

ЭКСПЕРИМЕНТ СО СОВМЕСТИМОЙ СКЕТЧ-ГРАММАТИКОЙ

COMPATIBLE SKETCH GRAMMAR EXPERIMENT

Аннотация. Наша статья описывает эксперимент переноса скетч-грамматики, первоначально созданной для словацких корпусов, на русский корпус, чтобы определить степень параллелизма между двумя скетчами. В рамках эксперимента мы создали и аннотировали русский веб-корпус среднего размера и обработали его с помощью системы Word Sketch Engine.

Abstract. Our paper describes an experiment of porting a sketch grammar originally created for Slovak corpora to a Russian corpus to determine the extent of parallelism in the resulting word sketches. In the framework of the experiment, we created and PoS-tagged a medium-size Russian corpus that was subsequently processed by the Word Sketch Engine.

1. Introduction

*Word sketches are one-page automatic, corpus-based summaries of a word's grammatical and collocational behaviour*¹. To create word sketches we need a PoS-tagged corpus, the Word Sketch Engine (WSE) software and a set of rules – the sketch grammar.

In our Institute, the WSE has been extensively used since autumn 2007 with several Slovak and Czech corpora. These corpora serve as a source of lexical evidence for our monolingual and bilingual lexicographic projects², as well as for other linguistic research activities. The following text will show the results of an experiment to include a new language – Russian – into our set of corpora.

¹ Kilgarriff A., Rychly P., Smrz P., Tugwell D. The Sketch Engine. Proc. EURALEX 2004, Lorient, France; 2004. P. 105–116.

² Benko V. Optimizing Word Sketches for a large-scale lexicographic project. http://videlectures.net/korpusi2010_benko_ows

2. The Corpus

To create a test web corpus, we decided to use the recently released open-source tool SpiderLing³ created at Masaryk University in Brno. SpiderLing has been optimized for collecting textual data from the web and contains an integrated language recognition module and a tool for boilerplate removal (jusText)⁴. The 100 input seed URLs were collected by means of Google. SpiderLing was started with three parallel processes (on a quad-core server) and, during its two runs with a total crawling time of just 14 hours, a total of 7.9 GB of raw text data was obtained.

3. PoS Tagging

The PoS tagging was performed by means of Helmut Schmid's TreeTagger⁵ with the Russian parameter file trained by Serge Sharoff using the MULTEXT-East tagset⁶. This step took 39 hours to complete. The resulting corpus was subsequently filtered to fix some errors in tokenization of punctuation, to segment the sentences, and to de-duplicate it on the document level. After this processing, the corpus contained 799 million tokens in 604 thousand documents. As the last processing step, the corpus was de-duplicated at the sentence level using the fingerprint method. The duplicate sentences were not deleted from the corpus but rather they were marked by a special flag to exclude them from the WSE processing. In total, there were 27.6% duplicate sentences containing 22.1% of the corpus tokens.

Before starting work on the sketch grammar, we decided to unify the tags for punctuation as TreeTagger assigned three different tags for this purpose.

³ *Suchomel V., Pomikálek J.* Efficient Web Crawling for Large Text Corpora. 7th Web as Corpus Workshop (WAC-7), Lyon, France; April 2012

⁴ *Pomikálek J.* 2011. Removing Boilerplate and Duplicate Content from Web Corpora. Ph.D. thesis, Masaryk University, 2011.

⁵ <http://www.ims.uni-stuttgart.de/forschung/ressourcen/werkzeuge/treetagger.html>

⁶ <http://corpus.leeds.ac.uk/mocky>

4. The Sketch Grammar

The sketch grammar used in our WSE installation differs from most “traditional” grammars for corpora stored at sketchengine.co.uk⁷ in several aspects:

- The rule names are not motivated syntactically (i.e., they do not indicate the syntactic relationship between the keyword and the collocate) but rather collocationally
- The right-hand or left-hand position of the collocate towards the keyword is indicated explicitly in the rule name
- The keyword’s PoS in the rule is not specified, i.e., it covers any PoS
- Recall is preferred over precision

The object names within the rules are governed by the following rules:

- The keyword is denoted by the *X* symbol
- The keyword’s grammatical attributes (mostly in unary rules) are indicated by lowercase abbreviation, e.g., *gen(X)* indicates the genitive case of keyword
- The collocate’s PoS is indicated by an abbreviation with a leading capital letter, e.g., *Aj X* indicates a left-hand adjective collocate
- *Y* indicates a collocate that is from any PoS category
- *Z* indicates a collocate from any PoS category not covered by the other «explicit» rules

This way of writing sketch grammar has (against the traditional one) several advantages that can be conveniently utilized within the lexicographic projects:

- The rules symmetrically cover collocational relationships among all word classes, i.e., not only those where a direct syntactic relationship can be seen
- The rule names are more user-friendly
- The rules also cover situations where an incorrect morphological tag has been assigned to the keyword

⁷ <http://trac.sketchengine.co.uk/wiki/SkE/CorpusQuerying>

5. Rule Name Summary

The core of our grammar consists of rules covering four basic autosemantic word classes. Taking into account our experience with early versions of the grammar, the rules for verbs (*Vb X/X Vb*) and adverbs (*Av X/X Av*) do not distinguish the left and right position of the respective collocate.

For nouns («substantives»), two separate rules take into account the position of the collocate (*Sb X, X Sb*). Similar situations can be found with adjectives (*Aj X, X Aj*) and prepositions (*Pp X, X Pp*).

The «catch all» rules for the remaining word classes (*Z X, X Z*) cover mostly numerals and pronouns, as well as some synsemantic word classes.

The remaining two binary (symmetric) rules map the relationship of coordination, i.e., the situation where a keyword and a collocate with an identical morphological tag are separated by a comma (*X/Y, X/Y*) or a conjunction (*X/Y Cj X/Y*).

The four trinary rules cover relationships among a keyword, collocate, and preposition in different positions (*Pp Y X, Pp X Y, Y Pp X*, and *X Pp Y*).

Our set of rules is complemented by unary rules showing the frequency distribution of the keyword's forms according to grammatical categories.

6. Compatible Sketch Grammars

If sketch grammars for corpora of different languages are written within the same paradigm, the resulting word sketches can be conveniently used in bilingual lexicography as a supplement to (or a partial replacement of non-existent) parallel corpora. This method is used in creating the new Czech-Slovak dictionary currently being compiled at our Institute. As Czech and Slovak are both linguistically and culturally closely related languages, the frequencies of collocates for a keyword are mostly very similar. Fig. 1 shows an example of such a “parallel” sketch derived from two comparable web corpora for the keyword *jazyk* «language». The second columns in both tables

show the raw frequency and the third ones the statistical value of salience⁸.

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

Fig. 1. Jazyk (cs) / jazyk (sk)

It can be seen that out of the 16 most frequent adjectival collocates on the source language (SL) side, 11 have their corresponding translation equivalents among the most frequent ones on the target language (TL) side (indicated by a line).

In compiling a bilingual dictionary, such parallel sketches can be effectively used for selection of appropriate examples. The author usually considers both the collocations appearing in both tables and also the most significant collocations that differ between the SL and the TL.

7. The Russian-Slovak Parallel Sketches

As the MULTTEXT-East tagset used in PoS-tagging of the Russian corpus is quite similar to the Slovak tagset⁹, porting the

⁸ <https://trac.sketchengine.co.uk/attachment/wiki/SkE/DocsIndex/ske-stat.pdf?format=raw>

sketch grammar was a fairly straightforward venture, partially repeating the work done by Maria Khokhlova back in 2010¹⁰.

Now, let us have a look at an example of the degree of parallelism that can be found between the Russian and Slovak word sketches. Within the limited size of this paper, we can only afford to present a small probe, covering just the already mentioned keyword *язык* «language». The Fig. 2 shows the situation of adjectival collocates.

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
иностраннный	<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

Fig. 2. Язык (ru) / jazyk (sk)

Though the number of parallelisms is smaller than that of Czech vs. Slovak, most SL collocates (9 out of 16) have TL translation equivalents. A similar situation can be observed with nouns in Fig. 3 (9 out of 16) and even better with verbs in Fig. 4 (12 out of 16). The final figure, Fig.5, shows some parallelisms seen within coordination collocates (8 out of 16).

⁹ <http://korpus.juls.savba.sk/morpho.html>

¹⁰ Khokhlova M. Building Russian Word Sketches as Models of Phrases. Proc. EURALEX, Leeuwarden, July 2010.

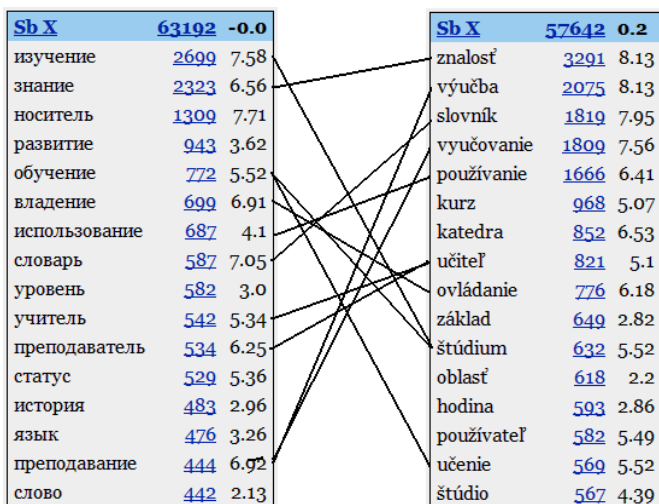


Fig. 3. Язык (ru) / jazyk (sk)

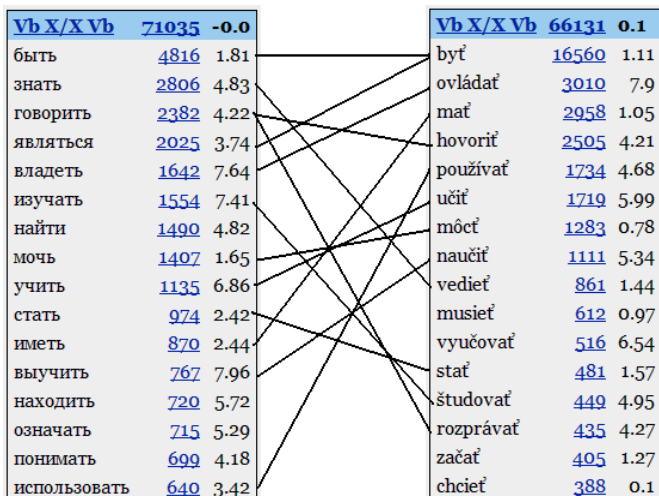


Fig. 4. Язык (ru) / jazyk (sk)

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

Fig. 5. Язык (ru) / jazyk (sk)

8. Conclusion and Further Work

Our experiment was targeted both at testing the new tool for creating web corpora, SpiderLing, and at finding out whether the degree of parallelism between Czech and Slovak word sketches can also be observed between Russian and Slovak, if a compatible sketch grammar is used. As all the tools necessary were either freely available or could be easily modified from our own tools, the whole experiment took approximately only 2 weeks to complete. SpiderLing proved to be extremely effective in gathering Russian web data. Initial analyses of the sketch data show that the degree of parallelism between Russian and Slovak is similar to that between Czech and Slovak and their use for a bilingual lexicographic project is feasible. Our experiment will continue in two directions. We want to increase the size of the Russian corpus to make it comparable with the largest Slovak corpus available. We also want to include other European languages into our corpus set, at least those covered by TreeTagger.