

Recent Czech Web Corpora

Vít Suchomel

Natural Language Processing Centre
Faculty of Informatics, Masaryk University
Botanická 68a, 602 00 Brno, Czech Republic
`xsuchom2@fi.muni.cz`

Abstract. This article introduces the largest Czech text corpus for language research – *czTenTen12* with 5.4 billion tokens. A brief comparison with other recent Czech corpora follows.

Key words: web corpora, Czech

1 Introduction

Algorithms in the field of natural language processing generally benefit from large language models. Many words and phrases occur rarely, therefore there is a need for very large text collections to research behaviour of words. [1] Furthermore, the quality of the data obtained from the web is also stressed. [2] Language scientists are increasingly turning to the web as a source of language data. [3] Nowadays, the web is the biggest, easily exploitable and the cheapest source of text data.

We decided to support corpora based research of Czech language by building a new Czech corpus from web documents. The aim was to apply successful data cleaning tools and label the words with grammatical categories.

2 Building a new Czech web corpus

CzTenTen12 is a new Czech web corpus built in 2012 using data obtained from the web in 2011. Several automatic cleaning and postprocessing techniques were applied to the raw data to achieve a good quality corpus for language research.

2.1 Crawling the web

We used web crawler *SpiderLing*¹, our previous work [4], to gather the data from the web. We started the crawl from 20000 seed URLs spanning over 8600 domains. The URLs were chosen using Corpus Factory [5], Czech Wikipedia and partially from older web corpus *czTenTen*. The crawl was restricted to the Czech national top level domain (.cz). 15 million documents of size 500 GB were downloaded in 24 days.

¹ <http://nlp.fi.muni.cz/trac/spiderling>

2.2 Postprocessing and tagging

The crawler performed character encoding detection and converted the data to UTF-8. The crawler detected language using character trigrams and filtered out texts in other languages than the focus language. Extra care had to be taken in case of Slovak which contains similar trigrams and unwanted texts may pass the filter. We prepared a list of Czech words not present in Slovak and a dual list of Slovak words not present in Czech. Using these lists, paragraphs containing three times more unique Slovak words than unique Czech words (0,4 %) or unique words mixed from both languages (6,4 %) were separated.

Boilerplate removal tool `justext`² was used to remove html markup, page navigation, very short paragraphs and other useless web content. The data was de-duplicated by removing exact and near duplicate paragraphs using tool `onion`³. Paragraphs containing more than 50 % seen 7-grams were dropped.

Paragraphs containing only words without diacritical marks were tagged for further use, e.g. when studying the informal language of the web.⁴ Parts containing more than 20 % of words not recognized by morphological analyzer *Desamb* were considered nonsense and removed from the corpus. The final size of the corpus reaches 5.4 billion tokens (4.4 billion words).

Czech morphological analyzer *Desamb* [6,7] was used to tag the corpus. The added information consists in the part of speech and other grammatical categories (where applicable): gender, number, case, aspect, modality and other.⁵

3 Comparison with other corpora

The following recent Czech corpora are used in the comparison:

- *SYN 2010* ... Czech national corpus – the SYN-series corpora up to 2010^{6,7},
- *czes2* (a web corpus from 2009),
- *czTenTen* (a web corpus from 2011),
- the Hector project corpus⁸ (*Hector*) [2].

3.1 Basic properties

According to [2], both *SYN* and *Hector* are deliberately balanced, e.g. the latter consists of 450 millions of words from news and magazines, 1 billion of words

² <http://nlp.fi.muni.cz/projects/justext>

³ <http://nlp.fi.muni.cz/projects/onion>

⁴ These texts come mostly from discussions or other informal websites (where people do not bother writing proper accent marks).

⁵ Reference of full tagset: <http://nlp.fi.muni.cz/projekty/ajka/tags.pdf>

⁶ <http://ucnk.ff.cuni.cz/english>

⁷ Although the full corpus is not publicly available, a wordlist with frequencies was enough to carry out measurements presented later on.

⁸ <http://hector.ms.mff.cuni.cz>

Table 1. Basic comparison of corpora. Only words consisting of letters are accounted in the word count. *Dictionary size* is the number of unique words with at least 5 occurrences. The *the-score* is the rank of word "the" in a list of words sorted by frequency from the most frequent one. The lower the value, the higher contamination by foreign words should be expected.

corpus	word count [10^6]	dictionary size [10^6]	the-score
SYN2010	1300	1.61	7896
czes2	367	1.03	42
czTenTen	1652	2.42	1023
Hector	2650	2.81	1184
czTenTen12	4439	4.16	1223

Table 2. Corpus distance measured for each couple of corpora. The lower the *distance score*, the more similar is the couple.

corpus	czes	czTenTen	Hector	czTenTen12
SYN2010	1.60	1.70	2.28	1.73
czes2		1.44	2.16	1.52
czTenTen			1.79	1.12
Hector				1.65

from blogs and 1.1 billion of words from discussions. The content of the new corpus was not controlled and a deeper analysis of content remains for further research.

Table 1 displays values of three metrics calculated for five corpora. We observe *czTenTen12* is the largest corpus with the largest dictionary. The *the-score* is a very simple metric offering a basic idea about contamination of the corpus by foreign (English) words. We observe *czes2* is the most polluted corpus and *SYN2010* is the most clean corpus in this measurement.

3.2 Corpora similarity

Table 2 shows a corpus comparison cross-table. The *distance score* calculation is based on relative corpus frequencies of 500 most frequent words in all corpora. The full method is described in [8]. We observe *czTenTen* and *czTenTen12* are very close. That can be explained by similar way of obtaining and processing the data and sharing a lot of documents. On the other hand, the balanced corpora are more distant.

Comparison of keywords (also based on the relative corpus frequency) in *czTenTen12* most different from *SYN2010* was published in [9]. We observe there are more discussions and blogs (informal words, verbs in 1st or 2nd person, pronouns, adverbs) and computer related words in the new unbalanced corpus. Comparing *czTenTen12* to *Hector*, we find the difference in presence of informal words too. Top *czTenTen12* related words in this comparison are quite formal: *již, lze, oblasti, společnosti, zařízení, této, roce, zde, mohou, rámci, projektu, těchto,*

<u>has_obj7</u>	<u>905</u>	<u>80.6</u>	<u>post dnem</u>	<u>564</u>	<u>23155.1</u>	<u>has subj</u>	<u>507</u>	<u>4.1</u>
dveře	<u>384</u>	8.09	nabytí	<u>272</u>	10.4	katolizace	<u>4</u>	7.88
žeň	<u>8</u>	6.91	účinnost	<u>99</u>	7.13	ožen	<u>3</u>	7.53
blaho	<u>4</u>	5.26	konání	<u>31</u>	6.33	bázeň	<u>4</u>	6.59
léto	<u>170</u>	4.92	splatnost	<u>5</u>	4.53	kočka	<u>18</u>	6.0
den	<u>154</u>	3.78	volba	<u>74</u>	3.76	nit	<u>5</u>	5.93
let	<u>35</u>	3.75	projednání	<u>4</u>	3.72	borec	<u>7</u>	5.63
mše	<u>3</u>	3.53	hlasování	<u>10</u>	3.61	mlha	<u>8</u>	5.53
týden	<u>38</u>	3.15	podání	<u>10</u>	3.38	dávno	<u>4</u>	4.42
měsíc	<u>28</u>	3.13	nástup	<u>5</u>	2.99	zázrak	<u>6</u>	4.37
dno	<u>7</u>	3.08	vznik	<u>6</u>	2.04	puk	<u>4</u>	4.02
rok	<u>51</u>	1.34	zahájení	<u>3</u>	1.66	kluk	<u>10</u>	3.71
			jednání	<u>7</u>	0.71	kousek	<u>8</u>	3.41
<u>coord</u>	<u>47</u>	<u>0.5</u>	<u>post na</u>	<u>9</u>	<u>0.2</u>	pán	<u>8</u>	3.31
tkát	<u>4</u>	8.74	kolovrátek	<u>4</u>	9.27	holka	<u>3</u>	3.0
						paní	<u>6</u>	2.7

Fig. 1. Word sketch for word *příst* in *czes2*. A part of grammatical relations is displayed. The number of hits of the word in the corpus is 5191.

systemu. They could belong to some project notes or contracts. The key words of the opposite direction are *no, holky, jo, xD, D, blog, teda, taky, já, dneska, sem, jdu, máš*, which leads to conclusion *Hector* contains more blogs and discussions (generally informal texts) than *czTenTen12*.

3.3 Word sketches – bigger is better

Figures 1 and 2 display word sketch for word *příst* in *czes2* and *czTenTen12* in SketchEngine⁹. As can be easily observed, the bigger corpus offers better words in relations with the head word. E.g. (*příst, blaho*) in relation *has_obj7* (which stands for a verb with a noun in instrumental case) is a quite common collocation in Czech. That is well reflected in the sketch for *czTenTen12* with 84 occurrences and the first place by saliency score in the relation table. However, the smaller corpus offers only 4 instances of this collocation. Relation *has_obj4* (which stands for a verb with a noun in accusative case) in *czes2* is very poor, while containing many well suiting words in the case of the bigger corpus: *len*,

⁹ <http://sketchengine.co.uk/>

<u>has subj</u>	<u>3653</u>	<u>-6.2</u>	<u>has obj7</u>	<u>2028</u>	<u>-27.1</u>	<u>coord</u>	<u>1056</u>	<u>-1.7</u>
mha	<u>50</u>	8.52	blaho	<u>84</u>	6.08	tkát	<u>63</u>	8.55
mlha	<u>129</u>	5.82	žeň	<u>28</u>	6.0	příst	<u>121</u>	7.44
kolovrátek	<u>10</u>	5.73	dveře	<u>700</u>	4.82	vrnět	<u>28</u>	7.24
přadlena	<u>5</u>	5.29	mše	<u>110</u>	4.8	lísat	<u>7</u>	6.52
kočka	<u>268</u>	5.18	slast	<u>6</u>	3.39	mazlit	<u>35</u>	5.99
rohožka	<u>7</u>	5.02	rozkoš	<u>6</u>	2.71	mňoukat	<u>7</u>	5.98
kocour	<u>34</u>	4.42	dno	<u>42</u>	2.57	tulit	<u>12</u>	5.81
len	<u>12</u>	4.41	spokojenost	<u>17</u>	2.0	otírat	<u>12</u>	4.96
hospodyně	<u>6</u>	3.99	ústrojí	<u>5</u>	1.9	přešlapovat	<u>5</u>	3.99
kotě	<u>19</u>	3.83	den	<u>546</u>	1.81	šít	<u>10</u>	3.53
nit	<u>17</u>	3.7	léto	<u>108</u>	0.85	hřát	<u>10</u>	3.52
pisatel	<u>8</u>	3.49	let	<u>20</u>	0.39	hladit	<u>15</u>	3.26
pavouk	<u>11</u>	3.23	svědek	<u>5</u>	0.27	nastavovat	<u>6</u>	1.21
chlápek	<u>8</u>	3.21				plést	<u>6</u>	1.06
motorek	<u>5</u>	3.14				prát	<u>5</u>	0.8

<u>has obj4</u>	<u>1041</u>	<u>-2.3</u>	<u>post na</u>	<u>285</u>	<u>-1.2</u>	<u>post pod</u>	<u>27</u>	<u>-2.7</u>
len	<u>43</u>	6.47	vřetánek	<u>5</u>	8.99	kapota	<u>5</u>	2.51
motůrek	<u>5</u>	6.16	kolovrátek	<u>42</u>	8.94			
příze	<u>18</u>	5.58	kolovrat	<u>27</u>	8.2	<u>post o</u>	<u>20</u>	<u>-0.3</u>
nit	<u>49</u>	5.32	volnoběh	<u>5</u>	5.07	stošest	<u>6</u>	7.64
kapsička	<u>18</u>	5.28	klín	<u>19</u>	3.62			
kotě	<u>40</u>	4.99	zip	<u>9</u>	2.18	<u>post do</u>	<u>66</u>	<u>-0.8</u>
nitka	<u>12</u>	4.69				zad	<u>14</u>	4.04
pavučina	<u>6</u>	3.6	<u>post v</u>	<u>73</u>	<u>-0.3</u>	ouško	<u>6</u>	3.37
zapínání	<u>8</u>	3.48	náručí	<u>5</u>	2.88	ucho	<u>13</u>	1.0
kotátko	<u>7</u>	3.07						

Fig. 2. Word sketch for word *příst* in *czTenTen12*. A part of grammatical relations is displayed. The number of hits of the word in the corpus is 28276, that is 5 times more frequent than in the smaller corpus.

příze, nit, nitka, pavučina.¹⁰ Furthermore, most of collocations with prepositions in prepositional relations *post_na*, *post_v*, *post_do* and other are present just in the word sketch of the bigger corpus: *na kolovratu, na klíně, v náručí, do ouška, pod kapotou, ostošest*. We conclude a bigger corpus is much more useful for language research based on collocations of words.

4 Conclusion and future work

This article introduced the largest Czech text corpus for language research. A basic comparison with other contemporary Czech corpora was made. Example work sketches were shown to support idea that bigger corpora are better.

The future plans for building web corpora of Slavonic languages include gathering resources in Polish and Croatian. Another interesting research opportunity is studying semantic topics automatically extracted from documents in the corpus. That would help us to know more about the content of the corpus and consequently of the Czech web.

Acknowledgements

This work has been partially supported by the Ministry of Education of CR within the LINDAT-Clarin project LM2010013, by the Ministry of the Interior of CR within the project VF20102014003 and by the Czech Science Foundation under the project P401/10/0792.

References

1. Pomikálek, J., Rychlý, P., Kilgarriff, A.: Scaling to billion-plus word corpora. *Advances in Computational Linguistics* **41** (2009) 3–13
2. Spoustová, J., Spousta, M.: A high-quality web corpus of czech. In: *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, Istanbul, Turkey, European Language Resources Association (ELRA) (may 2012)
3. Kilgarriff, A., Grefenstette, G.: Introduction to the special issue on the web as corpus. *Computational linguistics* **29**(3) (2003) 333–347
4. Suchomel, V., Pomikálek, J.: Efficient web crawling for large text corpora. In: *Proceedings of the Seventh Web as Corpus Workshop*, Lyon, France (2012)
5. Kilgarriff, A., Reddy, S., Pomikálek, J., Pvs, A.: A corpus factory for many languages. *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'10, Malta)* (2010)
6. Šmerk, P.: Unsupervised Learning of Rules for Morphological Disambiguation. In: *Lecture Notes in Artificial Intelligence* 3206, *Proceedings of Text, Speech and Dialogue 2004*, Berlin, Springer-Verlag (2004) 211–216
7. Jakubíček, M., Horák, A., Kovář, V.: Mining phrases from syntactic analysis. In: *Lecture Notes in Artificial Intelligence, Proceedings of Text, Speech and Dialogue 2009*, Plzeň, Czech Republic, Springer-Verlag (2009) 124–130

¹⁰ Presence of other words in this relation is caused by tagging mistakes or by putting them in a wrong relation.

8. Kilgarriff, A.: Comparing corpora. *International journal of corpus linguistics* 6(1) (2001) 97–133
9. Kilgarriff, A.: Getting to know your corpus. In: *Text, Speech and Dialogue*, Springer (2012) 3–15