\_a\\_logdice](https://wiki.korpus.cz/doku.php/pojmy:asociacni_miry?redirect=1#dice_a_logdice).



## Kandidáti na kolokace ?

Atribut: **word (lowercase)** V rozsahu od: **1** do: **1**

Minimální frekvence v korpusu: **5**

Minimální frekvence v daném rozsahu: **3**

Zobrazit funkce: **logDice** Třídít podle: **logDice**

Obrázek 14

### Collocation candidates

1 ▶

	Filter	lc	Freq ▼	MI	T-score	logDice
1.	p/n	pověstní	3	16.845	1.732	2.843
2.	p/n	dobírkový	4	16.623	2.000	3.258
3.	p/n	nedočkavější	3	14.693	1.732	2.840
4.	p/n	nejateističtějším	5	14.360	2.236	3.572
5.	p/n	pivařský	4	14.108	2.000	3.251
6.	p/n	rovnostáři	5	14.065	2.236	3.570
7.	p/n	rasisti	28	13.830	5.291	5.997
8.	p/n	bytostní	3	13.343	1.732	2.834
9.	p/n	švejci	5	13.326	2.236	3.563
10.	p/n	remcalové	3	13.301	1.732	2.833

Obrázek 15

## Kandidáti na kolokace

Strana   [Další >](#)

	<u>Poučet souvýskytů</u>	<u>Počet kandidátů</u>	<u>T-score</u>	<u>MI</u>	<u>logDice</u>
<a href="#">P</a>   <a href="#">N</a> <b>národem</b>	186	14,354	13.636	13.083	8.118
<a href="#">P</a>   <a href="#">N</a> <b>rasisti</b>	57	2,050	7.549	14.184	7.601
<a href="#">P</a>   <a href="#">N</a> <b>smějící</b>	40	2,646	6.323	13.305	7.003
<a href="#">P</a>   <a href="#">N</a> <b>rasisté</b>	28	1,818	5.290	13.332	6.610
<a href="#">P</a>   <a href="#">N</a> <b>mistři</b>	69	27,940	8.301	10.691	5.992
<a href="#">P</a>   <a href="#">N</a> <b>přeborníci</b>	16	790	3.999	13.727	5.970
<a href="#">P</a>   <a href="#">N</a> <b>národ</b>	230	114,572	15.154	10.392	5.947
<a href="#">P</a>   <a href="#">N</a> <b>zbabělci</b>	16	3,195	3.998	11.711	5.605
<a href="#">P</a>   <a href="#">N</a> <b>líní</b>	25	9,953	4.997	10.716	5.546
<a href="#">P</a>   <a href="#">N</a> <b>burani</b>	11	2,002	3.315	11.845	5.234

### Funkce nástroje Sketch Engine navíc

Na závěr bychom rádi ilustrovali některé funkce, jimiž disponuje nástroj Sketch Engine navíc oproti rozhraní KontText.

První je funkce **Word Sketch/Slovní profil** sloužící k vizualizaci frekvenčně uspořádaných gramaticky definovaných relací, do kterých vstupuje klíčové slovo v daném korpusu. Na screenshotu (obrázek 16) vidíme přehled Word Sketch k lemmatu *Čech* z korpusu czTenTen12. Povšimněme si například sloves, o nichž lze předpokládat, že substantivum plní roli subjektu, a která následují příslušné substantivum.

Obrázek 16

**čech** (noun)  
Czech Web 2012 (czTenTen12 v9) frekvence = [451,932](#) (89.14 v milionu)

<u>a_modifier</u>	<u>is_subj_of</u>	<u>gen_2</u>	<u>post_verb</u>	<u>prec_prep</u>
11.23	16.76	10.34	8.86	7.30
<b>volyňský +</b> <a href="#">989</a> 9.26 volyňských Čechů	<b>utratit +</b> <a href="#">422</a> 7.39 <b>jezdit +</b> <a href="#">782</a> 7.18 Češi jezdí	<b>svatopluk +</b> <a href="#">1,257</a> 9.69 Svatopluka Čecha	<b>nakupovat +</b> <a href="#">254</a> 7.19 Češi nakupují	<b>mezi +</b> <a href="#">3,939</a> 6.22 mezi Čechy a
<b>žijící +</b> <a href="#">910</a> 7.83 Čechů žijících	<b>krást +</b> <a href="#">331</a> 6.98 Češi kradou	<b>praotec +</b> <a href="#">912</a> 9.26 praotce Čecha	<b>jezdit +</b> <a href="#">656</a> 7.11 Češi jezdí	<b>vůči +</b> <a href="#">468</a> 6.08 vůči Čechům
<b>jižní +</b> <a href="#">2,815</a> 7.78 jižní Čechy	<b>umět +</b> <a href="#">653</a> 6.85 <b>prohrát +</b> <a href="#">373</a> 6.72 Češi prohráli	<b>petr +</b> <a href="#">2,676</a> 8.77 Petra Čecha	<b>krást +</b> <a href="#">196</a> 7.09 Češi kradou	<b>kromě +</b> <a href="#">395</a> 5.07 . Kromě Čechů
<b>rodilý +</b> <a href="#">389</a> 7.61 rodilý Čech	<b>prohrát +</b> <a href="#">373</a> 6.72 Češi prohráli	<b>procent +</b> <a href="#">1,914</a> 8.71 procent Čechů	<b>kupovat +</b> <a href="#">264</a> 6.88 Češi kupují	<b>proti +</b> <a href="#">996</a> 4.87 proti Čechům
<b>průměrný +</b> <a href="#">1,807</a> 7.45 průměrný Čech	<b>nakupovat +</b> <a href="#">302</a> 6.71 Češi nakupují	<b>soužití +</b> <a href="#">628</a> 8.57 soužití Čechů a Němců	<b>věřit +</b> <a href="#">459</a> 6.82 <b>utratit +</b> <a href="#">153</a> 6.74 Češi utratí za	<b>vedle +</b> <a href="#">171</a> 4.63 . Vedle Čechů
<b>etnický +</b> <a href="#">423</a> 7.26 etnických Čechů	<b>vyhrát +</b> <a href="#">785</a> 6.70 Češi vyhráli	<b>třetina +</b> <a href="#">1,072</a> 8.33 dvě třetiny Čechů	<b>umět +</b> <a href="#">614</a> 6.68 <b>utracet +</b> <a href="#">129</a> 6.56 Češi utrácejí	<b>oproti +</b> <a href="#">100</a> 4.43 oproti Čechům
<b>východní +</b> <a href="#">1,362</a> 7.20 Východní Čechy	<b>bojovat +</b> <a href="#">378</a> 6.56 Češi bojovali	<b>ringo +</b> <a href="#">373</a> 8.01 Františka Ringo Čecha	<b>miřit +</b> <a href="#">206</a> 6.43	<b>pro +</b> <a href="#">5,671</a> 4.41 pro Čechy
<b>hrdý +</b> <a href="#">279</a> 7.09	<b>věřit +</b> <a href="#">395</a> 6.51	<b>čtvrtina +</b> <a href="#">560</a> 7.99		<b>včetně +</b> <a href="#">255</a> 4.06

Nástroj má zabudována pravidla parciální syntaktické analýzy založené na morfologických značkách. Tak například na základě toho, že se v bezprostředním levém kontextu substantiva vyskytuje adjektivum, které se shoduje se substantivem v relevantních gramatických kategoriích, je vytvořen seznam a \_modifier (adjektivních modifikátorů) typických (s relevantí frekvencí) pro klíčové substantivum.

Druhou funkcí je **Thesaurus** (zobrazení podobných slov: na základě porovnání kontextů je vytvořen seznam a vizualizace slov, která mají podobné gramaticko-lexikální kontexty). Na screenshotu (obrázek 17) jsou slova vyskytující se v podobných kontextech jako slovo *Čech*.

Obrázek 17

**čech** (noun) Czech Web 2012 (czTenTen12 v9) frekvence = [451,932](#) (89.14 v milionu)

Lemma	Skóre	Frekvence
němec	0.480	252,884
američan	0.448	172,108
obyvatel	0.419	656,805
občan	0.414	818,691
student	0.412	847,292
politik	0.410	432,673
rom	0.401	175,093
rodič	0.395	949,755
národ	0.395	415,829
muž	0.394	1,616,463
podnikatel	0.389	289,755
učitel	0.389	456,974
hráč	0.388	1,580,578
člen	0.387	1,391,941
křesťan	0.385	194,918
žena	0.384	1,899,452
cizinec	0.381	151,574



Třetí funkcí je **Sketch rozdíl** (vizualizace shod a rozdílů kontextu dvojice slov). Na následujících screenshotech je srovnání frekvencí sloves stojících v postpozici za dvojicí *Čech/Němec* (Obrázek 18) a *Čech/Američan* (Obrázek 19). Sloveso v postpozici za slovními tvary lemmatu *Čech*, které se nikdy nevyskytuje v téže pozici spolu s lemmaty *Němec* nebo *Američan* je podle screenshotů *krást*. Naopak slovesa stojící (v analyzovaném korpusu) pouze v kontextu lemmat *Němec* nebo *Američan* (nikoli *Čech*) jsou *zaútočit* a *bombardovat*. Jediným společným slovesem vyskytujícím se v bezprostředním kontextu všech tří lemmat (*Čech/ Němec/ Američan*) je *milovat*. Takto na základě statistických šetření řízených korpusovými daty lze vytvářet objektivní obraz o tom, **co na sebe prozradíme, když o sobě a druhých píšeme**.<sup>8</sup>

Obrázek 18

<sup>8</sup> Zajímavé je porovnání Sketch rozdílů dalších analogických dvojic jako *Čech/Rus*, *Čech/Rom* a *Čech/Žid*.

post_verb	40,024	11,123	0.09	0.04
krást	<u>196</u>	0	7.1	--
utrácet	<u>129</u>	0	6.6	--
mířit	<u>206</u>	<u>13</u>	6.4	3.1
kupovat	<u>264</u>	<u>23</u>	6.9	4.1
utratit	<u>153</u>	<u>9</u>	6.7	4.0
nakupovat	<u>254</u>	<u>20</u>	7.2	4.5
vypít	<u>122</u>	<u>8</u>	6.4	3.9
jezdit	<u>656</u>	<u>100</u>	7.1	4.7
věřit	<u>459</u>	<u>78</u>	6.8	4.6
umět	<u>614</u>	<u>116</u>	6.7	4.5
cestovat	<u>130</u>	<u>14</u>	6.3	4.3
milovat	<u>279</u>	<u>52</u>	6.4	4.4
říkat	<u>275</u>	<u>289</u>	4.9	5.1
okupovat	<u>13</u>	<u>16</u>	3.3	5.1
česat	<u>12</u>	<u>16</u>	3.2	5.3
útočit	<u>15</u>	<u>29</u>	3.2	5.3
nazývat	<u>20</u>	<u>58</u>	3.2	5.5
prchat	<u>9</u>	<u>16</u>	2.8	5.2
střílet	<u>8</u>	<u>19</u>	2.5	5.1
obchodovat	<u>11</u>	<u>26</u>	3.0	5.7
obsazovat	0	<u>15</u>	--	5.1
kšeftovat	0	<u>12</u>	--	5.1
bombardovat	0	<u>13</u>	--	5.1
zaútočit	0	<u>18</u>	--	5.3
ustupovat	0	<u>26</u>	--	5.4

Obrázek 19

post_verb	40,024	15,569	0.09	0.09
krást	<u>196</u>	0	7.1	--
vypít	<u>122</u>	0	6.4	--
mířit	<u>206</u>	<u>13</u>	6.4	3.0
jezdit	<u>656</u>	<u>84</u>	7.1	4.4
kupovat	<u>264</u>	<u>27</u>	6.9	4.2
preferovat	<u>138</u>	<u>20</u>	6.3	4.3
nakupovat	<u>254</u>	<u>38</u>	7.2	5.3
umět	<u>614</u>	<u>152</u>	6.7	4.9
pít	<u>154</u>	<u>30</u>	6.3	4.7
cestovat	<u>130</u>	<u>26</u>	6.3	4.9
utratit	<u>153</u>	<u>35</u>	6.7	5.7
milovat	<u>279</u>	<u>108</u>	6.4	5.4
důvěřovat	<u>106</u>	<u>29</u>	6.3	5.5
utrácet	<u>129</u>	<u>39</u>	6.6	6.0
věřit	<u>459</u>	<u>300</u>	6.8	6.5
slavit	<u>75</u>	<u>69</u>	5.1	5.5
spotřebovat	<u>51</u>	<u>40</u>	5.0	5.6
přiznávat	<u>55</u>	<u>52</u>	5.0	5.7
konzumovat	<u>46</u>	<u>36</u>	5.1	5.8
tvrdit	<u>121</u>	<u>221</u>	4.6	5.7
domnívat	<u>34</u>	<u>38</u>	4.5	5.7
odejít	<u>10</u>	<u>54</u>	2.4	5.5
zaútočit	0	<u>28</u>	--	5.6
bombardovat	0	<u>34</u>	--	6.1

## Závěr

Ukázali jsme na ilustračních příkladech shodné rysy obou korpusových nástrojů použitelných pro práci s korpusy českého jazyka. Genetická příbuznost obou nástrojů je zřejmá. Uživatel se snadno zorientuje a může pracovat s oběma nástroji. Zvolí si jeden nebo druhý s ohledem na konkrétní korpusy, které manažery zpřístupňují. Oba nástroje mají velmi dobré instruktivní uživatelské příručky, takže je možné naučit se s nimi pracovat i bez pomoci učitele.

Pro studenty češtiny začátečníky je příznivé, že, manuály jsou dostupné i v angličtině.

**Sketch Engine** nabízí také nově nástroj **SkELL (Sketch Engine for Language Learning)**. Jedná se o jednoduchý nástroj pro studenty a učitele pro rychlé vyhledání užití jednotlivého slova/fráze v textech rodilých mluvčích (<https://cshell.sketchengine.co.uk/run.cgi/skell>). Nástroj disponuje funkcí pro hledání **příkladových vět (konkordancí)**, **slovních profilů (Words Sketch)** a **podobných slov (Thesaurus)**. Jde o pomůcku, která může nahradit výkladový slovník.

Také manažer **KonText** disponuje oproti **Sketch Engine** dalšími funkcemi a nástroji, kterým jsme zde nevěnovali patřičnou pozornost. Pro studenty češtiny jako cizího jazyka upozorníme na závěr alespoň na nástroj **Treq** k vyhledávání překladových ekvivalentů založený na paralelních korpusech češtiny a řady dalších jazyků. Tento nástroj může posloužit jako doplněk/náhrada překladového slovníku.<sup>9</sup>

## Bibliografie

Jan Hajič (2004): *Disambiguation of Rich Inflection (Computational Morphology of Czech)*. Vol. 1. Karolinum Charles University Press, Praha.

Milena Hnátková, Michal Křen, Pavel Procházka, Hana Skoumalová (2014): The SYN-series corpora of written Czech. In: *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, s. 160–164. Reykjavík: ELRA. [http://www.lrec-conf.org/proceedings/lrec2014/pdf/294\\_Paper.pdf](http://www.lrec-conf.org/proceedings/lrec2014/pdf/294_Paper.pdf)

Adam Kilgarriff, Pavel Rychlý, Pavel Smrž, David Tugwell (2004). Itri-04-08 the sketch engine. *Information Technology*.

Adam Kilgarriff, Vojtěch Kovář, Simon Krek, Irena Srdanovič, Carole Tiberius (2010): A quantitative evaluation of word sketches. *Proceedings of the 14th EURALEX International Congress*, s. 372–379.

Adam Kilgarriff, Vít Baisa, Jan Bušta, Miloš Jakubíček, Vojtěch Kovář, Jan Michelfeit, Pavel Rychlý, Vít Suchomel (2014): The Sketch Engine: ten years on. *Lexicography*, 1, s. 7–36.

Miloš Jakubíček, Adam Kilgarriff, Diana McCarthy, Pavel Rychlý (2010): Fast Syntactic Searching in Very Large Corpora for Many Languages. *PACLIC*, s. 741–747.

---

<sup>9</sup> Před nějakou dobou byl v českém tisku zveřejněn seznam slov, s jejichž významem si nevědí rady ani české děti. Víte, co znamenají slova: *bezbřehý, bytelný, balamutit, cudný, horlivý, hvozd, chmurný, chrabrý, jihnout, jímavý, jízlivý, kasat se, kazajka, klání, kloudný, komolit, konejšit, lačný, láteřit, ledabylý, lomožit, lpět, mamon, mdlý, mimoděk, nedůtklivě, nejapný, niterný, okounět, osočit, ostýchavý, otálet, otrapa, perný, pohnutka, pokoutně, ponurý, pookřát, pověra, potutelný, proradný, prchlivý, předpojatý, pyřit se, rmoutit, se, rozšafný, rusý, schlíplý, slídit, spílat, srdnatý, strádat, střenka, svérázný, svízel, sudí, šev, tklivý, trýznit, úděl, uhranout, unylý, úlisný, upejpat se, úporný, úskalí, uštěpačný, útlocitný, vesměs, vzývat, záhy, zakabonit se, zakolisat, záludný, zášť, zesinat, zevrubný, zmerčit, zpupný, ztepilý? Zkuste se na ně podívat zde: <http://treq.korpus.cz/index.php> a zde: <https://cshell.sketchengine.co.uk/run.cgi/skell#>.*

Tomáš Jelínek (2008): Nové značkování v Českém národním korpusu. In: *Naše řeč*, 91, 1, s. 13–20.

Vladimír Petkevič (2014): Problémy automatické morfologické disambiguace češtiny. In: *Naše řeč*, 97, 4, s. 194–207.

Pavel Rychlý (2008): A Lexicographer-Friendly Association Score. *Proc. 2nd Workshop on Recent Advances in Slavonic Natural Languages Processing, RASLAN*, 2, s. 6–9.

Drahomíra Spoustová, Jan Hajič, Jan Votrubec, Pavel Krbec, Pavel Květoň (2007): The Best of Two Worlds: Cooperation of Statistical and Rule-Based Taggers for Czech. In: *Proceedings of the Workshop on Balto-Slavonic Natural Language Processing. ACL 2007*, Praha, s. 67–74.

### **Elektronické korpusy a korpusové nástroje**

Křen, M. – Cvrček, V. – Čapka, T. – Čermáková, A. – Hnátková, M. – Chlumská, L. – Jelínek, T. – Kovářiková, D. – Petkevič, V. – Procházka, P. – Skoumalová, H. – Škrabal, M. – Truneček, P. – Vondříčka, P. – Zasina, A.: *Korpus SYN, verze 6 z 18. 12. 2017*. Ústav Českého národního korpusu FF UK, Praha 2017. Dostupný z WWW: <http://www.korpus.cz>

FI MU – *czTenTen12*. Centrum zpracování přirozeného jazyka FI MU, Brno. Dostupný z: WWW: <http://ske.fi.muni.cz/bonito>.

Nástroj *KonTetxt* dostupný z: WWW: [https://kontext.korpus.cz/first\\_form](https://kontext.korpus.cz/first_form).

Nástroj *Sketch Engine* dostupný z: WWW: <https://www.sketchengine.co.uk/>.

Nástroj *SkELL* dostupný z: WWW: <https://www.sketchengine.co.uk/skell/>.

Nástroj *Treq* dostupný z: WWW: <http://treq.korpus.cz/index.php>.