

# Sketching Words in Comparable Corpora

**Vladimír Benko**

[vladob@juls.savba.sk](mailto:vladob@juls.savba.sk)

Slovak Academy of Sciences  
Ľ. Štúr Institute of Linguistics

Comenius University in Bratislava  
UNESCO Chair for Translation Studies

**SKEW 4**  
Tallinn, 16 October 2013

# “Parallel” sketches

## jazyk (cs) / jazyk (sk) “language”

<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

# “Parallel” sketches

## jazyk (cs) / jazyk (sk) “language”

<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		hovoriť	<u>2505</u>	4.21
môct	<u>500</u>	0.36		používať	<u>1734</u>	4.68
ovládat	<u>456</u>	6.3		učit'	<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

# Aranea

**A family of (comparable) web corpora**

**Motivation: available corpora**

- **Do not cover all languages needed**
- **Are of different sizes**
- **Are mostly too big for classroom use**
- **Sketch grammars are too different**

# Aranea

***Araneum*** (pl. *aranaea*, n.) ... Latin expression for (cob)web

- **Slovak-Centric** (languages spoken and/or taught in Slovakia and the neighbouring countries)
- Crawled and pre-processed by *SpiderLing* at (approximately) the same time
- Language-independent filtration by the same tools
- Language-dependent filtration by the same methodology

# Aranea

- Compatible tokenization strategy
- PoS-tagged by (possibly) free tools (*Tree Tagger*, etc.)
- Sentence-segmented
- Sentence-level deduplicated, duplicate sentences marked
- Word sketches with **compatible sketch grammars**
- “Language-neutral” (Latin) names denoting the language and size

# Aranea

**Four sizes for each corpus planned**

**Maius** (greater) ... basic version, approx. 1 billion tokens

**Minus** (smaller) ... 10 % sample of Maius (for teaching purposes)

**Minumum** (minimal) ... 1 % sample of Maius (not accessible by general users, used for toolchain and sketch grammar experiments)

**Maximum** (maximal) ... as much as we can get

# Aranea

**Four Aranea family members available by now**

**Araneum Russicum** Maius & Minus (Russian Web, Version 13.10)

**Araneum Francogallicum** Maius & Minus (French Web, Version 13.10)

**Araneum Germanicum** Maius & Minus (German Web, Version 13.10)

**Araneum Hispanicum** Maius & Minus (Spanish Web, Version 13.11, will be published soon)



# Compatible sketch grammars

- **Common set of rules** for all languages
- **Fixed order of tables** in word sketches
- Rule names represent **collocational relationships** (i.e. not syntactic)
- **Syntactic functions** of keywords and/or collocates **not indicated** (e.g., for nouns, we do not speak about “subjects”, “objects”, or “modifiers”, we just indicate the “left-hand” and “right-hand” collocates)

# Compatible sketch grammars

- **Word class (PoS) of keyword is not indicated** (each rule works for any PoS)
- **Recall preferred over precision** (some tables may not be relevant for all word classes, output may contain lots of noise)

# Compatible sketch grammars

## Symbols in rule names

**X** ... keyword (of any PoS, except for punctuation)

**Y** ... collocate of any PoS, except for conjunction, preposition and punctuation

**Z** ... collocate of PoS not covered by explicit rules (“catch all” rules)

# Compatible sketch grammars

## Symbols in rule names

**Nn, Vb, Aj, Av, Pp, Cj ... collocate of indicated PoS**

***Aa(X), aaa(X) ... (in unary rules) PoS categories and subcategories of keyword***

# Compatible sketch grammars

## Summary of rule names

### Binary rules

**Nn X; X Nn; Aj X; X Aj; Pp X; X Pp; Z X; X Z ... side-sensitive rules**

**Vb X/X Vb; Av X/X Av ... side-insensitive rules**

### Symmetric rules

**X/Y, X/Y; X/Y Cj X/Y ... coordination**

# Compatible sketch grammars

## Summary of rule names

### Trinary rules

**Y Pp X; Y Pp X; Pp Y X; Pp X Y ...  
prepositions**

### Unary rules

**Ar(X); Nn(X); ... PoS categories**

**msc(X); fem(X); ... PoS subcategories**

**No dual rules**

# “Parallel” sketches

## vin (fr) / Wein (de) “wine”

<u>X</u> <u>Aj</u>	<u>22268</u>	<u>-0.4</u>		<u>Aj</u> <u>X</u>	<u>24029</u>	<u>-0.9</u>
blanc	<u>3016</u>	5.78		gut	<u>2555</u>	2.99
rouge	<u>2301</u>	5.81		deutsch	<u>475</u>	1.92
chaud	<u>634</u>	4.11		groß	<u>387</u>	0.65
sec	<u>454</u>	4.38		neu	<u>347</u>	0.26
issu	<u>433</u>	3.36		passend	<u>329</u>	3.31
biologique	<u>424</u>	3.54		edel	<u>297</u>	4.35
français	<u>375</u>	1.37		erlesen	<u>282</u>	4.71
grand	<u>318</u>	-0.1		alt	<u>277</u>	1.63
bon	<u>311</u>	0.4		hervorragend	<u>274</u>	3.61
doux	<u>302</u>	3.18		verschieden	<u>245</u>	1.16
jaune	<u>292</u>	3.58		jung	<u>237</u>	1.96
nouveau	<u>287</u>	-0.29		hochwertig	<u>218</u>	3.08
rosé	<u>260</u>	4.05		erst	<u>207</u>	-0.3
fin	<u>237</u>	0.51		trocken	<u>201</u>	3.54
naturel	<u>211</u>	1.77		italienisch	<u>184</u>	3.44
autre	<u>187</u>	-1.43		spanisch	<u>178</u>	3.58

# “Parallel” sketches

## vélo (fr) / Rad (de) “bike”

<u>X Aj</u>	<u>6994</u>	<u>-0.2</u>		<u>Aj X</u>	<u>7120</u>	<u>-0.5</u>
électrique	<u>833</u>	5.02		gut	<u>261</u>	-0.28
elliptique	<u>357</u>	5.24		neu	<u>243</u>	-0.23
libre	<u>214</u>	2.16		eigen	<u>234</u>	0.83
autre	<u>127</u>	-1.97		normal	<u>136</u>	2.5
urbain	<u>115</u>	2.0		alt	<u>116</u>	0.45
vert	<u>112</u>	1.61		erst	<u>113</u>	-1.15
petit	<u>98</u>	-1.29		hochwertig	<u>106</u>	2.32
bon	<u>94</u>	-1.29		einfach	<u>86</u>	-0.59
disponible	<u>88</u>	0.81		sicher	<u>85</u>	0.7
classique	<u>86</u>	1.48		elektrisch	<u>82</u>	2.75
grand	<u>72</u>	-2.22		richtig	<u>75</u>	0.12
pliant	<u>70</u>	2.93		klein	<u>75</u>	-0.85
unique	<u>56</u>	0.39		passend	<u>71</u>	1.31
même	<u>55</u>	-2.88		herkömmlich	<u>68</u>	2.62
neuf	<u>54</u>	1.04		groß	<u>66</u>	-1.87
couché	<u>52</u>	2.49		kostenlos	<u>63</u>	0.91



# “Parallel” sketches

américain (fr) / американский (ru) “American”

<u>Sb X</u>	<u>122056</u>	<u>-0.0</u>		<u>X Sb</u>	<u>118159</u>	<u>-0.1</u>
président	<u>2063</u>	5.38		компания	<u>2942</u>	4.83
gouvernement	<u>1648</u>	4.6		ученый	<u>1676</u>	6.38
armée	<u>1543</u>	6.03		президент	<u>1186</u>	4.63
société	<u>1510</u>	3.83		рынок	<u>1143</u>	4.22
dollar	<u>1206</u>	5.87		доллар	<u>1008</u>	5.25
marché	<u>1136</u>	3.82		экономика	<u>1002</u>	5.19
soldat	<u>874</u>	5.92		войско	<u>908</u>	5.41
série	<u>837</u>	4.42		общество	<u>875</u>	4.13
économie	<u>829</u>	4.32		система	<u>773</u>	2.85
continent	<u>828</u>	6.19		правительство	<u>728</u>	4.12
film	<u>812</u>	3.67		власть	<u>728</u>	3.13
université	<u>810</u>	4.0		армия	<u>716</u>	4.65
entreprise	<u>801</u>	2.5		политика	<u>712</u>	4.31
autorité	<u>800</u>	4.51		штат	<u>705</u>	5.54
politique	<u>784</u>	2.48		военный	<u>679</u>	4.02
Etat	<u>774</u>	4.23		солдат	<u>664</u>	5.73

# “Parallel” sketches

rouge (fr) / красный (ru) “red”

X/Y Cj X/Y	19498	0.0		X/Y Cj X/Y	6332	-0.1
blanc	<a href="#">2700</a>	5.78		белый	<a href="#">1099</a>	5.41
noir	<a href="#">2088</a>	5.37		черный	<a href="#">594</a>	5.05
bleu	<a href="#">1992</a>	6.37		синий	<a href="#">468</a>	6.83
vert	<a href="#">1233</a>	5.12		зеленый	<a href="#">386</a>	5.47
jaune	<a href="#">1047</a>	6.43		желтый	<a href="#">361</a>	6.48
rouge	<a href="#">387</a>	3.57		красный	<a href="#">292</a>	3.78
rose	<a href="#">346</a>	4.63		розовый	<a href="#">209</a>	6.02
orange	<a href="#">336</a>	5.99		оранжевый	<a href="#">171</a>	6.24
rosé	<a href="#">301</a>	7.67		коричневый	<a href="#">94</a>	5.64
autre	<a href="#">266</a>	-0.7		голубой	<a href="#">91</a>	4.51
petit	<a href="#">228</a>	-0.31		золотой	<a href="#">69</a>	2.71
gris	<a href="#">180</a>	4.15		серый	<a href="#">65</a>	3.66
violet	<a href="#">168</a>	5.71		фиолетовый	<a href="#">62</a>	5.99
or	<a href="#">133</a>	2.32		советский	<a href="#">45</a>	0.03
brun	<a href="#">120</a>	4.12		темный	<a href="#">29</a>	1.78
grand	<a href="#">95</a>	-2.02		бордовый	<a href="#">28</a>	6.36

# Compatible sketch grammars

## Corpora with compatible sketch grammars

- **supplement to** (or instead of non-existent) **parallel corpora**
- **teaching** foreign languages & translation studies
- **contrastive** linguistic **studies**
- collocationally-based sketch grammar as **general-purpose grammar** and as a **convenient starting point for developing language-, tagset- & purpose-specific grammar**