

CJDSL001 Korpusová lingvistika (1)

Klára Osolsobě

osolsobe@phil.muni.cz

Experimentální a počítačová lingvistika

O čem budeme mluvit

- Krátký historický exkurz
- Definice korpusu v moderním slova smyslu
- Dva metodologické přístupy k vytěžování korpusu
- Dva pohledy na korpus (lingvista a informatik)
- Filologie a korpusy
- Výuka jazyků a korpusy

Krátký historický exkurz

- Myšlenka korpusu
- Korpusová lingvistika – empirická disciplína
- Data a introspekce
- Technický pokrok

Konkordance – KWIC (Key Word In Context)

Výskytů: 213 | i.p.m. 1,76 (vztaženo k celému "omezeni/syn2015") | ARF 88,32 | Výsledek je promíchán

1 / 6

Výběr řádků: základní | Atributy:

<input type="checkbox"/>	Brno Business	se ho 4379 lidí . Vyplynulo z něj , že	Češi jsou	se svými šéfy poměrně spokojeni . Šéfů , kteří se
<input type="checkbox"/>	Reflex	jsou hlavní autority , pak je to svědectvím , že	Češi jsou	národ-nenárod , jenž se vždycky sjednotí na nějakém kýči ,
<input type="checkbox"/>	Maxim	bezpečnost . Žádný slovenský záchranář si nezapomene rýpnout , že	Češi jsou	blázni . Devět z deseti mrtvých v Tatrách jsou Češi
<input type="checkbox"/>	Týden	, který máme dobýt obchvatem shora , horskými hřebenovkami .	Češi jsou	na Tematině ostatně dávnými hosty , hrad pokořil už Přemysl
<input type="checkbox"/>	Sport	. Česká vzpoura v Letošní NHL je v polovině a	Češi jsou	vidět ! Po několika hubenějších sezonách je tu konečně vydařený
<input type="checkbox"/>	Reflex	vnímáno hlavně jako neustálý konflikt dvou zneprátených národů , kde	Češi jsou	vždy ti malí , ale spravedliví , utlačovaní , ale
<input type="checkbox"/>	Mladá fronta Dnes	by to ryzí zázrak . Matematicky to spočítat nelze .	Češi jsou	odkázání na pomoc shůry . Kromě toho , že sami
<input type="checkbox"/>	Lidé a země	předpisu rakouské armády z roku 1749 . " Proslulí nadšenci	Češi jsou	v Evropě obecně známi jako jedni z nejlepších „ vojáků
<input type="checkbox"/>	Reflex	, že mají zhasínat , pokud jdou z místnosti .	Češi jsou	národem s maloměstským myšlením a jsou závistiví . Pozorujete to
<input type="checkbox"/>	Reflex	koncem dubna byl projednán návrh zákona o její regulaci .	Češi jsou	totiž vůči prostituci značně liberální . A nejen vůči ní
<input type="checkbox"/>	Mladá fronta Dnes	Alsasané ve Francii skoro dodnes . Řekla bych , že	Češi jsou	dnes s vyrovnáváním se s touto částí naší minulosti dál
<input type="checkbox"/>	Týdeník Květy	v hrubém domácím produktu na hlavu je přímo propastný :	Češi jsou	35x bohatší ! Za zmínku určitě také stojí , že
<input type="checkbox"/>	Parlamentní listy	Češi leniví a s tím klesá jejich ochota recyklovat »	Češi jsou	ochotni odnášet vyřazené světelné zdroje maximálně do vzdálenosti 7
<input type="checkbox"/>	Mladá fronta Dnes	tomu není žádný důvod , nic se nemění , ale	Češi jsou	velmi finančně ngramotní . Nepřehnal se to v Česku s
<input type="checkbox"/>	Reflex	naprosto zřejmé , že dává přednost startování auta klikou ?	Češi jsou	někdy trochu záhadní . Taxikář se mě znovu zeptal na
<input type="checkbox"/>	Rodinný dům	několik hollywoodských komedií o Santa Clausovi , testující , zda	Češi jsou	již převychování , takže lze očekávat sérii sebevražd několika program
<input type="checkbox"/>	Respekt	, že oslepl právě po vodce této společnosti . „	Češi jsou	nejen světovými přeborníky v pití piva , ale podle Světové
<input type="checkbox"/>	Metro	na horách , v Paříži nebo překvapivě doma . „	Češi jsou	tedy spokojení nebo bez fantazie , " vyplývá z průzkumu
<input type="checkbox"/>	Parlamentní listy	Dánsko , Švédsko , Německo , Velká Británie) ,	Češi jsou	se 17 % polepšených kuřáků v dolní polovině tabulky .
<input type="checkbox"/>	Lidové noviny	poklesu spotřebitelské poptávky . Není důvodem i to , že	Češi jsou	příliš opatrní a skeptičtí ohledně stavu země a jejího dalšího

PSJČ

psjc Příruční slovník jazyka českého

vroucí adj. **vařící** **horký** **vřelý** Vzala čajovou konvičku, hodila tam bezový květ a nalila na to vroucí vody. **Vrch.** Vroucí vodu napouštěla na odvar čaje ve sklenicích. **Mach.** Vařil na lihovém kahanu lék. Odbíhaje od vroucího plecháčku k oknu své pracovny, pečlivě pozoroval noční oblohu. **John.** Jakási jiná, vroucí, tmavá krev byla nalita v žilách těchto lidí. **R. Svob.** Z nedovřených rtů sálal jí vroucí dech. **Čap. Ch. D Zř. vařící se, kolotající, vířící.** Vlny bičované a rozryté tisícerými křižujícími se jizvami a vroucí tisícerými víry. **Pašek. D citově opravdový; hluboký, vřelý.** Mluvil tlumeně, vroucím hlasem. **Jir.** U příležitosti pátého výročí osvobození Československa posílám Vám své vroucí přátelské pozdravy. **R. právo.** Dívka zvedla kvapně skloněnou hlavu, vroucí pohled utkvěl na okamžik na mladém, statečném muži. **Jir.** Bůh vroucí modlitby její neoslyšel. **Arb.** Přitiskla vroucí políbení na jeho ruku. **Schulz.** Zasloužíš si takto někdy ještě matčin vroucí dík. **Zey.** Má vroucí duše žije jen tobě, přírodo. **Vrch.** Ty, od jakživa tak vroucí ctitel všeho venkovského, jsi nevěren svým zásadám? **Just.** Byl vroucím přívržencem josefinského osvícenství. **J. Vlč.** Je tu mnoho hluboce procítěného a opravdu vroucího. **Vrch.**

SEU (SURVEY OF ENGLISH USAGE)



BROWN CORPUS

W. Nelson Francis - Henry Kucera

- 1964
- 1. elektronicky zpracovaný korpus
- 1 milion slovních tvarů
- britská a americká angličtina
- pečlivý výběr textů
- vzorky

Definice korpusu v moderním slova smyslu

- Elektronické uložení
- Elektronická přístupnost
- Definovaný obsah (ČEHO) a rozsah (KOLIK)
- Standardní anotace – metada a interpretace jazykových jednotek
- Rychlost, spolehlivost a opakovatelnost vyhledávání a kvantifikace nalezeného

Dva metodologické přístupy k vytěžování korpusu

- Corpus based / korpusem ověřovaný, na korpusu založený výzkum
- Pravidlo/výjimka – otevřený/uzavřený seznam, frekvence
- Corpus driven / korpusem inspirovaný výzkum, korpusem řízený výzkum
- Výzkum kolokací /lexical bundles

Maskulina mají v češtině v gen. pl. koncovku –ů (pánů, hradů, mužů, strojů, předsedů, soudců). Z tohoto pravidla existují výjimky. Které? Kolik?

[lc!=".*ů" & tag="NN[MI]P2.*"]

za nímž sedí u svého náčiní rybáři . Ostatním vyjma	rozhodčích/rozhodčí/NNMP2-----A-----	je tam po dobu soutěže zakázán vstup . " Letos
ský Beroun - Policisté obvinili osmadvacetiletou chovatelku	koní/kůň/NNMP2-----A-----	z Moravského Berouna na Bruntálsku z týrání zvířat . Podle
v Mosambiku za posledních 50 let postihly zatím asi milion	lidí/člověk/NNMP2-----A-----	. Oficiálně je registrováno kolem 200 obětí , skutečný počet
a Pospěcha spokojen . " Potřebovali bychom však takových	lidí/člověk/NNMP2-----A-----	víc , " poznamenal . Radní chtějí ve Zlíně zrušit
je prý marné . Tištěné ploše dominují názory starostů a	radních/radní/NNMP2-----A-----	, odlišným sdělením je věnováno jen 2,6 % této plochy
, jako 71 . hráče v pořadí . V dresu	Flyers/flyers/NNIP2-----A-----	odehrál 64 zápasů s bilancí 11 gólů a 26 přihrávek
Regio Taxis Bohemia . Ten by zahrnoval asi dvacet tisíc	lidí/člověk/NNMP2-----A-----	, to znamená pětaticet obcí od Benátek nad Jizerou až
místě tragédie policejní komisař Emmanuel Adebayo . Z pěti	lidí/člověk/NNMP2-----A-----	, kteří byli nejbliž místu výbuchu , zbyly jen zčernalé
oho ze zkušených ruských letců , podle něhož smrt desítek	lidí/člověk/NNMP2-----A-----	zavinila lidská hamižnost . Pilot Vjačeslav Achremenko potv
. " Přijeli lidé z tréninkových skupin z Plzně ,	Klatov/klatovy/NNIP2-----A-----	, Přeštic , Rokycan a Třemošné , " uvedl šéf
něj u nás na každých 100 osob připadá téměř 7	nemocných/nemocný/NNMP2-----A-----	a délka nemocenské je mnohdy delší než 30 dní .
Tu pořádá Občanské sdružení Omega . Jméno projekční	lidí/člověk/NNMP2-----A-----	něže a stále obdívá hodaš . Ve středu převzeme první

Celkem: 5056 (102 str.)

	Filtr	word	Frekvence
1.	p/ n	lidí	2 515 537
2.	p/ n	peněz	647 968
3.	p/ n	dní	465 084
4.	p/ n	obyvatel	243 036
5.	p/ n	přátel	125 885
6.	p/ n	koní	102 894
7.	p/ n	hostí	92 521
8.	p/ n	rozhodčích	78 016

	Filtr	word	Frekvence
1.	p/ n	milionů	1 400 417
2.	p/ n	metrů	806 199
3.	p/ n	bodů	725 571
4.	p/ n	hráčů	588 016
5.	p/ n	kilometrů	504 610
6.	p/ n	dolarů	491 535
7.	p/ n	měsíců	475 891
8.	p/ n	milionů	474 446
9.	p/ n	mužů	453 972

Jaké je mínění o Češích?

Kolokace na pozici 1-3 vpravo od KWIC <Češi jsou> seřazené podle míry **MI-score**

	Filtr	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
11.	p/n	národem	277	13.259	16.642	8.583
12.	p/n	šetřiví	4	13.163	2.000	3.243
13.	p/n	recesisti	3	13.108	1.732	2.832
14.	p/n	náruživými	4	12.899	2.000	3.240
15.	p/n	nejateističtější	4	12.723	2.000	3.237
16.	p/n	smějící	36	12.708	5.999	6.242
17.	p/n	švejkové	3	12.631	1.732	2.827
18.	p/n	přeborníci	23	12.495	4.795	5.642
19.	p/n	studení	9	12.494	2.999	4.372
20.	p/n	slovani	3	12.336	1.732	2.823
21.	p/n	nejsprostší	6	12.329	2.449	3.801
22.	p/n	rasisté	21	12.310	4.582	5.505

Češi jsou národem ...

	Filtr	lc	Freq ▼	MI	T-score	logDice
1.	p/n	kutilů	17	17.397	4.123	8.125
2.	p/n	pivařů	19	17.338	4.359	8.094
3.	p/n	houbařů	31	17.186	5.568	8.023
4.	p/n	chatařů	18	16.397	4.243	7.234
5.	p/n	pejskařů	10	15.917	3.162	6.725
6.	p/n	chalupářů	12	15.894	3.464	6.725
7.	p/n	vášnivých	7	15.585	2.646	6.376
8.	p/n	sázkařů	5	15.019	2.236	5.817
9.	p/n	zahrádkářů	5	12.695	2.236	3.607
10.	p/n	milovníků	3	11.531	1.731	2.452
11.	p/n	muzikantů	3	10.551	1.731	1.483
12.	p/n	lyžařů	3	10.513	1.731	1.445

SketchEngine (učo+sekundární heslo)

The screenshot displays the SketchEngine web interface. At the top left is the SketchEngine logo. To its right is a search bar with a dropdown arrow and a magnifying glass icon. In the top right corner, it says "words: 7 % / 1,000,". On the left side, there is a navigation menu with the following items: Home, + Create corpus, + WebBootCaT, + Upload TMX or XLS, Parallel corpora, Compare corpora, My jobs, Advanced features (highlighted with a white box), Corpus templates, Sketch grammars, Subcorpus definitions, User groups, and Subscription overview. The main content area is titled "Log in" and "Masaryk University staff and students". Below this, there is a button for "SSO authentication" with a question mark icon. Underneath, the section "External account holders" contains a login form with fields for "User name" (containing "osolsobe") and "Password" (represented by dots). A "Log in" button is located at the bottom right of the form.

SketchEngine

- Nástroj disponuje dalšími funkcemi zpracování jednotek (slovních tvarů/lemmat) v korpusech
- Slovní profily (wordsketches) – gramatická kombinovatelnost slov
- Zobrazování slov na základě podobností ve výskytu (thesaurus)

Funkce Word Sketch

- Umožňuje vytvářet vizualizace frekvenčně uspořádaných gramaticky definovaných relací, do kterých vstupuje klíčové slovo v daném korpusu
- 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)

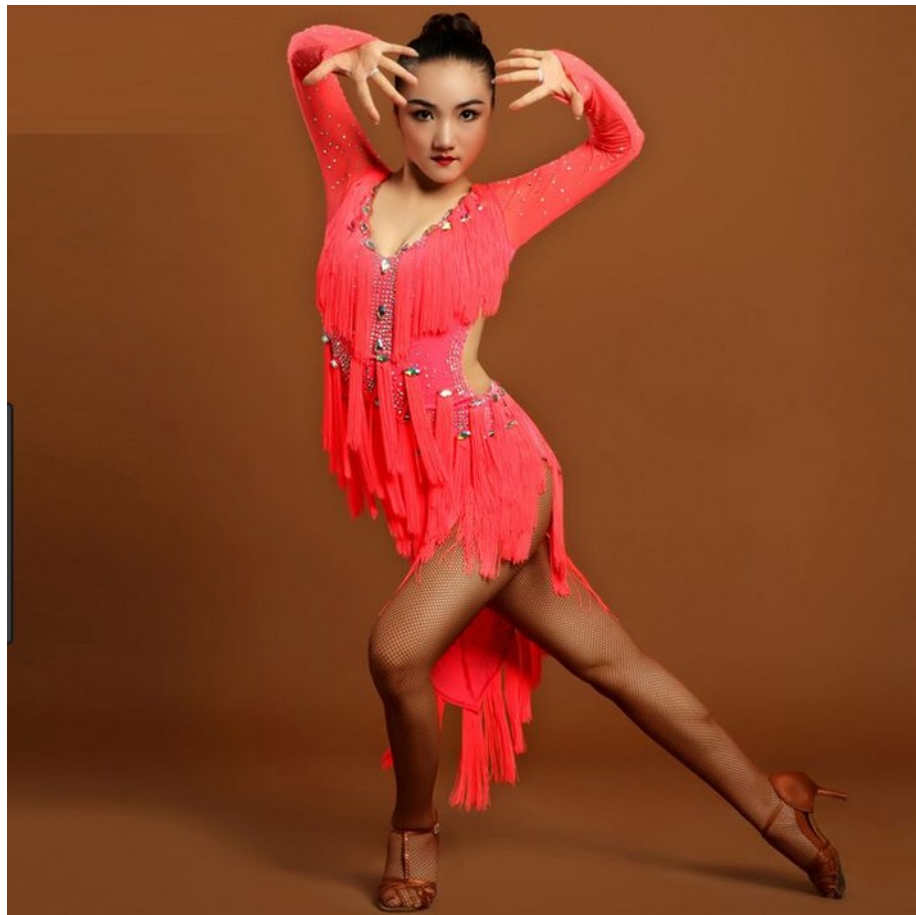
Word sketch *latina*

latina (*noun*)
Czech Web 2012 (czTenTen12 v9) frekvence = 26,837 (5.29 v milionu)

<u>a_modifier</u>		<u>prec_prep</u>		<u>gen_2</u>		<u>coord</u>		<u>is_obj4_of</u>	
	7.87		32.51		8.79		7.88		6.78
myslivecký +	<u>159</u> 7.84	vedle	<u>58</u> 3.23	učebnice +	<u>131</u> 7.95	řečtina +	<u>746</u> 12.75	sekat +	<u>868</u> 11.64
myslivecká latina		vedle latiny ,		učebnice latiny		latiny a řečtiny		sekat latinu	
vulgární	<u>72</u> 7.16	namísto	<u>6</u> 2.69	výslovnost	<u>30</u> 7.86	němčina +	<u>228</u> 9.44	tancovat	<u>25</u> 7.91
vulgární latiny		z +	<u>2,627</u> 2.60	výslovnost latiny		hebrejšтина	<u>45</u> 8.96	tancovala latinu	
ciceronský	<u>9</u> 7.11	z latiny		kódování	<u>42</u> 7.55	staroslověština	<u>33</u> 8.93	učít	<u>24</u> 7.74
středověký +	<u>278</u> 6.71	místo	<u>36</u> 2.39	kódování Latin		francouzština +	<u>114</u> 8.85	se učil latinu	
středověké latiny		kromě	<u>56</u> 2.35	quartier	<u>9</u> 6.96	starořečtina	<u>30</u> 8.83	tančit	<u>31</u> 7.38
akvaristický	<u>14</u> 6.70	kromě latiny a		lingua	<u>7</u> 6.58	italština	<u>50</u> 8.64	tančit latinu	
akvaristické latiny		včetně	<u>43</u> 1.56	znalost +	<u>236</u> 6.43	čeština +	<u>113</u> 8.34	učít +	<u>113</u> 6.20
scholastický	<u>8</u> 6.33	včetně latiny a		znalost latiny		funk	<u>26</u> 8.02	učit latinu	
námořnický	<u>15</u> 6.33	do +	<u>994</u> 1.23	výuka +	<u>148</u> 6.35	funku a latiny		vytlačovat	<u>9</u> 6.11
zkomolený	<u>8</u> 6.29	do latiny		výuky latiny		aramejšтина	<u>17</u> 7.92	vyučovat	<u>25</u> 5.97
hovorový	<u>22</u> 6.19	oproti	<u>9</u> 1.21	profesor +	<u>117</u> 6.26	aramejšтинě a latině		ovládat	<u>98</u> 5.67
hovorová latina		prostřednictvím	<u>13</u> 1.09	profesor latiny		španělština	<u>34</u> 7.68	ovládal latinu	
humanistický	<u>18</u> 5.64	prostřednictvím latiny		profesorka	<u>16</u> 6.16	ruština	<u>35</u> 7.58	vytlačit	<u>11</u> 5.60
archaický	<u>10</u> 5.23	v +	<u>3,495</u> 1.08	profesorka latiny		jazz	<u>34</u> 7.54	vytlačila latinu	
rybářský	<u>61</u> 5.21			znalec	<u>30</u> 6.06	jazzu a latiny		studovat +	<u>105</u> 5.42

Zrada v podobě homonymie (paronymie)

jazz/tancovat/tačit



Funkce Thesaurus (zobrazení podobných slov)

- Na základě porovnání kontextů je vytvořen seznam a vizualizace slov, která mají podobné (gramaticko-lexikální) kontexty

latina (noun) Czech Web 2012 (czTenTen12 v9) frekvence = [26,837](#) (5.29 v milionu)

Lemma	Skóre	Frekvence
francouzština	0.367	31,593
ruština	0.349	25,527
němčina	0.336	62,767
řečtina	0.320	10,714
španělština	0.316	19,214
arabština	0.281	8,834
italština	0.277	11,852
hebrejština	0.270	7,170
čínština	0.266	8,672
poľština	0.262	10,041
slovenština	0.247	14,334
japonština	0.236	7,430
angličtina	0.228	212,669
čeština	0.192	227,991
romština	0.188	4,026
matematika	0.183	88,091
maďarština	0.183	4,228
dějepis	0.183	26,800
portugalština	0.160	3,558
zeměpis	0.156	14,562
švédština	0.151	2,735



Sketch rozdíl (vizualizace kontextu dvojice): *čeština/latina*

- Společné kontexty (bíle)
- Kontexty typické pro každý
- člen dvojice (zeleně
- a červeně podbarvené)

a_modifier	2,113	29,737	0.08	0.13
myslivecký	<u>159</u>	0	7.8	--
ciceronský	<u>9</u>	0	7.1	--
akvaristický	<u>14</u>	0	6.7	--
scholastický	<u>8</u>	0	6.3	--
námořnický	<u>15</u>	0	6.3	--
rybářský	<u>61</u>	0	5.2	--
vulgární	<u>72</u>	<u>25</u>	7.2	4.2
středověký	<u>278</u>	<u>58</u>	6.7	4.1
znalý	<u>22</u>	<u>23</u>	5.1	4.0
zkomolený	<u>8</u>	<u>42</u>	6.3	5.5
vytříbený	<u>12</u>	<u>42</u>	5.2	5.2
humanistický	<u>18</u>	<u>101</u>	5.6	6.4
archaický	<u>10</u>	<u>100</u>	5.2	6.5
hovorový	<u>22</u>	<u>854</u>	6.2	9.5
mluvený	<u>23</u>	<u>581</u>	4.8	8.4
srozumitelný	<u>8</u>	<u>269</u>	2.9	7.1
spisovný	<u>23</u>	<u>5,797</u>	4.7	11.7
demo	0	<u>112</u>	--	6.5
bezchybný	0	<u>149</u>	--	6.6
jadrný	0	<u>125</u>	--	6.9
psaný	0	<u>359</u>	--	7.1
přeložený	0	<u>230</u>	--	7.5
obecný	0	<u>1,973</u>	--	7.7
znakovaný	0	<u>462</u>	--	9.0
lámaný	0	<u>554</u>	--	9.1

Dva pohledy na korpus (lingvista a informatik)

- Nástroje NLP a korpusy
- Konverzní programy, vertikál, tokenizér
- Korpusové manažery
- Automatické analyzátory
- Lingvistické interpretace v korpusech

Vyhledání slovního tvaru *jít*

" Můj partner na rakovinu zemřel , protože se bál	jít	k lékaři . Z vlastní zkušenosti tedy vím , jak
, " uvedl Malášek . Nevyloučil možnost , že může	jít	o vojáka nebo policistu . K druhé loupeži došlo jen
pádu . ^ * S výjimkou případů , kdy má	jít	k veterináři a vy se ji marně snažíte napěchovat do
mrazu . (Lidovci už ztrácejí i Moravu . Zkusí	jít	do měst a zezelenat Jihomoravští lidovci získali v komunálních volbách
patriotům se změna názvu zřejmě líbit nebude , nemusí však	jít	o řešení trvalé . Vše nyní i v budoucnosti závisí
. Tyto potíže se léčí laserem . Po operaci může	jít	pacient hned domů Prof. MUDr. Pavel Kuchynka , CSc. ,
v žádném případě . Je pravda , že jsme mohli	jít	na Švédsko a Kanadu . A nebudu lhát , že
v jedenáctimilionové zemi . ATÉNY „ Nevíme , kam máme	jít	, přijeli jsme lodí z Lesbosu , moje žena je
více Dominiků , než bývá obvyklé . A nemuselo hned	jít	o zapálené hokejové fanoušky . Stačí , že se jméno
unie přijala směrnici o nedovolené podpoře podnikání , která neumožňuje	jít	s nabídkou pod určitou minimální částku , " vysvětlil mluvčí
řadových odborářů , s jakým argumentačním mandátem k němu mají	jít	. Doslova v odborářském šoku jsem po té , co
trasy k nejzajímavějším místům parku . Můžete použít lodě ,	jít	pěšky nebo jet džípem a bahnem se dobrodit až k

Vyhledání lemmatu *jít* (KWIC+lemma+tag)

v Ledči . Pokud dva tři dny pořádně prší ,	jde/jít/VB-S---3P-AA---I	začít pár kilometrů výše proti proudu a zdolat Stvořidla .
to přece dělal baron Prášil ! Americký kaskadér David Smith	jde/jít/VB-S---3P-AA---I	v jeho stopách - o víkendu tenhle bláznivý kousek předvedl
takové se mezi nimi možná občas najdou - , ale	jde/jít/VB-S---3P-AA---I	o spory , jež jdou nutně ruku v ruce s
ale také ceny za téměř třicet tisíc korun . "	Jde/jít/VB-S---3P-AA---I	opravdu o nálož cen kvalitních značek , " řekl pořadatel
příplácet , když chce státní úřad . Zvlášť , když	jde/jít/VB-S---3P-AA---I	o peníze , které tak jako tak od státu předtím
v baru , bowling . Něco , co lidi přiměje	jít/jít/Vf-----A---I	tam , a ne trávit čas jinak . Vesecko nabízí
s onou šaškárnou . Upevňuje se klamně přesvědčení , že	jde/jít/VB-S---3P-AA---I	o replikaci bytostí . Že je to něco jako nesmrtelnost
rání firemního mobilního telefonu a automobilu k soukromým účelům	jdou/jít/VB-P---3P-AA---I	ruku v ruce s rostoucími ekonomickými výsledky společnosti . Zkrátka
filmy promítají . „ Máme rozpis , podle kterého měla	jít/jít/Vf-----A---I	kopie filmu od nás do Sezimova Ústí . Bohužel jsem
byl opilý . Tehdy mi řekl , že když to	nepůjde/jít/VB-S---3F-NA---I	po dobrém , půjde to po zlém , " vypověděla
na Liberec ? No , dlouho to vypadalo , že	jdeme/jít/VB-P---1P-AA---I	na Karlovy Vary , protože Litvínov v posledním kole vyrovnal
zjevně mladší a nezdravě agresivní . " Jsem ženská a	jdou/jít/VB-S---1P-AA---I	z práce . " odpověděla jsem poněkud mimo . "
v koupelně a přistihla je tam Kate , když si	šla/jít/VpFS---3R-AA---I	upravit nalíčení . Pak ale Pippa začala chodit s bankéřem
neexistující exoty a kdo všechno o tom podvůdku věděl a	nešel/jít/VpMS---3R-NA---I	to oznámit na policii . Ztratit paměť je , milí
, protože se jedná o klasické inflační peníze , které	jdou/jít/VB-P---3P-AA---I	v podstatě do stavebních firem . Mnohem lepší podle mne
= hrůza ! Jaké máte zkušenosti s českými silnicemi ?	Nejde/jít/VB-S---3P-NA---I	mi dost dobře do hlavy stav českých dálnic , které
v síti STEP () . HLEdEJtE spoLEčnou cEstu Pokud	jde/jít/VB-S---3P-AA---I	o předkládání projektů na odstraňování starých ekologických zátěží , nejbližší
. Celý dospělý život jsem měla největší přání - aby	šel/jít/VpMS---3R-AA---I	bolševik od válu . A to se mi splnilo .

desambiguace

- Pánové, *nežeňte se*
- Nemluv a *rožni*.
- Jan *je osel*.

Víceznačné tvary

- *nežeňte/(ne)hnat/V*
- *nežeňte/(ne)ženit/V*
- *se/se/P*
- *se/s/R*
- *rožni/rožnit/V*
- *rožni/rozžehnout/V*
- *rožni/rožeň/N*
- *je/být/V*
- *je/on/P*
- *osel/osel/N*
- *osel/osít/V*

hnát/ženit

pracujete za byt a stravu . Ale nevdávejte se a	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	jen kvůli tomu , že chcete uniknout z dor
sexu	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	s učitelkou . Potkáte - li ženu , která krác
Shrnutí těchto vědeckých poznatků : Nejezděte na kole a	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	a nevdávejte) .
legendární Knoflenkou .	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	KDO S KÝM KDE
Takže se mějte a smějte (a	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	do ženění
tvůrci :	<u>Nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	HEJNICE
vážný důvod , proč tam chce jet .	<u>Nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	Pot
takovýchto " žertíků " neustále a iv jiném prostředí .	<u>Nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	ZA SLEVOU , ŽÁDNÁ NENÍ
Uvědomte si , že vaše zásoba energie je omezená .	<u>Nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	Jedn
Udo Pollmer říká : " Vyhodte z bytu televizi a	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	z jedné činnosti do druhé .
ventilaci - či pootevřete e dveře . Když skončíte ,	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	7 . V
to , co vyděláte hned dál investujte .	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	. Svatba je počátkem tukového obalovář
požehnaného Čechova , jenž pravil , cituji : „	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	sprchovat - studenou vodou , ale zůstaň
se mohli vrátit k zaříkadlům , která přivolávají lásku .	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	do nejistých investic
na cizí . Sebastian Roch Chamfort	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	Nezkoušim
Bojíteli se samoty ,	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	, pokud se bojíte osamělosti .
Uvědomte si , že vaše zásoba energie je omezená .	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	"
světě vidět a chovat vlastní dítě .	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	však hned do vážných známostí , nejdřív
Nepodléhejte trendům a	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	! Anton Pavlovič Čechov
vědění je chromá . Věda bez náboženství je slepá .	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	Bůh st
růst zajistí .	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	z jedné činnosti do druhé .
20.00 hodin kino v Hejnicích . Motto filmu zní :	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	7 . V
	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	na sál jen proto , že kamarádi byli . Zept
	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	za detaily . Nechť je vám v myšlenkách i
	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	do velkých zakázek bezhlavě . Získavejte
	<u>nežeňte/hnát/Vi-P---2--N---I se/se/P7--4-----</u>	do ženění . Až na svatební cestě Eddie z

být/on

Není to nic moc , ale Millicent to nepozná .	Je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	. Znova - que voulez-vous ? " Prozpěvovaly kvulevu ,
Titulní strana Šípu , která o premiérovi hlásala , že	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	, či ta , kde trenér Brückner drží před začátkem
a profous považuje za svou povinnost říci mu , že	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	. Mužstvo jde na těžké cvičení do hor , ale
Nebudu dělat to , co řekne Manuel . „ Manuel	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	, Pere ! Osel je zvíře a ty přece nebudeš
mějský jarmark . V roce 1616 následovala hra Dábel	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	, v níž ani dábel nestačí na zlo lidí ,
ocourskovský sedlák osít dvě míry pole krupicí . Kdo	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	, nechť se přihlásí . V říčce , do které
rozesmálo všechny . Mě tedy nejvíc překvapilo , že	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	chytřejší než kůň . Po dotazech , které jsme měli
ly smím použít přejemnělého výrazu - jaký ten druhý	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	. Ale ne tak Evženka . Ta , když vypravuje
de se před drátěným plotem něco pohybovalo . " To	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	! " zajásala Julie a už se k němu rozeběhla
: A co dále , Baltazare , kde hlavním hrdinou	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	nebo u S. Paradžanova : Barva granátového jablka , kde
osla a pro jistotu pod obraz napsal : " Toto	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	. " Bedlivější analýzou vyprávěných anekdot lze zjistit , že
to na tebe narařčil , ty osle ! Román Gaius	<u>je/oni/PPFP4--3----- osel/ose1/NNMS1-----A-----</u>	ukazuje , že Henry Winterfeld psal už před půlstoletím lepší
Trutnovska z obce Mladé Buky , že v ulici Tmavá	je/být/VB-S---3P-AA---I osel/ose1/NNMS1-----A-----	. „ Policisté díky místní znalosti vypátrali jeho majitele a

rožeň/rožnit/rozžehnout

NY . , O. LI . V . Y ,	ROŽNI/rožeň/NNIP7-----A-----	, K. I . Š , Pa . Le .
to vím jistě . Z nářečí znám třeba slovo „	rožni/rožeň/NNIP7-----A-----	“ . Ester Hozáková , 6 let , FrýdekMístek :
k paní Haleové , „ že s těmi hodinami ,	rožni/rožeň/NNIP7-----A-----	a planetárium máte doma malé muzeum . “ „ Jestli
jste naučila vy jeho ? Různé moravismy jako šufánek nebo	rožni/rožeň/NNIP7-----A-----	. Je nevěra tak zajímavá , že o ní teď
nás chodí dívat místní . Je tam obrovský krb s	rožni/rožeň/NNIP7-----A-----	, ve kterém se peče drůbež - od bažanta přes
je jen krátký úryvek z nabídky mas na grilu i	rožni/rožeň/NNIP7-----A-----	, které Farma nabízí . Mezi omáčkami , přílohami a
pro služebnictvo sestoupil do kuchyní , kde skomíral oheň pod	rožni/rožeň/NNIP7-----A-----	bez masa . Na zemi vedle stolu na porcování masa
pražsky , což rozhodně popírám ! Děti říkají fajne ,	rožni/rožeň/NNIP7-----A-----	místo rozsvít , štrampliky místo punčocháče . Opakují slova ,
, ale také teplá kuchyně s grily , rošty a	rožni/rožeň/NNIP7-----A-----	, kde gastronomické poklady před našima očima do zlatova dozrávaly
místnosti Dechem řádky staletí s krápníky sazí z ohňů pod	rožni/rožeň/NNIP7-----A-----	praset či skopců hovoří centrální černá kuchyně , kam by
sklad plný všelikého harampádí a jen mohutné topeniště krbu s	rožni/rožeň/NNIP7-----A-----	a naproti stojící sporák byly důkazem , že stojím v
to pyšný . V Praze pořád používám slovo žufánek nebo	rožni/rožeň/NNIP7-----A-----	. A nerozuměj mi , “ prozradil . Když Okamura
, ale také teplá kuchyně s grily , rošty a	rožni/rožeň/NNIP7-----A-----	, kde gastronomické poklady před našima očima dozlatova dozrávaly .

Filologie a korpusy

- Obecné a specializované korpusy
- Příklady z českého prostředí
- Tvorba vlastního korpusu

ÚČNK <http://ucnk.ff.cuni.cz/cs/>

- akademický projekt 1994
- systematicky mapuje češtinu i další jazyky
- po bezplatné registraci otevřeny všem zájemcům

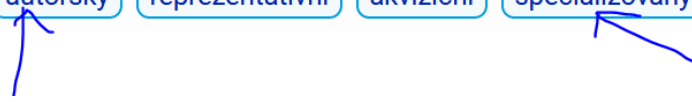
Korpusy ÚČNK

https://kontext.korpus.cz/first_form?corpname=omezeni%2Fsyn2015



Můj seznam | **Všechny korpusy**

Zrušit výběr řada SYN řada ORAL InterCorp synchronní diachronní mluvený psaný webový současná verze
starší verze čeština cizojazyčný paralelní srovnatelný autorský reprezentativní akviziční specializovaný



Specializovaný - příklad

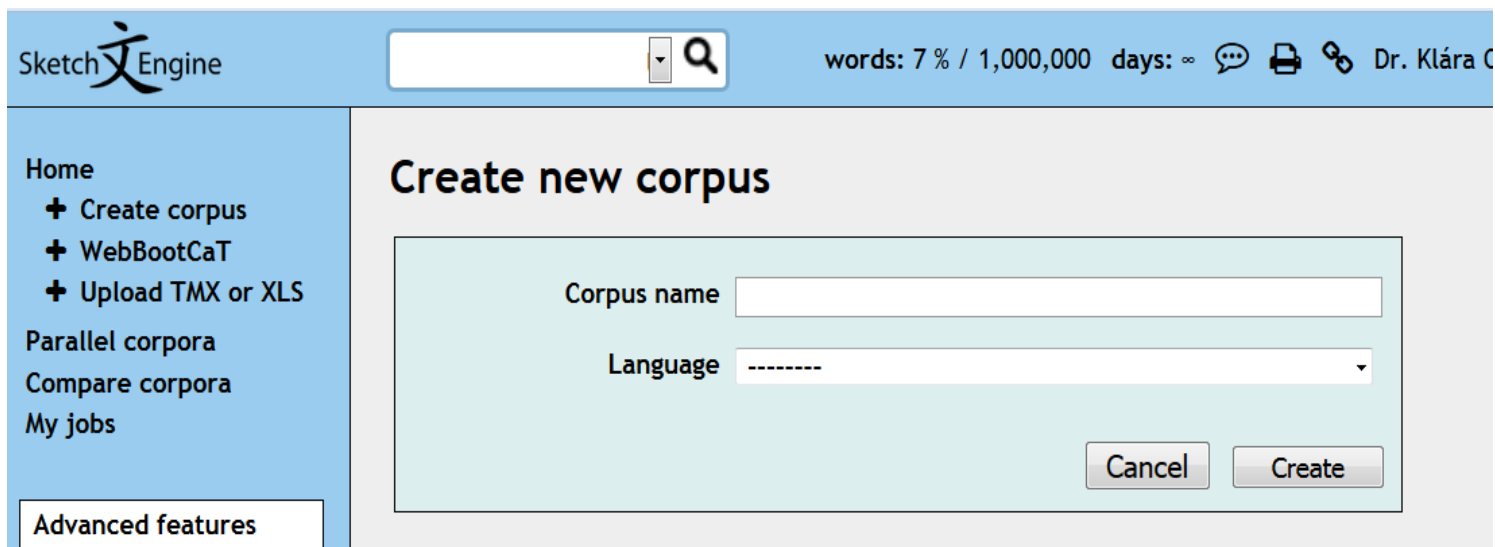
- Korpus Jerome je **jednojazyčný srovnatelný korpus** (*monolingual comparable corpus*) speciálně sestavený pro zkoumání překladové češtiny (tedy textů přeložených do češtiny z jiných jazyků) v porovnání s češtinou nepřekladovou (původní česky psanou).

autorský

- Karel Čapek
- Bohumil Hrabal
- Korespondence Karla Havlíčka Borovského

Tvorba vlastního korpusu - DIY

- <https://nlp.fi.muni.cz/cs/JakVytvoritKorpus1>
- https://ske.fi.muni.cz/auth/create_corpus/



The screenshot shows the Sketch Engine web interface. At the top, there is a search bar and a status bar displaying "words: 7 % / 1,000,000 days: ∞" along with icons for chat, print, and share, and the name "Dr. Klára C". On the left side, a navigation menu includes "Home", "Create corpus", "WebBootCaT", "Upload TMX or XLS", "Parallel corpora", "Compare corpora", and "My jobs". The main content area is titled "Create new corpus" and contains a form with two input fields: "Corpus name" and "Language". Below the form are "Cancel" and "Create" buttons. A box labeled "Advanced features" is visible at the bottom left of the interface.

Výuka jazyků a korpusy

- Metoda DDL (Tim Johnes)
- Žákovské korpusy (Learner Corpora)
- Učebnicové korpusy

Žákovské korpusy

- Texty mluvčích L2
- Výzkum interlanguage
- Zpětná vazba na základě analýzy chyb
- Sledování růstu jazykových kompetencí

Učebnicový korpus Učko

1	Adamovičová, A. – Ivanovová, D. (2007): <i>Basic Czech I, II</i> . Praha: Karolinum.
2	Bischofová, J. et al. (2007): <i>Čeština pro cizince a azylanty B1</i> . Brno: SOZE.
3	Bořilová, P. – Holá, L. (2011): <i>Česky krok za krokem 2</i> . Praha: Akropolis.
4	Cvejnová, J. (2011): <i>Česky, prosím</i> . Praha.
5	Čechová, E. – Remediosová, H. (2005): <i>Chcete mluvit česky? 1, 2</i> Liberec: Harry Putz.
6	Froulíková, L. (2008): <i>Adam a Eva v Českém ráji</i> . Praha: Academia.
7	Hádková, M. (2005): <i>Čeština pro cizince a azylanty A1, A2</i> . Brno: SOZE.
8	Holá, L. (2006): <i>New Czech Step by Step</i> . Praha: Akropolis.
9	Holá, L. (2010): <i>Čeština Express 1, 2</i> . Praha: Akropolis.
10	Kestráňková, M. et al. (2010): <i>Čeština pro cizince B1</i> . Brno: Cpress.
11	Matula, O. (2007): <i>Český den. Kurz českého jazyka pro azylanty navazující na Manuál pro učitele českého jazyka pro cizince bez znalosti latinky</i> . Praha: Člověk v tísni o.p.s., Projekt Varianty.
12	Nekovářová, A. (2006): <i>Čeština pro život - 15 moderních konverzačních témat</i> . Praha: Akropolis.
13	Parolková, O. (2004): <i>Czech for foreigner s 1, 2</i> . Praha: Bohemika.
14	Pintarová, M. – Režková, I. (1995): <i>Communicative Czech. Elementary Czech. Intermediate Czech</i> . Praha: Univerzita Karlova.
15	Rigerová, K. (2000): <i>Czech for Everyone</i> . Praha: K. Rigerová.
16	Štindl, O. (2008): <i>Easy Czech. Elementary</i> . Praha: Akronym.
17	Štindlová, B. (2008): <i>Česky v Česku 1, 2</i> . Praha: Ústav jazykové a odborné přípravy UK: Akropolis.

chodit do/na (porovnání obecného a učebnicového korpusu)

Celkem: 17391 (348 str.)

	Filtr	word	lc [lowercase word]	Frekvence
1.	p/ n	do	školy	21 199
2.	p/ n	do	práce	15 196
3.	p/ n	na	procházky	4 640
4.	p/ n	do	kina	3 661
5.	p/ n	do	kostela	2 940
6.	p/ n	do	posilovny	2 785
7.	p/ n	do	školky	2 693
8.	p/ n	na	fotbal	2 153
9.	p/ n	do	divadla	2 104
10.	p/ n	do	lesa	2 063
11.	p/ n	do	zaměstnání	1 808
12.	p/ n	na	tréninky	1 629
13.	p/ n	do	hospody	1 540
14.	p/ n	na	pivo	1 534
15.	p/ n	na	houby	1 380

	lemma	word	Frekvence
P N	chodit na	procházky	5
P N	chodit do	kina	4
P N	chodit do	školy	3
P N	chodit do	divadla	3
P N	chodit na	oběd	2
P N	chodit na	dlouhé	2
P N	chodit na	balet	2
P N	chodit do	kostela	2
P N	chodit na	různé	1
P N	chodit na	ryby	1
P N	chodit na	lekce	1
P N	chodit na	koncerty	1
P N	chodit na	houby	1
P N	chodit na	hodinu	1
P N	chodit na	dlouhou	1
P N	chodit na	brigádu	1
P N	chodit do	přírody	1
P N	chodit do	práce	1
P N	chodit do	lesa	1
P N	chodit do	kin	1
P N	chodit do	hospodv	1

Závěr

- Historie a současnost – technický pokrok a metodologické přístupy
- Rychlost- spolehlivost – opakovatelnost experimentu
- Zdroje nespolehlivosti
- Co je k dispozici a co si mohu sám udělat
- Na co nezbyl čas

Děkuji vám za pozornost