

ЭКСПЕРИМЕНТ С СОВМЕСТИМОЙ СКЕТЧ-ГРАММАТИКОЙ

Владимир Бенко

vladob@juls.savba.sk

Словацкая академия наук
Институт языкознания им. Людовита Штура
Братислава

«Корпусная лингвистика»
Санкт-Петербург, июнь 2013

Word Sketches

Word sketches are one-page automatic, corpus-based summaries of a word's grammatical and collocational behaviour.

[Kilgarriff et al., Euralex 2004]

Word Sketches

*Word sketches are **one-page** automatic, corpus-based summaries of a word's **grammatical** and collocational behaviour.*

[Kilgarriff et al., Euralex 2004]

Word Sketches

*Word sketches are **one-page** automatic, corpus-based summaries of a word's **grammatical** and collocational behaviour.*

[Kilgarriff et al., Euralex 2004]

*Word sketch ... «эскиз слова» ...
коллокационный профиль*

Коллокационные профили

- Корпус текстов с **морфологической** (или морфосинтактической) **разметкой**
- Программа **Word Sketch Engine**
- Набор правил (**скетч-грамматика**) созданных с учетом нужд пользователей корпуса

Скетч-грамматика

- Основана на языке CQL
(Corpus Query Language)

[tag="A.*"] []{0,3} [tag="N.*"]

Скетч-грамматика

- Основана на языке CQL
(Corpus Query Language)

2: [tag="A.*"] [] {0,3} **1:** [tag="N.*"]

Скетч-грамматика

- Основана на языке CQL
(Corpus Query Language)

=a_modifier

2: [tag="A.*"] [] {0,3} 1: [tag="N.*"]

Скетч-грамматика

- **Ключевое слово:**
словарь

<u>a_modifier</u>	<u>260656</u>	<u>0.3</u>
толковый	<u>43173</u>	11.95
энциклопедический	<u>35926</u>	11.8
англо-русский	<u>4412</u>	9.05
немецко-русский	<u>3250</u>	8.62
этимологический	<u>3213</u>	8.61
русско-немецкий	<u>3159</u>	8.6
русско-английский	<u>3147</u>	8.59
биографический	<u>3347</u>	8.48
Англо-русский	<u>2411</u>	8.22
орфографический	<u>2271</u>	8.07
терминологический	<u>2184</u>	8.01
краткий	<u>4906</u>	7.72
философский	<u>3550</u>	7.58
толково- фразеологический	<u>1074</u>	7.07

Скетч-грамматика

Наш подход

- **Названия правил не означают синтаксические отношения, а коллокационные**
- **Позиция ключевого слова и коллоката показана явно**
- **Часть речи ключевого слова не определяется**

Скетч-грамматика

Наш подход

$A_v X/X A_v$

$V_b X/X V_b$

$A_j X, X A_j$

$S_b X, X S_b$

$Z X, X Z$

$P_p X, X P_p$

$X/Y C_j X/Y$

$P_p X Y, X P_p Y$

$P_p Y X, Y P_p X$

$A_j(X), V_b(X)...$

$nom(X), pl(x)...$

Применение скетчей в Институте языкознания им. ЛШ

- **Словарь современного словацкого языка (8-томный, 2 тома уже опубликованы)**
- **Коллокационный словарь (1-ый том – Коллокации с существительными именами – подготовлен к печати)**
- **Чешско-словацкий словарь (2-томный, работы пока в начале)**

Параллельные скетчи

jazyk (cs) / jazyk (sk)

<u>Aj X</u>	<u>35946</u>	<u>-2.1</u>		<u>Aj X</u>	<u>116650</u>	<u>1.2</u>
cizí	<u>3970</u>	8.53		slovenský	<u>15377</u>	6.14
český	<u>3863</u>	4.16		cudzí	<u>14845</u>	9.42
anglický	<u>2251</u>	7.72		anglický	<u>11514</u>	9.3
programovací	<u>1760</u>	9.77		štátny	<u>4449</u>	6.04
německý	<u>1241</u>	5.21		nemecký	<u>4297</u>	7.16
jiný	<u>901</u>	3.09		spisovný	<u>3354</u>	9.09
světový	<u>878</u>	4.48		materinský	<u>2707</u>	9.06
další	<u>616</u>	1.76		vyučovací	<u>2377</u>	8.28
úřední	<u>564</u>	6.94		úradný	<u>1862</u>	7.78
mateřský	<u>543</u>	6.25		maďarský	<u>1745</u>	6.23
rodný	<u>486</u>	6.72		programovací	<u>1707</u>	8.54
zlý	<u>440</u>	5.01		český	<u>1690</u>	5.12
různý	<u>407</u>	3.03		svetový	<u>1673</u>	5.01
sněhový	<u>404</u>	6.81		ruský	<u>1604</u>	6.13
ruský	<u>386</u>	4.01		rodný	<u>1473</u>	7.03
spisovný	<u>337</u>	7.82		slovanský	<u>1409</u>	7.46

Параллельные скетчи

jazyk (cs) / jazyk (sk)

<u>Vb X/X Vb</u>	<u>29725</u>	<u>-0.4</u>		<u>Vb X/X Vb</u>	<u>66131</u>	<u>0.1</u>
být	<u>6112</u>	0.99	—	byť	<u>16560</u>	1.11
mluvit	<u>954</u>	5.22	—	ovládať	<u>3010</u>	7.9
mít	<u>862</u>	0.39	—	mať	<u>2958</u>	1.05
učit	<u>550</u>	5.98	—	hovoríť	<u>2505</u>	4.21
mocet	<u>500</u>	0.36	—	používať	<u>1734</u>	4.68
ovládat	<u>456</u>	6.3	—	učíť	<u>1719</u>	5.99
používat	<u>442</u>	3.92	—	môcť	<u>1283</u>	0.78
umět	<u>407</u>	4.56	—	naučiť	<u>1111</u>	5.34
začít	<u>405</u>	2.34	—	vedieť	<u>861</u>	1.44
naučit	<u>384</u>	5.57	—	musieť	<u>612</u>	0.97
hovořit	<u>364</u>	4.55	—	vyučovať	<u>516</u>	6.54
tvrdit	<u>272</u>	3.35	—	stať	<u>481</u>	1.57
muset	<u>266</u>	0.6	—	študovať	<u>449</u>	4.95
znát	<u>233</u>	3.48	—	rozprávať	<u>435</u>	4.27
tvořit	<u>189</u>	3.2	—	začať	<u>405</u>	1.27
studovat	<u>171</u>	4.58	—	chcieť	<u>388</u>	0.1

Параллельные скетчи

jazyk (cs) / jazyk (sk)

<u>Sb X</u>	<u>20405</u>	<u>-0.4</u>		<u>Sb X</u>	<u>57642</u>	<u>0.2</u>
výuka	<u>1488</u>	8.0		znalosť	<u>3291</u>	8.13
znalosť	<u>1410</u>	7.39		výučba	<u>2075</u>	8.13
studium	<u>442</u>	4.82		slovník	<u>1819</u>	7.95
mluvení	<u>351</u>	8.6		vyučovanie	<u>1809</u>	7.56
slovník	<u>310</u>	6.43		používanie	<u>1666</u>	6.41
špička	<u>306</u>	5.72		kurz	<u>968</u>	5.07
dar	<u>269</u>	5.47		katedra	<u>852</u>	6.53
podpora	<u>239</u>	2.93		učiteľ	<u>821</u>	5.1
používání	<u>206</u>	5.23		ovládanie	<u>776</u>	6.18
učitel	<u>196</u>	4.52		základ	<u>649</u>	2.82
základ	<u>172</u>	2.29		štúdium	<u>632</u>	5.52
pře	<u>151</u>	-0.15		oblasť	<u>618</u>	2.2
katedra	<u>146</u>	5.85		hodina	<u>593</u>	2.86
překladač	<u>145</u>	6.96		používateľ	<u>582</u>	5.49
možnost	<u>143</u>	1.35		učenie	<u>569</u>	5.52
ústav	<u>136</u>	3.17		štúdio	<u>567</u>	4.39

Русский веб-корпус

- Веб-краулер **SpiderLing**, оптимизированный для скачивания текстов (модуль для определения языка документа, удаление шаблонов)
- Морфологическая разметка: **TreeTagger** & **MULTEXT-East** тегсет
- Дедупликация документов методом фингерпринтов (**дубликаты удалены**)

Русский веб-корпус

- **Сегментация предложений**
- **Дедупликация предложений
(дубликаты обозначены)**
- **Унификация тегов для пунктуации**
- **Перенос словацкой скетч-грамматики
на русский тегсет MULTEXT-East**
- **Обработка с помощью WSE**

Русский веб-корпус

- **7,9 Гбайт сырого текста**
- **18,9 Гбайт текста в вертикальном формате (после морфологической разметки)**
- **604 тысяч документов**
- **41,3 млн предложений**
- **699 млн токенов (словоформ, включая пунктуацию)**

Русский веб-корпус

Дедупликация

- Удалено 15,9 тысяч документов (2,7%)
- Выявлено дублетами и обозначено 11,1 млн предложений (27,6%)
содержащих 155 млн токенов (22,1%)

Параллельные скетчи

язык (ru) / jazyk (sk)

Aj X	124234	-0.2		Aj X	116650	1.2
русский	<u>27343</u>	8.14		slovenský	<u>15377</u>	6.14
английский	<u>8991</u>	8.93		cudzí	<u>14845</u>	9.42
иностранн	<u>5202</u>	7.89		anglický	<u>11514</u>	9.3
родной	<u>4013</u>	8.11		štátny	<u>4449</u>	6.04
общий	<u>3003</u>	5.48		nemecký	<u>4297</u>	7.16
украинский	<u>2718</u>	6.99		spisovný	<u>3354</u>	9.09
государственный	<u>1999</u>	4.92		materinský	<u>2707</u>	9.06
литературный	<u>1996</u>	7.6		vyučovací	<u>2377</u>	8.28
немецкий	<u>1978</u>	6.57		úradný	<u>1862</u>	7.78
современный	<u>1880</u>	4.98		maďarský	<u>1745</u>	6.23
французский	<u>1750</u>	6.75		programovací	<u>1707</u>	8.54
разный	<u>1494</u>	4.52		český	<u>1690</u>	5.12
греческий	<u>1187</u>	6.83		svetový	<u>1673</u>	5.01
арабский	<u>1040</u>	6.86		ruský	<u>1604</u>	6.13
славянский	<u>918</u>	6.89		rodný	<u>1473</u>	7.03
официальный	<u>915</u>	5.14		slovanský	<u>1409</u>	7.46

Параллельные скетчи

язык (ru) / jazyk (sk)

Sb X	63192	-0.0		Sb X	57642	0.2
изучение	<u>2699</u>	7.58		znalosť	<u>3291</u>	8.13
знание	<u>2323</u>	6.56		výučba	<u>2075</u>	8.13
носитель	<u>1309</u>	7.71		slovník	<u>1819</u>	7.95
развитие	<u>943</u>	3.62		vyučovanie	<u>1809</u>	7.56
обучение	<u>772</u>	5.52		používanie	<u>1666</u>	6.41
владение	<u>699</u>	6.91		kurz	<u>968</u>	5.07
использование	<u>687</u>	4.1		katedra	<u>852</u>	6.53
словарь	<u>587</u>	7.05		učiteľ	<u>821</u>	5.1
уровень	<u>582</u>	3.0		ovládanie	<u>776</u>	6.18
учитель	<u>542</u>	5.34		základ	<u>649</u>	2.82
преподаватель	<u>534</u>	6.25		štúdium	<u>632</u>	5.52
статус	<u>529</u>	5.36		oblasť	<u>618</u>	2.2
история	<u>483</u>	2.96		hodina	<u>593</u>	2.86
язык	<u>476</u>	3.26		používateľ	<u>582</u>	5.49
преподавание	<u>444</u>	6.92		učenie	<u>569</u>	5.52
слово	<u>442</u>	2.13		štúdio	<u>567</u>	4.39

Параллельные скетчи

язык (ru) / jazyk (sk)

Vb X/X Vb	71035	-0.0		Vb X/X Vb	66131	0.1
быть	<u>4816</u>	1.81		byť	<u>16560</u>	1.11
знать	<u>2806</u>	4.83		ovládať	<u>3010</u>	7.9
говорить	<u>2382</u>	4.22		mať	<u>2958</u>	1.05
являться	<u>2025</u>	3.74		hovoríť	<u>2505</u>	4.21
владеть	<u>1642</u>	7.64		používať	<u>1734</u>	4.68
изучать	<u>1554</u>	7.41		učíť	<u>1719</u>	5.99
найти	<u>1490</u>	4.82		môcť	<u>1283</u>	0.78
мочь	<u>1407</u>	1.65		naučiť	<u>1111</u>	5.34
учить	<u>1135</u>	6.86		vedieť	<u>861</u>	1.44
стать	<u>974</u>	2.42		musieť	<u>612</u>	0.97
иметь	<u>870</u>	2.44		vyučovať	<u>516</u>	6.54
выучить	<u>767</u>	7.96		stať	<u>481</u>	1.57
находить	<u>720</u>	5.72		študovať	<u>449</u>	4.95
означать	<u>715</u>	5.29		rozprávať	<u>435</u>	4.27
понимать	<u>699</u>	4.18		začať	<u>405</u>	1.27
использовать	<u>640</u>	3.42		chcieť	<u>388</u>	0.1

Параллельные скетчи

язык (ru) / jazyk (sk)

<u>X/Y</u> Cj <u>X/Y</u>	<u>12694</u>	<u>-0.1</u>		<u>X/Y</u> Cj <u>X/Y</u>	<u>18183</u>	<u>0.4</u>
язык	<u>1375</u>	4.82		literatúra	<u>2931</u>	7.65
культура	<u>1305</u>	5.3		jazyk	<u>1535</u>	4.97
литература	<u>1014</u>	6.31		kultúra	<u>1361</u>	5.25
мышление	<u>228</u>	5.03		slovenčina	<u>513</u>	5.87
губа	<u>215</u>	5.31		matematika	<u>456</u>	6.39
меньшинство	<u>192</u>	5.99		angličtina	<u>236</u>	4.91
речь	<u>177</u>	2.48		reč	<u>179</u>	3.37
стиль	<u>170</u>	3.29		štýl	<u>175</u>	2.91
математика	<u>161</u>	5.62		národ	<u>174</u>	2.71
народ	<u>158</u>	1.38		predmet	<u>130</u>	1.93
диалект	<u>131</u>	6.98		pera	<u>120</u>	4.53
история	<u>130</u>	1.08		komunikácia	<u>119</u>	2.0
система	<u>115</u>	0.14		nárečie	<u>112</u>	5.74
средство	<u>111</u>	0.77		výchova	<u>107</u>	2.57
религия	<u>96</u>	3.06		ústa	<u>101</u>	3.13
обычай	<u>94</u>	4.59		slovo	<u>100</u>	0.09

Заключение

- Программа **SpiderLing** является очень эффективным и удобным средством для создания веб-корпусов

Заключение

- Программа SpiderLing является очень эффективным и удобным средством для создания веб-корпусов
- Создание **совместимых скетч-грамматик** осуществимо

Заключение

- Программа SpiderLing является очень эффективным и удобным средством для создания веб-корпусов
- Создание совместимых скетч-грамматик осуществимо
- Параллелизм между скетчами наблюдается также на **русском** и даже на **французском** языках

Заключение

- Программа SpiderLing является очень эффективным и удобным средством для создания веб-корпусов
- Создание совместимых скетч-грамматик осуществимо
- Параллелизм между скетчами наблюдается также на русском и даже на французском языках
и **это любопытно**, не правда ли?

Песочница

<http://sketch.juls.savba.sk/sandbox>

ЛОГИН: [spb](#)

пароль: [spb2013,,](#)

будет действителен по [14 июля 2013](#)

e-mail для справок

vladob@juls.savba.sk