

# The slWaC 2.0 Corpus of the Slovene Web

Tomaz Erjavec<sup>†</sup> and Nikola Ljubešić\*

<sup>†</sup>Dept. of Knowledge Technologies, Jožef Stefan Institute  
Jamova cesta 39 1000 Ljubljana  
tomaz.erjavec@ijs.si

\*Faculty of Humanities and Social Sciences, University of Zagreb  
Ivana Lučića 3, 10000 Zagreb, Croatia  
nikola.ljubestic@ffzg.hr

## Abstract

Web corpora have become an attractive source of linguistic content, as they can be made automatically, contain varied text types of contemporary language, and are quite large. This paper introduces version 2 of slWaC, a web corpus of Slovene containing 1.2 billion tokens. The corpus extends the first version of slWaC with new materials and updates the corpus compilation pipeline. The paper describes the process of corpus compilation with a focus on near-duplicate removal, presents the linguistic annotation, format and accessibility of the corpus via web concordancers, and then investigates the content of the corpus using frequency profiling, by comparing its lemma and part-of-speech annotations with the first version of slWaC and with KRES, the reference balanced corpus of Slovene.

## Korpus slovenskega spleta slWaC 2.0

Korpusi besedil zajetih s spleta so postali popularen vir jezikovnih vsebin, saj jih lahko zgradimo avtomatsko, vsebujejo pester nabor sodobnih besedilnih zvrsti in so zelo veliki. Prispevek predstavi drugo različico korpusa slWaC, spletnega korpusa slovenščine, ki vsebuje 1,2 milijarde pojavnic. Korpus dopolnjuje prvo različico slWaC z novimi besedili, pridobljenimi z izboljšanimi orodji za zajem. V prispevku opišemo proces izdelave korpusa s poudarkom na odstranjevanju podobnih vsebin, predstavimo jezikoslovno označevanje, format korpusa in njegovo dostopnost preko konkordančnika. Nato raziščemo vsebino korpusa s pomočjo frekvenčnega profila, kjer leme in oblikoskladenjske oznake druge različice korpusa slWaC primerjamo s prvo ter z referenčnim in uravnoteženim korpusom slovenščine KRES.

## 1. Introduction

With the advent of the web, a vast new source of linguistic information has emerged. The exploitation of this resource has especially gained momentum with the WaCky initiative (Baroni et al., 2009), which has popularised the concept of "Web as Corpus". It has also made available tools for compiling such corpora and produced large WaC corpora for a number of major European languages. Now such corpora are also being built for the so called smaller languages, such as Norwegian (Guevara, 2010), Czech (Spoustová et al., 2010) and Serbian (Ljubešić, 2014), moving the concept of a "large corpus" for smaller languages up to the 1 billion token frontier. As Web corpus acquisition is much less controlled than that for traditional corpora, the necessity of analysing their content gains in significance. The linguistic quality of the content is mostly explored through word lists and collocates (Baroni et al., 2009) while the content itself is explored using unsupervised methods, such as clustering and topic modelling (Sharoff, 2010).

For Slovene, a web corpus has already been built (Ljubešić and Erjavec, 2011). However, the first version of slWaC (hereafter slWaC<sub>1</sub>) was rather small, as it contained only 380 million words. Furthermore, it contained domains from the Slovene top-level domain (TLD) only, i.e. only URLs ending with ".si" were harvested. In the meantime, hrWaC, the Croatian web corpus had already moved to version 2, touching the 2 billion token mark, and web corpora for Serbian and Bosnian were built as well (Ljubešić, 2014), all of them passing the size of slWaC<sub>1</sub>, making it high time to move forward also with slWaC.

This paper presents version 2 of slWaC (hereafter

slWaC<sub>2</sub>) which tries to overcome the limitations of slWaC<sub>1</sub>: it extends it with a new crawl, which also includes well known Slovene web domains from other TLDs, and introduces a new pipeline for corpus collection and cleaning, resulting in a corpus of 1.2 billion tokens with removed near-duplicate documents and flagged near-duplicate paragraphs.

The rest of the paper is structured as follows: Section 2 presents the corpus construction pipeline, Section 3 introduces the linguistic annotation of the corpus, its format and its availability for on-line concordancing, Section 4 investigates the content of the corpus, by comparing it to slWaC<sub>1</sub> and to the KRES balanced corpus of Slovene, while Section 5 gives some conclusions and directions for future work.

## 2. Corpus construction

### 2.1. Crawling

For performing the new crawl we used the SpiderLing crawler<sup>1</sup> with its associated tools for guessing the character encoding of a web page, its content extraction (boilerplate removal), language identification and near-duplicate removal (Suchomel and Pomikálek, 2012). The SpiderLing crawler has two predefined size ratio thresholds that control when a low-yield-rate web domain (concerning new text) is to be abandoned; we used the lower one which is recommended for smaller languages. As seed URLs we used the home pages of web domains obtained during the construction of slWaC<sub>1</sub> and additionally 30 well known Slovene web domains, which are outside the .si TLD.

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

The crawl was run for 21 days, with 8 cores used for document processing, which includes guessing the text encoding, text extraction, language identification and physical duplicate removal, i.e. removing copies of identical pages which appear under different URLs. After the first 14 days there was a significant decrease in computational load, showing that most of the domains had been already harvested and that the process of exhaustively collecting textual data from the extended Slovene TLD was almost finished.

After completing the crawling process, which already included document preprocessing, we merged the new crawl with slWaC<sub>1</sub>. We added the old dataset to the end of the new one, thereby giving priority to new data in the following process of near-duplicate removal. It should be noted that the corpus can, in cases when the content has changed, contain two texts with the same URL but with different crawl dates.

## 2.2. Near duplicate removal

We performed near-duplicate identification both on the document and the paragraph level using the onion tool<sup>2</sup> with its default settings, i.e. by calculating 5-gram overlap and using the 0.5 duplicate content threshold. We removed the document-level near-duplicates entirely from the corpus, while keeping paragraph-level near-duplicates, labelling them with a binary attribute on the <p> element. This means that the corpus still contains the (near)duplicate paragraphs, which is advantageous for showing contiguous text from web pages, but if, say, language modelling for statistical machine translation were to be performed (Ljubešić and Toral, 2014), near-duplicate paragraphs can easily be removed.

The resulting size of the corpus (in millions of tokens) after each of the three duplicate removal stages is given in Table 1. We compare those numbers to the ones obtained on the Croatian, Bosnian and Serbian domains (Ljubešić, 2014), showing that the second versions of the corpora (hrWaC and slWaC), which merge two crawls obtained with different tools and were collected three years apart, show a smaller level of reduction (around 30%) at each step of near-duplicate removal, while the first versions of corpora (bsWaC and srWaC), obtained with SpiderLing only and in one crawl, suffer more data loss in this process (around 35-40%).

	PHYS	DOCN	PARN	R1	R2
<b>slWaC 2</b>	1,806	<b>1,258</b>	<b>895</b>	0.31	0.29
hrWaC 2	2,686	1,910	1,340	0.29	0.30
bsWaC 1	722	429	288	0.41	0.33
srWaC 1	1,554	894	557	0.42	0.37

Table 1: Sizes of the web corpora in millions of tokens after removing physical duplicates (PHY), document near-duplicates (DOCN) and paragraph near-duplicates (PARN), with the reduction ratio (R1 and R2) after the DOCN and subsequent PARN steps.

## 2.3. Linguistic annotation

slWaC<sub>2</sub> was tagged and lemmatised with ToTaLe (Erjavec et al., 2005) trained on JOS corpus data (Erjavec and Krek, 2008). However, it should be noted that ToTaLe had been slightly updated, so in particular the tokenisation of slWaC<sub>1</sub> and slWaC<sub>2</sub> at times differs. The morphosyntactic descriptions (MSDs) that the words of the corpus are annotated with follow the JOS MSD specifications, however, these do not define a tag for punctuation. As practical experience has shown this to be a problem, we have introduced a punctuation category and MSD, named “Z” in English and “U” in Slovene.

## 3. Overview of the corpus

### 3.1. Size of the corpus

Table 2 gives the size of slWaC<sub>2</sub>, for the included slWaC<sub>1</sub> from 2011 and the new additions in 2014, and together. For each of the counted elements we also give the size of the complete corpus, i.e. after removing document near duplicates (DOCN from Table 1), and for the corpus which has also paragraph near duplicates removed (PARN).

slWaC <sub>2</sub>	2011	2014	All
Domains	25,536	22,062	37,759
URLs	1,528,352	1,295,349	2,795,386
Pars	7,535,453	18,303,123	25,838,576
(PARN)	6,325,075	10,329,692	16,654,767
Sents	22,615,610	50,693,747	73,309,357
(PARN)	19,001,653	31,560,289	50,561,942
Words	360,273,022	718,332,186	1,078,605,208
(PARN)	301,547,669	465,780,456	767,328,125
Tokens	421,178,853	837,727,874	1,258,906,727
(PARN)	352,474,874	542,912,192	895,387,066

Table 2: Size of the slWaC 2.0 corpus.

Starting with the number of domains, it can be seen that the new crawl produced less domains than the first one, due to a large number (of the complete space of URLs) of static domains being removed in the physical deduplication stage (PHY). Nevertheless, the complete corpus has, in comparison to slWaC<sub>1</sub>, about 12,000 new domains. Observing the URLs, we note that the new crawl gave somewhat less URLs than the old one, and that there is little overlap between the two, i.e. about 1%: 28,315 URLs are the same from both crawls, which means that their content has changed in the last three years (and are then in the corpus distinguished by having a different crawl date).

Starting the the number of paragraphs we give both the numbers for DOCN and PARN, with the reduction having been already expressed in Table 1, i.e. 29%. For paragraphs, sentences, words and tokens, the complete corpus is simply the sum of the items for each of the two crawls. The most important numbers are the sizes of the complete corpus in tokens, i.e. 1.25 billion words for the DOCN and 900 million for PARN, which makes the corpus almost as large as the largest corpus of Slovene to date, i.e. Gigafida.

<sup>2</sup><https://code.google.com/p/onion/>

### 3.2. Corpus format

The annotated corpus is stored in the so called vertical format, used by many concordancing engines. This is an XML-like format in that it has opening and closing or empty (structural) XML tags, but the tokens themselves are written one per line, with the first (tab separated) column giving the token (word or punctuation) itself, the second (in our case) its lemma (or, for punctuation, again the token), the third its MSD in English and the fourth the MSD in Slovene, as illustrated by Figure 1.

```
<text domain="www.cupradan.si"
  url="http://www.cupradan.si/"
  crawled="2014">
<gap extent="1000+"/>
<p type="text" duplicate="0">
<s>
*      *      Z      U
Izmed  izmed  Sg      Dr
vseh   ves     Pg-mpg  Zc-mmr
</g/>
,      ,      Z      U
ki     ki     Cs      Vd
boste  biti   Va-f2p-n Gp-pdm-n
delili deliti Vmpp-pm  Ggnd-mm
video  video  Ncmsan   Sometn
...

```

Figure 1: Vertical format of the annotated sIWaC<sub>2</sub>.

The example also shows a few other features of the encoding. Each text is given its URL, the domain of this URL and the year (2011 or 2014) on which it was crawled. Boilerplate removal often deletes linguistically uninteresting texts from the start (and end) of the document, which is marked by the empty gap element, which also gives the approximate extent of the text removed. The paragraphs are marked by their type, which can be “heading” or “text”, while the “duplicate” attribute tells whether the paragraph is a (near) duplicate of some other paragraph in the corpus, in which case its value is “1”, and “0” otherwise. Finally, we also have the empty “glue” element g, which can be used to remove the space between two adjacent tokens in displaying the corpus.

### 3.3. Availability

The corpus is mounted under the noSketchEngine concordancer (Rychlý, 2007) installed at nl.ijs.si/noske. The concordancer allows for complex searches in the corpus, from concordances taking into account various filters, to frequency lexica over regular expressions.

We also make the corpus available for download, but not directly, mainly due to question of personal data protection. Namely, the corpus contains most of the Slovene Web, at least in the .si domain, so it also contains a lot of personal names with accompanying text. This is not such a problem with the concordancer, as similar results on Web-accessible personal names can be also obtained by searching through Google or Najdi.si. However, being able to analyse the complete downloaded corpus enables much

more powerful information extraction methods to be used, potentially leading to abuse of personal data. This is why we make the corpus available for research only, and require a short explanation of the use it will be put to. However, we (will) make available the metadata of the corpus, in particular the list of URLs included in it, which enables other to make their own corpus on this basis.

## 4. Comparative corpus analysis

This section investigates how different the sIWaC<sub>2</sub> corpus is from its predecessor, sIWaC<sub>1</sub> and from the KRES balanced reference corpus of Slovene (Logar et al., 2012). For this we used the method of frequency profiling, introduced by (Rayson and Garside, 2000). We first made a frequency lexicon of the annotation under investigation (lemma or grammatical description) for sIWaC<sub>2</sub> and the corpus it was compared with, and then for each item in this lexicon computed its log-likelihood (LL). The formula takes into account the two frequencies of the element as well as the sizes of the two corpora which are being compared; the greater LL is, the more the item is specific for one of the corpora. To illustrate, we give in Table 3 the first 15 lemmas with their LL score and their frequency per million words in sIWaC<sub>1</sub> and sIWaC<sub>2</sub>, with the larger frequency in bold.

Lemma	LL	sIWaC <sub>1</sub> pm	sIWaC <sub>2</sub> pm
člen	30,366	0.131	<b>0.282</b>
foto	23,092	0.018	<b>0.081</b>
m2	22,826	0	<b>0.033</b>
biti	22,767	<b>76,984</b>	74,493
°	21,447	0.001	<b>0.036</b>
3d	17,738	0	<b>0.026</b>
spoštovan	11,177	0.019	<b>0.059</b>
2x	11,092	0	<b>0.016</b>
tožnik	9,909	0.008	<b>0.036</b>
odstotek	9,265	<b>0.515</b>	0.393
co2	9,090	0	<b>0.013</b>
amandma	8,992	0.007	<b>0.031</b>
hvala	8,954	0.106	<b>0.173</b>
1x	8,505	0	<b>0.012</b>
ekspr	8,373	0	<b>0.012</b>

Table 3: The first 15 lemmas with highest log-likelihood scores and their frequency per million words for the comparison of the old and new version of sIWaC

As can be noted, most of these highest LL lemmas are more prominent in sIWaC<sub>2</sub>; only “biti” (*to be*) and “odstotek” (*percent*) are more frequent in sIWaC<sub>1</sub>. Furthermore, quite a few lemmas have frequency 0 in sIWaC<sub>1</sub>. This is indicative of a difference in annotation between the two corpora: as mentioned, the tokenisation module of ToTaLe had been somewhat improved lately, which is evidenced in the fact that strings, such as “m2” and “3d” were wrongly split into two tokens in sIWaC<sub>1</sub> but are kept as one in sIWaC<sub>2</sub>. It is a characteristic of LL scores that they show such divergences, which should ideally be fixed, to arrive at uniform annotation of the resources.

#### 4.1. Lemma comparison with slWaC

The motivation behind comparing the previous and current version of slWaC was primarily to investigate what kind of text types are better represented in the new (or old) version of the corpus. Apart from the already mentioned differences in tokenisation, slWaC<sub>2</sub> is more prominent in three types of lemmas (texts). First, there are legal texts, (characterised by lemmas such as “člen” (*article*), “odstavek” (*paragraph*), “amandma” (*amendment*) “tožnik” (*plaintiff*)), which come predominantly from governmental domains, e.g. for “člen” mostly from uradni-list.si (official gazette), dz-rs.si (parliament), sodisce.si (courts). Second are texts that address the reader (or, say, parliamentary speaker) directly (“spoštovan” (*honoured*), “pozdravljen” (*hello*)). For “spoštovan”, the most highly ranked domains are, again, the parliament, i.e. dz-rs.si, followed by vizita.si (medical help page of commercial POP.TV), delo.si (main Slovene daily newspaper), in the latter two mostly from user forums. The corpus is thus more representative in text-rich domains whose content changes rapidly and that contain user-generated content. Third, the list contains two interesting “lemmas” with very high LL scores. The first is “ekspr” (only 19 in slWaC<sub>1</sub> but more than 9,000 in slWaC<sub>2</sub>), which is the (badly tokenised) abbreviation “ekspr.” meaning “expressive”. It turns out that practically the only domain that uses this abbreviation is bos.zrc-sazu.si, i.e. the portal serving the monolingual Slovene dictionary SSKJ, which was newly harvested in slWaC<sub>2</sub>. Similarly, the word “ino” (less than 500 in slWaC<sub>1</sub> but more than 7,000 in slWaC<sub>2</sub>) turns out to be the historical form of “in” (*and*). Practically the only domain containing this word (6,000x) is nl.ijs.si, which now hosts a large library of old Slovene books. The new slWaC thus contains some extensive new types of texts coming from previously unharvested domains or domains that have had large amounts of new content added. Finally, it is worth mentioning that the first slWaC<sub>2</sub> proper noun appears only at position 36 in the LL list, and is “bratušek” with almost 6,000 occurrences, referring to Alenka Bratušek, the former PM of Slovenia.

It is also instructive to see which lemmas are now less specific against slWaC<sub>1</sub>. Interestingly, the greatest drop in frequency concerns the auxiliary verb “biti” (*to be*). As all texts contain this lemma, it is difficult to analyse where this difference comes from, but our hypothesis is that legal texts, of which there are now significantly more, are more likely to use the present tense and passive constructions, which are made without the auxiliary. Among function words, there are less particles “pa” (*but*), used more in informal texts and less of “da” (*that*), used to introduce relative clauses. One verb is much less used, “dejati” (*say, formal register*), indicating a drop in the proportion of news items, where reporting on what a certain person said is quite frequent. Most of the list of course consists of nouns: in slWaC<sub>2</sub> there is relatively less written about “odstotek” (*percent*), “delnica” (*share*), “milion” (*million*), “premier” (*prime minister*), “predsednik” (*president*), “dolar” (*dollar*), “zda” (*USA*), again indicating less news and also the shifting of major news topics. Also, “evro” (*Euro*) is used less, but then the Euro symbol

is used a lot more.

#### 4.2. Lemma comparison with KRES

With slWaC<sub>2</sub>, as with Web corpora in general, it is an interesting question of how representative and balanced they are. The easiest approach towards an answer is a comparison with “traditional” reference corpora, and such experiments have been already performed, e.g. between the British Web corpus ukWaC and BNC, the British National Corpus (Baroni et al., 2009). The comparisons have shown that while web corpora are different from classical corpora, which contain mostly printed sources, the differences are in general not great and so they can function as modern-day reference corpora.

We made a comparison between slWaC<sub>2</sub> and KRES (Logar et al., 2012), which is the balanced reference corpus of Slovene with 100 million words, sampled from Gigafida, the representative corpus of contemporary Slovene. Gigafida (*ibid*) contains texts from 1990 to 2012. The comparison shows that, as with slWaC<sub>1</sub>, some of the differences are due to the different linguistic analyses. As mentioned, slWaC<sub>2</sub> was processed with ToTaLe, while KRES used the Obeliks tokeniser, tagger and lemmatiser (Grčar et al., 2012), and the two disagree in some lemmatisations, the most prominent being “veliko/več” (*much*), “mogoče/mogoč” (*possible*), “edini/edin” (*only*), “desni/desen” (*right*), “levi/lev” (*left*), “volitve/volitev” (*elections*), as well as some differences in tokenisation, e.g. “le-ta” and “d.o.o.” as one token or three.

Real linguistic differences concern mostly two types of lemmas. The first are highly ranked non-content words such as “pa, tudi, ter, naš” (*but, also, and, our*), which most likely show the bias of slWaC to informal writing. The second are content lemmas, which fall into several groups: “spleten” (*Web*), “podjetje” (*company*), “tekma, ekipa” (*match, team*), “volitve” (*elections*), “sistem, uporabnik, aplikacija” (*system user, applications*), and “blog”, i.e. slWaC has more commercial, sports, political and computer related texts, and, of course, texts specific to the web (blogs).

Conversely, KRES shows more lemmas to do with legal texts, such as “člen, odstavek, zakon” (*article, paragraph, law*), so that even with slWaC<sub>2</sub> having more texts of this type than slWaC<sub>1</sub>, it still has much less than KRES. KRES also has many more of two highly specific lemmas: “tolar” (*former Slovene currency*) shows that KRES is by now already dated, while “wallander”, the hero of a series of detective novels, shows that KRES – at least in this instance – has too much text from a single source, in this case a book series.

#### 4.3. Grammatical comparison with KRES

Apart from lemmas, it is also interesting to compare how the distribution of morphosyntactic categories of slWaC<sub>2</sub> differs from that of KRES. To this end we calculated six LL comparison scores, for uni-, bi- and tri-grams of part-of-speech (PoS) and of complete morphosyntactic descriptions (MSDs).

The uni-gram PoS LL scores show that slWaC has significantly more adjectives, unknown words, conjunctions,

prepositions and particles, in this order. However, it has much less punctuation and numerals, and slightly less interjections. Esp. with unknown words and punctuation the differences might be, at least partially, an artefact of different annotation programs. For the others, the results show that s1WaC tends more towards informal, user generated language, although this conclusion is somewhat offset by the fact that it has less interjections. However, tagging interjections is notoriously imprecise, and the difference here might also be due to different taggers used. Conversely, KRES with its numerals shows a preponderance of newspaper texts, which tend to use lots of dates, times, amounts, and sports scores.

PoS bi-grams again highlight the different annotation tools used. The most prominent combination in s1WaC is a numeral followed by an abbreviation, e.g. “90 EUR, 206 kW, 298,80 m<sup>2</sup>” but this difference is due to the fact that in s1WaC “EUR”, “kW” etc. are treated as abbreviations, whereas they are common nouns in KRES. The same reasoning applies to combinations with punctuation. However, there are also legitimate combinations in the top scoring LL PoS bi-grams: s1WaC has more noun + verb, adjective + noun and verb + adjective combinations, while KRES has more numeral + numeral, numeral + noun and verb + verb combinations. Scores for PoS tri-grams give little new information: apart from annotation differences, the most prominent s1WaC combination is noun + noun + verb, which are mostly name + surname + predicate, e.g. “Oto Pestner naredil”, while the most prominent for KRES is a sequence of three numerals.

As for MSDs, the differences in unigrams in favour of s1WaC<sub>2</sub> are greatest for the three unknown word types that KRES doesn't use (Xf: foreign word, Xp: program mistake and Xt: typo), followed by general adverbs in the positive degree, coordinating conjunctions, present tense first person auxiliary verb in the plural (“smo”) and animate common masculine singular noun in the accusative, i.e. the object of a sentence, e.g. “otroka”. Conversely, KRES has much more punctuation, digits, common masculine and feminine singular nouns in the nominative (i.e. subjects) and general adverbs in comparative and superlative degrees. Bigrams show that s1WaC has many more general adjective + common noun combinations in various genders and cases, while KRES has many more combinations with digits. The space of MSD trigrams is very large, and, if we discount the combinations appearing as a result of different annotations, does not show very interesting differences.

## 5. Conclusion

The paper presented a new version of the Slovene Web corpus, which is almost three times larger than its initial version and is made available through a powerful and freely accessible concordancer. During the construction process we focused on the content reductions obtained through near-duplicate removal, showing that both reductions to document and paragraph level remove a similar amount of content. We also compared the content of the s1WaC<sub>2</sub> corpus to the s1WaC<sub>1</sub> corpus and to the reference corpus KRES via frequency profiling on lemmas and grammatical descriptions. This comparison showed that the new version

of the corpus has significantly more legal texts and specific text types, such as a dictionary and a library of historical books and (comparatively) less news. In the lemma comparison with KRES it has less legal texts but more user generated content and more commercial, sports, political and computer related texts, while the comparison of grammatical categories also shows a bias to informal writing as well as against newspaper items. But maybe the most surprising (although, in retrospect, quite logical) insight of the comparison using frequency profiling is that it is a very good tool to detect even slight differences in the processing pipelines used for the compared corpora, which then lead to significant differences in the (token, lemma and MSD) vocabularies.

There are several directions that our future work could take. First, by constructing the second version of two out of four existing web corpora of South Slavic languages, two ideas have emerged: one is to build a multilingual corpus consisting of all South Slavic languages, and the second to develop a monitor corpus which would be automatically extended with new crawls in predefined time frames. The second direction is in the annotation of the corpus, where more effort should be invested in developing a gold standard processing pipeline, which could then be used to re-annotate the Slovene corpora in a unified manner. In addition, given that the Web contains a significant portion of user generated content containing non-standard language, the annotation pipeline should be extended by introducing a standardisation (normalisation) step on word-forms, similar to our approach to modernisation of historical Slovene words (Scherrer and Erjavec, 2013), which would then give better lemmas and MSDs, allowing for easier exploration of Web corpora.

## 6. References

- Marco Baroni, Silvia Bernardini, Adriano Ferraresi, and Eros Zanchetta. 2009. The WaCky wide web: a collection of very large linguistically processed web-crawled corpora. *Language Resources and Evaluation*, 43(3):209–226.
- Tomaž Erjavec and Simon Krek. 2008. The JOS Morphosyntactically Tagged Corpus of Slovene. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco.
- Tomaž Erjavec, Camelia Ignat, Bruno Pouliquen, and Ralf Steinberger. 2005. Massive multilingual corpus compilation: Acquis Communautaire and ToTaLe. *Archives of Control Sciences*, 15(3):253–264.
- Miha Grčar, Simon Krek, and Kaja Dobrovoljc. 2012. Obeliks: statistični oblikoskladenjski označevalnik in lematizator za slovenski jezik. In *Zbornik Osme konference Jezikovne tehnologije*, Ljubljana. Jožef Stefan Institute.
- Emiliano Guevara. 2010. NoWaC: A Large Web-based Corpus for Norwegian. In *Proceedings of the NAACL HLT 2010 Sixth Web As Corpus Workshop*, WAC-6 '10, pages 1–7.
- Nikola Ljubešić and Antonio Toral. 2014. caWaC - a Web Corpus of Catalan and its Application to Language Mod-

- eling and Machine Translation. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, Reykjavik, Iceland, may. European Language Resources Association (ELRA).
- Nikola Ljubešić and Tomaž Erjavec. 2011. hrWaC and slWac: Compiling Web Corpora for Croatian and Slovene. In Ivan Habernal and Václav Matousek, editors, *Text, Speech and Dialogue - 14th International Conference, TSD 2011, Pilsen, Czech Republic, September 1-5, 2011. Proceedings*, Lecture Notes in Computer Science, pages 395–402. Springer.
- Nikola Ljubešić. 2014. {bs,hr,sr}WaC: Web corpora of Bosnian, Croatian and Serbian. In *Proceedings of the WAC-9 Workshop*.
- Nataša Logar, Miha Grčar, Marko Brakus, Tomaž Erjavec, Špela Arhar Holdt, and Simon Krek. 2012. *Korpusi slovenskega jezika Gigafida, KRES, ccGigafida in ccKRES: gradnja, vsebina, uporaba*. Zbirka Sporazumevanje. Trojina, zavod za uporabno slovenistiko: Fakulteta za družbene vede, Ljubljana.
- Paul Rayson and Roger Garside. 2000. Comparing Corpora Using Frequency Profiling. In *Proceedings of the Workshop on Comparing Corpora*, pages 1–6. Association for Computational Linguistics.
- Pavel Rychlý. 2007. Manatee/bonito – a modular corpus manager. *1st Workshop on Recent Advances in Slavonic Natural Language Processing*, pages 65–70.
- Yves Scherrer and Tomaž Erjavec. 2013. Modernizing historical Slovene words with character-based SMT. In *BSNLP 2013 - 4th Biennial Workshop on Balto-Slavic Natural Language Processing*, Sofia, Bulgaria.
- Serge Sharoff. 2010. Analysing Similarities and Differences between Corpora. In *Proceedings of the Seventh Conference on Language Technologies*, pages 5–11, Ljubljana. Jožef Stefan Institute.
- Drahomíra Spoustová, Miroslav Spousta, and Pavel Pecina. 2010. Building a Web Corpus of Czech. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta. European Language Resources Association (ELRA).
- Vít Suchomel and Jan Pomikálek. 2012. Efficient Web Crawling for Large Text Corpora. In Serge Sharoff Adam Kilgarriff, editor, *Proceedings of the seventh Web as Corpus Workshop (WAC7)*, pages 39–43, Lyon.