
Compatible Sketch Grammars for Comparable Corpora

Vladimír Benko

Slovak Academy of Sciences, L. Štúr Institute of Linguistics

Comenius University in Bratislava, UNESCO Chair in Translation Studies

vladob@juls.savba.sk

Abstract

Our paper describes an on-going experiment aimed at creating a family of billion-token web corpora that could to a large extent deserve the designation “comparable”: corpora are of the same size, data gathered by crawling the web at (approximately) the same time, containing similar web-specific domains, genres and registers, further pre-processed, filtered and deduplicated by the same tools, morphologically annotated by (possibly) the same tagger and made available via Sketch Engine. To overcome the problem of great differences in the existing sketch grammars for the respective languages, a set of “compatible” sketch grammars have been written that will aid contrastive linguistic research and bilingual lexicographic projects. The sketch grammars use a uniform set of rules for all word categories (parts of speech) and the resulting set of tables is displayed in a fixed order in all languages.

Keywords: comparable web corpora; sketch grammars; bilingual lexicography

1 Introduction

Ten years after its introduction to the lexicographic community at the Lorient Euralex Congress (Kilgarriff et al., 2004), Sketch Engine (*SkE*) has become a standard tool in numerous lexicographic projects, as well as in various areas of corpus-based linguistic research. Sketch grammars for corpora in many languages have been written (cf. References). Recently published open-source tools for efficient web crawling (Suchomel a Pomikálek, 2012) stimulate the building of very large web corpora, the analysis of which is hardly imaginable without a powerful summarisation machine such as *SkE*. Newly implemented *SkE* support for parallel and comparable corpora (Kovář, 2013) facilitate its use in the area of bilingual lexicography and contrastive linguistic research.

In bilingual and multilingual linguistic work with *SkE*, we often encounter the problem of sketch grammars defining the collocational profiles of a headword and its translation equivalent for the respective languages. Those sketch grammars have often been created for different purposes, having in mind different user requirements, with resulting word sketches being rather disparate, making its use for contrastive linguistics problematic. Our paper suggests an alternative approach to the creation of sketch grammars, within the framework of which the respective grammars can be made compatible for (almost) all languages.

2 The Aranea project

2.1 Why new corpora?

Besides our interest in testing the new corpus-building tools, the motive for starting a new corpus project was the lack of suitable corpora that could be used by students of foreign languages and translation studies at our university. The existing web corpora families that are available for download do not cover all the languages needed. As for corpora stored at the *SkE* web site¹, they (1) are not available for download, (2) are mostly too large for classroom use², and (3) have too different sketch grammars, which makes them difficult to use in a mixed-language classroom.

We expect that a set of corpora for several languages of equal size and built by a standardized methodology can not only be used for teaching purposes, but also in other areas of linguistic research (contrastive studies) and in lexicography (both mono- and bilingual).

2.2 The names

For our corpora, we have decided to use “language-neutral” Latin names denoting the language of the texts and the corpus size. The whole corpus family is called *Aranea*, and the respective members bear the appropriate language name, e.g. *Araneum Anglicum*, *Araneum Francogallicum*, *Araneum Russicum* for English, French, and Russian, respectively, etc.

Each corpus exists in several editions, differing by their sizes. The basic (medium-sized) version, *Maius* (“greater”), contains approximately 1.2 billion tokens (i.e., over 1 billion words). This size can be reached relatively quickly for all participating languages, and for the “large” ones with plenty of web data available, it usually takes just one or two days of download time. The 10% random sample of *Maius*, called *Minus* (“smaller”), is to be used for teaching purposes (e.g. for lessons in the framework of Corpus Linguistics programmes for students of foreign languages and translation studies). A 1% sample, *Minimum* (“minimal”), is not intended to be used directly by the end users, and is utilized in debugging the processing pipelines and tuning the sketch grammars. And lastly, the largest *Maximum* (“maximal”) edition will contain as much data as can be downloaded from the web for the particular language, and its size is mostly determined by the configuration of the server.

1 <http://www.sketchengine.co.uk/>

2 According to our experience, the ideal corpus for teaching corpus linguistics is about BNC-sized, i.e. it contains some 100 million tokens. As it is not easy to prevent students from invoking search operations taking several minutes to evaluate, billion-plus token corpora proved to be quite unsuitable for teaching purposes.

2.3 Web crawling

The source data acquisition is being performed by means of *SpiderLing*³, a web crawler optimized for collecting textual data from the web. The system contains an integrated character encoding (*chared.py*) and language recognition (*trigrams.py*) module, as well as a tool for boilerplate removal (*justext*). The input seed URLs (some 1,000 for each language) have initially been harvested by BootCAT⁴ (Baroni and Bernardini; 2004).

Several input parameters of the crawling process are to be set by the user, most notably the language name, a file containing sample text in the respective language (to produce a model for language recognition), a language similarity threshold (a value between 0 and 1; default 0.5), the number of parallel processes, and the crawling time.

In our processing, we usually crawled the web in 24-hour slots (the process could then be re-started) with all other values set to defaults. The only exception was crawling for Slovak and Czech, where we used 7-day slots, as the process was much slower here. The language similarity threshold also had to be changed in case of Slovak and Czech. As these languages are fairly similar, the trigram method did not seem to be able to distinguish between them sufficiently. We have therefore increased the similarity threshold value to 0.65 (saving many “good” documents, and causing many “wrong” ones to pass the filter) and removed the unwanted texts by subsequent filtration based on character frequencies .

2.4 Post-download processing

Besides the basic filtration aimed to remove texts with incorrect or misinterpreted character encoding, missing diacritics and texts with non-standard proportion of punctuation and uppercase characters, the main processing operation in this phase is tokenization. As the tokenization strategy has to be compatible with that of the corpus used to train the tagger, we decided to use the tokenizers supplied with Tree Tagger and TaKIPI for the respective languages. In the future, we would like to make use of the *unitok.py* universal tokenizing program developed at Masaryk University in Brno (Jakubíček; 2014).

2.5 Deduplication

The whole procedure (Benko; 2013) consists of three stages. The first stage detects near-duplicate documents by means of the Onion (Pomikálek; 2012) utility (similarity threshold 0.95), and the duplicate documents are deleted. The second stage deduplicates the remaining text at the paragraph level using the same procedure and settings. The tokens of the duplicate paragraphs, however, are not deleted but rather they are marked to make them “invisible” during corpus searches, while they can be displayed

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

4 <http://bootcat.sslmit.unibo.it/>

as context at the boundary of non-duplicate and duplicate text. In the last stage, we make use of our own tool based on the fingerprint method (with ignoring punctuation, special graphics characters and digits) to deduplicate the text at the sentence level. The tokens of duplicate sentences are marked similarly to the previous stage. This last step can “clean up” the duplicities among the short segments that fail to be detected as duplicates by Onion.

As deduplication is beyond the scope of our paper, we only mention here that the process has typically removed some 20–45% of tokens in the *Maius* versions of our corpora

2.6 Morpho-syntactic annotation

For languages covered by the parameter files of Tree Tagger (Schmid; 1994), this tagger has been used to annotate the respective corpora. For Polish, the TaKIPI (Piasecki; 2007), and for Czech, the Morče (Hajič; 2004) taggers were used, respectively. The question of tools for tagging Hungarian and Ukrainian data has not been resolved yet.

2.7 Tagging-related filtration

To improve the precision of tag assignments, a series of pre- and post-tagging filters have been developed that fix issues introduced by Unicode encoding of the source text⁵. The filtration fixes known tagger issues for the respective languages, namely the misassigned tags for many punctuation and special graphic characters (that are often tagged as nouns, adjectives, or abbreviations, and sometimes even as verbs with subcategories). For some languages, an additional tag with masked subcategories for gender and number is created, that is later used by some rules within the respective sketch grammars.

2.8 Current state of the project

At present, eight language versions of the *Aranea* corpus family have been created, containing both *Maius* and *Minus* editions as follows (in chronological order): *Araneum Russicum* (Russian), *Araneum Francogallicum* (French), *Araneum Germanicum* (German), *Araneum Hispanicum* (Spanish), *Araneum Polonicum* (Polish), *Araneum Anglicum* (English), *Araneum Nederlandicum* (Dutch), and *Araneum Slovaccum* (Slovak).

5 As an example we can point out the problem of the “apostrophe” character in French texts. As much as 8 different Unicode characters representing the apostrophe (with just two of them being “canonical”) can be found in the texts collected from the web. As the Tree Tagger French parameter file originated in the pre-Unicode era, even one of the two “canonical” representations would not be processed (i.e., tokenized and lemmatized) properly without special measures, and tokens like “l” and “d”, that belong to the most frequent ones, would be mistagged.

The crawling has also been done for *Araneum Bohemicum* (Czech). This data is now being pre-processed to be ready for annotation that will be performed by the Institute of Theoretical and Computational Linguistics at the Faculty of Arts of Charles University in Prague.⁶

The first stage of our project will be completed by *Araneum Hungaricum* (Hungarian), *Araneum Italicum* (Italian), and *Araneum Ukrainicum* (Ukrainian). With the exception of the last mentioned, we expect to complete the whole venture by the end of 2014.

For all of the languages mentioned, sketch grammars have been written and at least two rounds of testing have been performed for each corpus. The procedure involved is described in the following section.

3 Sketch grammars

A sketch grammar⁷ is a set of rules based on the CQL (Corpus Query Language⁸) used by the Sketch Engine to generate the respective collocation profiles (“word sketches”) for all lexical units (lemmas) in a corpus. The word sketches are pre-computed in advance, which makes the system user-friendly and very fast.

A sketch grammar rule consists of (1) an optional comment indicated by hash “#” character, (2) the rule type marked by an asterisk “*”, (3) the rule name preceded by the equal sign “=”, and (4) a list of CQL expressions. For example, a rule describing the relationship between two nouns (in English using the Penn Treebank tagset) might look as follows:

```
# noun followed by another noun
*DUAL
=modifier/modified
      2: [tag="NN.*"] 1: [tag="NN.*"]
```

The “1:” label denotes the “keyword”, i.e. the lemma the word sketch is created for, and the “2:” label marks the lemma of the collocate. The “*DUAL” keyword indicates that the rule is to be used twice, the second time with swapped labels, i.e. exchanging the positions of the keyword and the collocate. The text following the slash “/” character will be used as a name for the second use of the rule.

In reality, the rules usually look slightly more complex to indicate that “intermediate” words may be present between a keyword and a collocate, or in the vicinity of them.

6 Besides Ukrainian, Czech is the only language within the *Aranea* project with no free tagging tool available.

7 <https://www.sketchengine.co.uk/documentation/wiki/SkE/GrammarWriting>

8 <https://www.sketchengine.co.uk/documentation/wiki/SkE/CorpusQuerying>

3.1 What's in a name

Unlike Juliet Capulet⁹, we believe that the name is often really important, and the sketch grammar rule name is exactly such a case. On one hand, it is the only component of the sketch grammar that is not predetermined, and thus can be “virtually anything”. On the other hand, the name is the only clue for the user about the contents of the respective word sketch tables, and therefore should be as informative as possible. It has, however, to be very short as the name is displayed in the heading of the respective word sketch table within a only a limited space available. Rule names longer than 10–12 characters would increase the table widths, and the resulting word sketches could possibly not fit the screen.

Most sketch grammars used for corpora available at the *SkE* site follow the naming conventions introduced by A. Kilgarriff in the first English and French sketch grammars. These rule names are motivated syntactically, i.e. they denote the syntactic function of the collocate, with that of the keyword being implied. For example the rule name:

=modifier/modified

is representing two rule names with readings as follows: “collocate is a modifier of the keyword”, and “collocate is modified by the keyword”, respectively.

The syntactically motivated rules are transparent and user-friendly for description of basic relationships between subjects, object, modifiers/attributes, and verbs/predicates, but in more complex cases this strategy is not easily applicable. The nature of the problems can be observed in the Czech sketch grammar written by P. Smrž (Kilgarriff et al.; 2004). Some examples of rule names are as follows:

is_subj_of/has_subj
 is_obj4_of/has_obj4
 prec_prep
 gen1/gen2

As it can be seen, it is not really easy for the user the figure out “who is who” in the keyword – collocate – syntactic function “puzzle”. Moreover, rule names like “prec_verb” do not denote any syntactic functions but rather just describe collocational relationships.

There are two notable deviations from the “traditional” rule name conventions in the sketch grammars. In the grammar for the Slovenian FidaPLUS corpus¹⁰, S. Krek (Krek; 2006) uses rule names containing (among other features) Slovenian “case questions”. For example, the “*koga-česa*” name means

9 Juliet: “*What’s in a name? that which we call a rose / By any other name would smell as sweet*” (William Shakespeare: *Romeo and Juliet*, Act II, Scene 2).

10 <http://www.sketchengine.co.uk/documentation/wiki/Corpora/FidaPLUS>

that only collocates of the keyword that are in genitive case are displayed, with the syntactic function of the collocate being implied.

The second notable exception is the sketch grammar written by P. Whilelock (2010) for the Oxford English Corpus¹¹ (OEC) where the rule names not only name the syntactic function, but also the PoS of the keyword and the collocate and their mutual position within the collocation. For example, the “V* ADJ” rule name stands for verb modified by an adjective, with asterisk indicating the keyword.

3.2 Sketch grammar for Slovak corpora

In our Institute, the *SkE* 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 projects, as well as for other linguistic research activities.

The sketch grammar used in our *SkE* installation has been optimized for a lexicographic use, and differs from most “traditional” grammars for corpora stored at the *SkE* web site 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 number of rules and the order of resulting tables is fixed
- 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

3.3 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.

11 <http://www.sketchengine.co.uk/documentation/wiki/Corpora/OEC>

For nouns, 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), prepositions (Pp X; X Pp) and for immediate autosemantic collocates (Y X; X Y). 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 compatible morphological tags 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 and subcategories..

The compatible grammars

In creating sketch grammars for a group of languages, it is convenient not to use the “native” tagsets for the respective languages, but rather to use a common symbolic notation. This can be done, e.g., by means of a macro processor (such as m4). We have, however, decided to adopt a different approach and to create a simple universal tagset (“*Araneum Universal Tagset*” – AUT) similar to that of the *Universal PoS Tagset*¹² (UPT; Petrov et al., 2011), and to map all the respective tagsets into this tagset at the source vertical data level, i.e. to create a new layer of annotation. The AUT contains 11 tags for “traditional” part of speech categories, 7 additional tags for other elements, and one tag to indicate errors in the mapping process.

aTag	PoS	aTag	PoS	aTag	PoS
Dt	determiner/article	Pp	preposition	Xx	other (content word)
Nn	noun	Cj	conjunction	Xy	other other (function word)
Aj	adjective	Ij	interjection	Yy	unknown/foreign/alien
Pn	pronoun	Pt	particle	Zz	punctuation
Nm	numeral	Ab	abbreviation/acronym	Er	mapping error
Vb	verb	Sy	symbol		
Av	adverb	Nb	number		

Table 1: Araneum Universal Tagset (AUT).

The compatible sketch grammar using AUT consists of three sections. The first part (AUT-based) contains unary rules showing PoS category distribution for a particular lemma. The second part is

12 The AUT PoS tags for the eleven „traditional“ word classes directly correspond with those of UPT, with the difference being just in the names as we wanted to keep the names of the PoS categories identical with those used in the sketch grammar rule names introduced before the UPT tagset has been published. The additional 7 categories accommodate information provided by the respective “native” tagsets that is being ignored by UPT. For example, the “Xx” (other: content word) tag is assigned to participles in Slovak that have a category of their own in the SNK Slovak tagset.

tagset-dependent and contains unary rules showing PoS subcategories provided by the respective tagset. Due to differences in the depth of the morpho-syntactic annotation, the number of subcategories varies among the languages. With verbs, e.g., we have just 5 subcategories for Spanish, while more than 20 for Polish. The final third section (*AUT*-based) covers the collocational relationships of the respective keyword by means of binary, symmetric and trinary rules.

The compatible sketch grammar is basically identical for all the languages with one important exception: the number of intermediate tokens between a keyword and a collocate is increased by one for languages having articles in their language system.

4 Discussion and conclusion

A collocationally-based sketch grammar has (against a traditional one) several advantages. It can symmetrically cover all relationships between keywords and collocates of all word classes (parts of speech). As the PoS category is not tested for the keyword, a word sketch can be created even in cases of incorrectly assigned tags. If the same (compatible) sketch grammar is used with corpora for two or more languages, the resulting word sketches can be conveniently used in contrastive linguistic research, as well as within bilingual lexicographic projects.

The disadvantage of our approach is that not all tables for some words represent linguistically relevant relationships, and they may contain a lot of noise. We believe, however, that having a fixed number of tables gives the user a clear overview, and he or she can easily ignore the irrelevant data.

In the Appendix, we present the word sketches for the lemma “without” created by means of compatible sketch grammars from four Aranea web corpora..

5 Further work

In the near future, we plan to carry out activities within several tracks. Firstly, we would like to improve the quality of the Aranea corpus data itself (by means of better filtration, normalization and deduplication), as well as its morpho-syntactic annotation by means of long-term evaluation of the resulting word sketches. Secondly, we want to include new languages into our Aranea corpus family and to write the respective corpus grammars, at least for the languages taught at Slovak universities. And finally, we plan to tune the global parameters of the compatible sketch grammars, as well as provide language-specific improvements so that the bilingual word sketches provide more relevant results.

6 References

- Baroni, M. – Bernardini, S. (2004). BootCaT: Bootstrapping corpora and terms from the web. In: *Proceedings of LREC'2004*. Lisbon: ELRA ,2004.
- Benko, V. (2013). Data Deduplication in Slovak Corpora. In: *Slovko 2013: Natural Language Processing, Corpus Linguistics, E-learning*. Katarína Gajdošová (Ed.), Adriána Žáková. Lüdenscheid: RAM-Verlag, 2013, pp. 27–39.
- Hong, J. F. – Huang, C. R. (2007). Using Chinese Gigaword Corpus and Chinese Word Sketch in linguistic research. URL: <http://www.ling.sinica.edu.tw/eip/FILES/publish/2007.7.18.93102662.8243902.pdf>.
- Jakubíček, M. (2014). Personal communication.
- Kilgarriff, A. et al. (2004). The Sketch Engine. In: G. Williams and S.Vessier (eds.), *Proceedings of the eleventh EURALEX International Congress EURALEX 2004 Lorient, France, July 6–10, 2004*. Lorient: Université de Bretagne-Sud, pp. 105–116.
- Khokhlova, M. (2010). Building Russian Word Sketches as Models of Phrases. In: *Proc. EURALEX 2010, Leeuwarden*, July 2010.
- Kovář, V. (2013). New features in the Sketch Engine interface. Part 1. In: *SKEW-4 Workshop*, Tallinn, October 2013. URL: https://www.sketchengine.co.uk/documentation/raw-attachment/wiki/SKEW-4/Program/ske_interface_part1.pdf
- Krek, S. – Kilgarriff, A. (2006). Slovene Word Sketches. In: *Proceedings of 5th Slovenian and 1st international Language Technologies Conference 2006*. October 9th – 10th 2006. Jožef Stefan Institute, Ljubljana, Slovenia
- Macoveiciuc, M. – Kilgarriff, A. (2010). The RoWaC Corpus and Romanian Word Sketches. In: *Multilinguality and Interoperability in Language Processing with Emphasis on Romanian*. Edited by Dan Tufis and Corina Forascu. Romanian Academy.
- Petrov, S. – Das, D. – McDonald, R. (2012). A Universal Part-of-Speech Tagset. In: *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*. Istanbul: ELRA, May 2012.
- Piasecki, M. (2007). Polish tagger TaKIPI: Rule based construction and optimisation. *TASK QUARTERLY* 11, No 1–2, 151–167
- Radziszewski, A. – Kilgarriff, A. – Lew, R. (2011). Polish Word Sketches In: Zygmunt Vetulani (ed.) *Human Language Technologies as a Challenge for Computer Science and Linguistics. Proceedings of the 5th Language & Technology Conference*. Poznań : Fundacja Uniwersytetu im. A. Mickiewicza.
- Schmid, H. (1994): Probabilistic Part-of-Speech Tagging Using Decision Trees. In: *Proceedings of International Conference on New Methods in Language Processing*, Manchester, UK.
- Srdanović, E. I. – Erjavec. T. – Kilgarriff, A. (2008). A web corpus and word-sketches for Japanese. In: *Journal of Natural Language Processing* 15/2, 137– 159. (reprinted in *Information and Media Technologies* 3/3, 2008, 529– 551)
- Suchomel, V. – Pomikálek, J. (2012). Efficient Web Crawling for Large Text Corpora. In: *7th Web as Corpus Workshop (WAC-7)*, Lyon, 2012.
- Tiberius, C. – Kilgarriff, A. (2009). The Sketch Engine for Dutch with the ANW corpus. In: *Fons Verborum, Festschrift for Fons Moerdijk*. Instituut voor Nederlandse Lexicologie, Leiden, The Netherlands.
- Whitelock, P. (2010). Personal communication.

Appendix

To demonstrate the compatible word sketches, we present screen shots for the preposition “without” in four languages (English, French, German, and Russian). Prepositions belong to word classes that are usually either not covered by the respective traditional sketch grammars at all, or that produce a limited number of output word sketch tables only.

For all languages involved, we can observe the typical binary collocations with noun and verbs. The collocations with adjectives usually form a multi-word expression that is not fully displayed in the word sketches, but many of those can be easily recognized even without going into the actual concordances.

Note: Due to the longer adjectives in Russian, the interesting table with verbal collocates did not fit onto the screen.

without (*non-verb*) Araneum Anglicum Maius (En Web 1.2.01) 1,20 G freq = **413459** (344.5 per million) Click on collocates in boldface to get multi word sketches.

X	X/Y, X/Y 31449 -0.0	X/Y Cj X/Y 23389 -0.0	YX 291638 -0.1	XY 379129 -0.1
Pp(X) 413459 -0.5	except 141 3.61	within 753 3.35	complete 1130 4.76	limitation 4001 7.46
	whereas 36 2.92	with 7024 2.62	reproduce 355 4.63	doubt 4084 7.08
	unless 111 2.83	between 237 1.13	viagra 318 4.54	notice 3171 6.71
	albeit 11 2.65	under 111 0.41	survive 662 4.45	hesitation 1133 6.47
	although 97 1.76	beyond 28 0.39	incomplete 190 4.08	regard 2352 6.21
	though 160 1.64	against 80 0.17	function 323 4.05	prescription 1401 6.16
	despite 42 1.57	except 12 0.09	live 2861 4.0	express 1024 6.13
	upon 118 1.25	whether 44 -0.0	impossible 385 4.0	prior 2708 6.1
	like 402 1.21	from 687 -0.06	exist 1153 3.98	permission 1676 6.09
	via 45 1.2	at 662 -0.08	possible 1451 3.88	Borders 837 6.08
	because 352 1.2	toward 17 -0.13	cialis 164 3.8	compromise 1176 6.06
	till 14 1.17	upon 45 -0.13	detention 184 3.79	delay 1265 5.99
	within 166 1.17	behind 25 -0.14	proceed 288 3.75	prejudice 920 5.93
	unto 16 1.16	through 145 -0.17	reporter 213 3.58	fear 2381 5.87
	while 253 1.11	into 210 -0.18	die 672 3.39	consent 1235 5.83
	under 169 1.01	for 1374 -0.23	even 2829 3.38	exception 1332 5.7

Nn X 292904 -0.1	X Nn 401566 -0.1	Aj X 75340 -0.0	X Aj 110090 -0.1	Vb X/X Vb 704673 -0.1
viagra 467 5.09	consent 4488 7.65	incomplete 195 5.4	express 1335 7.7	compromise 1273 5.53
cialis 255 4.43	permission 4964 7.61	impossible 737 5.4	written 2518 7.36	ca 4366 5.46
detention 223 4.07	limitation 4115 7.46	complete 1228 5.13	prior 3214 6.69	sacrifice 1016 5.21
reporter 231 3.7	notice 4362 7.14	meaningless 82 4.38	undue 371 6.38	leave 5606 4.96
life 3435 3.61	doubt 4205 7.09	generic 158 4.33	parental 419 5.95	worry 1337 4.91
buy 129 3.56	hesitation 1212 6.49	possible 1826 4.31	proper 1536 5.95	could 10736 4.89
sentence 259 3.4	delay 1674 6.34	useless 106 4.24	adequate 731 5.95	resort 701 4.83
prescription 177 3.38	regard 2478 6.26	worthless 65 4.15	explicit 451 5.81	can 29062 4.82
taxation 116 3.35	prescription 1463 6.17	indefinite 56 4.11	further 2508 5.2	survive 1130 4.72
imprisonment 111 3.29	Borders 818 5.97	lonely 58 3.58	foregoing 174 5.16	allow 4815 4.66
sex 362 3.28	authorization 989 5.94	consecutive 68 3.51	slight 357 5.06	depart 694 4.63
Viagra 91 3.26	fear 2517 5.92	usable 45 3.49	formal 483 4.85	live 4724 4.62
party 829 3.24	prejudice 947 5.91	able 1160 3.46	added 238 4.83	imagine 1196 4.57
nothing 650 3.18	exception 1407 5.74	difficult 466 3.39	excessive 208 4.63	lose 2519 4.55
anything 655 3.14	ado 628 5.65	unthinkable 30 3.37	conscious 294 4.62	go 12313 4.53
day 2643 3.11	warning 1087 5.63	inconceivable 28 3.33	additional 1159 4.61	would 13830 4.45

Av X/X Av 112453 -0.1	ZX 226890 -0.0	XZ 304244 -0.1	Pp X 92448 -0.0	X Pp 99334 -0.0
whatsoever 340 5.61	those 2787 4.04	any 27935 6.89	albeit 112 4.95	into 2483 3.37
ever 3020 5.13	any 3658 3.97	them 5997 4.18	onto 173 3.62	except 139 3.35
indefinitely 128 4.83	themselves 630 3.63	a 49871 3.78	into 2579 3.42	whether 379 3.03
freely 262 4.82	off 1294 3.44	him 2262 3.62	through 1679 3.35	of 38204 2.97
necessarily 454 4.82	i.e. 147 3.42	an 7055 3.48	although 240 2.98	from 5422 2.92
overly 155 4.72	another 911 3.38	some 2749 3.32	than 1688 2.82	unless 119 2.75
even 6980 4.71	them 3308 3.33	it 13245 3.23	across 305 2.73	through 1091 2.72
anywhere 360 4.45	itself 387 3.05	their 6214 3.21	because 971 2.64	upon 330 2.67
too 2634 4.42	yourself 311 3.0	its 3091 3.17	like 1071 2.61	about 2781 2.66
properly 346 4.36	it 11199 2.99	yourself 354 3.1	from 4267 2.58	like 1093 2.63
actually 1420 4.2	these 1877 2.96	the 79951 3.08	except 77 2.52	for 9121 2.5
explicitly 142 4.17	this 6489 2.93	your 5936 3.08	on 6479 2.5	until 310 2.44
easily 607 4.11	himself 318 2.88	her 2447 3.0	for 9024 2.48	behind 140 2.23
overboard 67 4.08	him 1310 2.85	either 373 2.93	at 3684 2.39	on 5230 2.19
completely 592 4.08	no 1735 2.84	whom 248 2.85	upon 269 2.39	at 3079 2.13
physically 161 4.06	to 40462 2.79	these 1714 2.82	whereas 33 2.32	till 35 2.08

ohne Araneum Germanicum Maius (De Web 1.2.01) 1,20 G freq = **592115** (493.4 per million) Click on collocates in boldface to get multi word sketches.

X	X/Y, X/Y	34121 0.2	X/YC X/Y	41201 0.2	YX	323276 0.3	XY	483395 0.4
Pp(X) 505192 3.0	dass	1892 4.32	einschließlich	57 3.19	Kredit	2342 6.4	vorherig	4624 7.62
Cj(X) 86923 0.9	ob	270 3.32	seitens	40 2.8	jederzeit	1222 5.36	Schufa	3481 7.57
	oder	1064 3.25	mit	8640 2.61	Reporter	414 5.2	Zweifel	3621 7.26
	d.h.	21 3.05	inklusive	34 2.33	gänzlich	475 5.11	Rücksicht	2860 7.24
	wie	659 2.85	zwischen	382 1.97	ganz	8310 4.75	Weiteres	1810 6.82
	und	5054 2.84	trotz	71 1.8	Fahren	306 4.63	Einschränkung	2002 6.49
	einschließlich	39 2.69	mittels	40 1.8	völlig	929 4.54	Angabe	4455 6.46
	bzw.	54 2.23	binnen	12 1.79	Fass	236 4.34	gesondert	1607 6.35
	inklusive	30 2.2	außer	30 1.55	Abnehmen	236 4.32	Zustimmung	2145 6.3
	außer	45 2.17	samt	12 1.54	Girokonto	288 4.28	Abzug	1372 6.19
	weil	61 2.1	aufs	22 1.53	Handy	539 4.26	Gewähr	1324 6.19
	als	464 2.06	außerhalb	39 1.51	Kreditkarte	348 4.24	Problem	6810 6.14
	infolge	12 1.72	wider	11 1.47	Rechnung	486 4.13	Umweg	1192 6.11
	indem	12 1.6	ans	18 1.4	Leben	2437 4.12	Unterbrechung	1176 6.11
	aufs	20 1.43	ob	71 1.38	allerdings	1516 4.09	ausdrücklich	2164 6.1
	sondern	55 1.37	gegen	269 1.36	Pfanne	226 4.03	Behinderung	1596 5.93

Nn X	311082 0.3	X Nn	606564 0.5	Aj X	89134 0.2	X Aj	200309 0.5	Yb X/X Vb
Kredit	2584 6.56	Zustimmung	5620 7.5	gedruckt	973 7.25	vorherig	4938 8.37	auskommen
Publikation	1004 5.87	Schufa	3453 7.29	gänzlich	485 6.2	gesondert	1704 7.27	gestatten stat
Reporter	415 5.26	Zweifel	4051 7.21	völlig	1022 4.96	ausdrücklich	3061 7.06	funktionieren
Fahren	337 4.82	Ankündigung	3067 7.07	anwaltlichen	83 4.73	nennenswert	811 6.66	verändern
Rechnung	739 4.76	Rücksicht	3041 7.07	gesamt	1581 4.62	schriftlich	2599 6.55	verlaufen
Handy	708 4.67	Aufwand	3653 6.93	erhoben	133 4.52	lästig	639 6.18	verlieren
Girokonto	349 4.55	Einschränkung	2510 6.59	Prepaid	112 4.5	erkennbar	856 6.09	leben
Kreditkarte	402 4.48	Genehmigung	2522 6.55	personenbezogen	248 4.5	störend	448 5.8	laufen
Abnehmen	247 4.43	Weiteres	1818 6.52	komplett	820 4.38	weit	11902 5.58	zögern
Fass	239 4.41	Angabe	4874 6.49	selbstverständlich	433 4.38	zusätzlich	3407 5.58	dürfen
Leben	2895 4.37	Grund Gründen	3484 6.45	viagra	61 4.27	möglich	4803 5.48	verlassen
Angebot	1875 4.2	Einwilligung	1781 6.31	vollkommen	226 4.26	ersichtlich	345 5.34	nachdenken
Tarif	439 4.19	Problem	7956 6.31	undenkbar	64 4.24	fremd	827 5.32	bleiben
Pfanne	242 4.17	Abzug	1758 6.29	kommerziell	153 4.22	unnötig	449 5.28	können
Abmahnung	267 4.13	Umweg	1353 6.01	verlinkten	97 4.19	aufwendig	523 5.21	kündigen
Baufinanzierung	197 4.05	Gewähr	1384 5.99	berufsmäßig	52 4.19	finanziell	1189 5.19	überstehen

Av X/X Av	173826 0.4	Z X	285861 0.3	X Z	356121 0.4	Pp X	67488 0.1	X Pp	111231 0.3
jemals	724 6.31	man	8344 4.02	jegliche	4363 7.76	dank	140 3.13	seitens	198 4.65
jederzeit	1696 6.07	14	690 3.93	irgendwelche	1088 6.08	pro	258 2.87	durchs	80 3.87
vorher	1598 5.93	diese	1417 3.91	dabei	6094 5.38	bei	5283 2.86	ans	124 3.79
jedoch	4386 5.15	nicht	20169 3.85	irgendeine	556 4.95	einschließlich	51 2.85	auf	15175 3.71
allerdings	2735 5.01	wer	1280 3.83	allzu	476 4.74	gegenüber	221 2.71	außer	138 3.48
ganz	9471 4.96	keine	4041 3.74	dafür	1743 4.61	inklusive	49 2.7	aufs	105 3.43
niemals	434 4.86	daher	740 3.57	jede	5801 4.4	trotz	135 2.67	von	16501 3.42
leider	1386 4.66	solche	1093 3.54	zu	52715 4.36	nach	2766 2.58	durch	3718 3.34
irgend	201 4.65	sie	6700 3.53	nichts	1282 4.16	innerhalb	209 2.49	binnen	56 3.32
kaum	1261 4.63	niemand	257 3.45	keine	4669 3.94	für	7039 2.48	über	3330 3.0
freilich	204 4.62	eine	37041 3.41	solche	1384 3.85	wegen	200 2.46	an	7594 2.92
also	3264 4.58	jemand	317 3.39	jedwede	159 3.75	seit	561 2.46	innerhalb	277 2.84
bisher	1131 4.56	was	2380 3.32	darauf	740 3.74	zeit	20 2.46	per	225 2.79
auch	35534 4.56	er	5165 3.25	jemand	424 3.71	von	8310 2.43	in	22140 2.69
sofort	1136 4.53	nichts	655 3.24	viele	4400 3.71	mittels	63 2.36	mit	8122 2.52
meistens	430 4.52	es	9643 3.23	darüber	676 3.69	binnen	20 2.23	gegen	606 2.5

SANS Araneum Francogallicum Maius (Fr Web 1.2.02) 1,23 G freq = **1023016** (829.5 per million) Click on collocates in boldface to get multi word sketches.

X	X/Y, X/Y 165883 -0.2	X/Y Cj X/Y 89170 -0.1	Y X 638024 -0.1	X Y 878596 -0.1
Pp(X) 1012928 -0.6	sauf 248 3.9	avec 11553 4.27	reporter 975 5.06	doute 85421 9.99
Cj(X) 10088 -0.0	jusque 599 3.29	hors 104 2.71	non 5532 4.86	cesse 20493 9.02
	malgré 217 3.13	soit 77 2.17	c'est-à-dire 804 4.49	oublier 15100 7.48
	hormis 51 3.03	malgré 78 1.8	fonctionner 1169 4.48	fil 7755 7.07
	hors 147 3.01	parce 36 1.58	organisme 1212 4.2	précédent 6972 6.97
	avec 4598 2.94	entre 519 1.56	répéter 637 4.16	préavis 3701 6.96
	depuis 869 2.89	sauf 41 1.52	rester 3613 4.16	autorisation 5545 6.84
	chez 541 2.89	en 5765 1.48	accepter 1380 4.15	autant 9143 6.81
	vers 626 2.88	sous 284 1.45	presque 971 4.14	relâche 3002 6.72
	pendant 439 2.83	pendant 162 1.45	vivre 2985 4.14	compter 12738 6.67
	sous 725 2.77	contre 250 1.37	consommer 669 4.14	frontière 4720 6.61
	dès 262 2.63	à 7766 1.37	interdire 864 4.01	conteste 2740 6.6
	à 18464 2.62	dès 98 1.29	mourir 943 3.98	faille 3165 6.6
	au 7379 2.59	du 9242 1.28	vétérinaire 364 3.94	consentement 2514 6.28
	devant 323 2.58	jusque 137 1.22	voiture 857 3.87	moindre 3757 6.26
	en 12377 2.58	de 23057 1.17	dérouler 816 3.85	limite 4816 6.24

Nn X 526473 -0.1	X Nn 969635 -0.1	Aj X 118195 -0.1	X Aj 167461 -0.1	Vb X/X 1329224 -0.1
reporter 920 5.16	doute 96148 10.11	vétérinaire 260 5.06	lucratif 2851 8.2	oublier 16163 7.31
organisme 1441 4.53	cesse 20510 8.93	impossible 736 5.01	préalable 4847 7.94	compter 14308 6.61
journée 1822 4.3	fil 7829 7.02	tierce 139 4.55	moindre 5688 7.7	parler 12885 5.91
voiture 1075 4.28	préavis 3832 6.89	modifiable 97 4.54	expresse 1296 7.5	soucier 2432 5.71
nuit 1242 4.16	autorisation 6013 6.88	gratuit 747 4.34	nul 2863 7.15	tarder 2551 5.71
Internet 1353 4.07	précédent 6605 6.82	accessible 517 4.29	apparent 1106 6.93	laisser 9961 5.61
réseau 2004 4.07	relâche 3004 6.59	possible 1898 4.26	frais 3179 6.59	attendre 7517 5.61
acceptation 367 4.02	faille 3200 6.49	utilisable 108 4.08	frontière 1676 5.97	perdre 5087 5.31
amour 1131 3.91	conteste 2760 6.47	remboursable 69 3.93	égal 1040 5.64	pouvoir 54251 5.31
jour 4495 3.88	limite 5331 6.32	sexuel 320 3.89	réel 2466 5.61	passer 13149 5.21
médecin 739 3.88	surprise 3308 6.29	correct 133 3.85	gras 748 5.6	hésiter 2779 5.11
licencierement 300 3.86	consentement 2647 6.24	partiel 173 3.8	fixe 640 5.5	savoir 12005 5.11
connexion 400 3.79	souci 3214 6.03	réalisable 63 3.71	excessif 420 5.44	rester 7627 5.01
monde 3438 3.77	hésitation 2068 6.02	estre 81 3.7	supplémentaire 1195 5.42	vivre 6437 5.01
sexe 522 3.74	arrêt 3353 6.0	immédiat 197 3.67	valable 526 5.37	bouger 1713 5.01
gens 1270 3.74	frontière 3229 5.98	inconcevable 53 3.66	explicite 298 5.17	regarder 3542 5.01

Av X/X Av 287079 -0.1	Z X 526005 -0.1	X Z 554614 -0.1	Pp X 258101 -0.1	X Pp 349901 -0.1
autant 10797 7.41	ceci 1004 4.96	aucun 36125 8.34	devant 691 3.61	jusque 1966 4.88
toutefois 3317 6.99	cela 6226 4.58	quoi 3522 5.73	chez 870 3.52	pour 27339 4.31
jamais 11860 6.7	la le 9272 4.47	OGM_ 963 5.49	vers 992 3.49	avec 9295 3.94
trop 11522 6.58	la 2222 4.32	la le 10965 4.71	du 42461 3.48	pendant 1000 3.89
forcément 1325 5.84	se 32382 4.28	la 2817 4.65	à 31897 3.41	sur 13933 3.83
nécessairement 855 5.74	un 77428 4.26	toi 1128 4.64	de 107821 3.4	à 42038 3.8
rien 6754 5.71	on 10142 4.24	eux 1697 4.46	pendant 677 3.39	depuis 1684 3.78
vraiment 3859 5.5	elle 9775 4.24	lequel 3094 4.41	après 1187 3.35	de 137573 3.75
non 7636 5.4	il 25892 4.1	lui 6975 4.33	dans 11797 3.18	au 15714 3.68
priori 564 5.33	votre 5454 4.07	y 7918 4.32	pour 11453 3.06	par 11444 3.67
réellement 926 5.33	vous 12695 4.07	leur 7251 4.22	avec 4790 2.99	devant 715 3.59
c'est-à-dire 989 5.31	tu 2004 3.99	se 30956 4.21	sur 7429 2.92	dans 15547 3.58
exprès 387 5.23	ce 38418 3.91	me 6916 4.14	au 9114 2.9	vers 1043 3.52
presque 1466 5.05	y 5689 3.85	son 19423 4.14	derrière 177 2.85	contre 1189 3.51
préalablement 416 5.02	leur 5436 3.81	que 19955 4.11	contre 726 2.83	envers 242 3.43
apparemment 419 4.92	qui 24786 3.8	votre 5045 3.96	en 14810 2.83	quant 240 3.42

X	X/Y, X/Y	16565	-0.2	X/Y Cj	9119	-0.1	YX	424569	-0.3	XY	620883	-0.3		
Pp(X)	749749	-2.3	безо	15	3.46	вне	76	2.84	обойтись	12297	8.55	исключение	9151	7.69
nom(X)	6477	-0.0	вне	90	3.05	у	526	0.2	обойтись	5964	8.15	сомнение	8080	7.63
gen(X)	737804	-1.3	кроме	230	2.46	помимо	12	-0.07	остаться	14808	6.92	весть	4171	7.49
dat(X)	1640	-0.0	ради	54	1.86	вместо	16	-0.28	пропасть	2590	6.68	попечения	2636	7.02
acc(X)	1003	-0.0	посредством	18	1.21	против	44	-0.28	невозможный	2621	6.53	преувеличение	2256	6.76
loc(X)	2266	-0.0	вместо	39	0.99	вокруг	14	-0.74	оставить	4941	6.24	особый	9784	6.71
ins(X)	559	-0.0	путем	24	0.98	для	445	-0.81	невозможно	2174	6.07	труд	8844	6.61
			помимо	24	0.91	со	90	-0.84	немыслимый	954	5.93	лишний	4748	6.58
			среди	74	0.65	до	163	-0.86	оставлять	1678	5.37	учет	6474	6.55
			выше	16	0.6	из-за	22	-0.91	пропавшими	485	5.2	малое	1960	6.51
			возле	11	0.56	от	294	-1.01	практически	2233	5.14	предварительный	2943	6.33
			от	782	0.4	после	78	-1.19	оставаться	3850	5.12	малейший	2032	6.25
			для	989	0.35	среди	18	-1.39	прожить	741	4.98	присмотр	1582	6.25
			из	811	0.33	из	237	-1.44	жить	4890	4.89	посредник	1786	6.21
			у	552	0.27	с	296	-2.38	вообще	2127	4.79	согласие	2584	6.19
			внутри	22	0.27				почти	1985	4.67	разбор	1567	6.16

Nn X	364089	-0.3	X Nn	844256	-0.4	Aj X	88968	-0.2	X Aj	194940	-0.3	Vb Xj
оставление	423	5.12	сомнение	9447	7.64	немыслимый	1008	7.55	малейший	2061	7.24	обойтись
кредит	1386	4.88	исключение	9375	7.52	пропавшими	471	7.31	предварительный	3438	7.19	обойтись
репортер	312	4.68	весть	4173	7.12	невозможный	2767	7.25	видимый	1687	7.12	остаток
дым	377	4.51	труд	10846	6.8	пропавшим	176	5.94	лишний	4940	7.08	оставление
участок	1432	4.4	согласие	4377	6.69	худой	324	5.78	особый	10017	6.96	пропуск
кофе	432	4.26	попечения	2737	6.66	мыслимый	187	5.55	посторонний	1570	6.9	оставление
чай	565	4.17	преувеличение	2693	6.56	предпринимательский	249	5.14	уважительный	761	6.6	прожечь
лицо	2825	4.14	ведомо	2464	6.49	минеральный	263	4.66	должный	1241	6.53	жизнь
секс	396	4.02	учет	6771	6.48	неполный	156	4.63	невозможный	1714	6.32	оставление
отпуск	412	4.01	разрешение	4377	6.35	исковый	124	4.43	дополнительный	4464	6.13	предложение
наличные	219	3.98	ограничение	3996	6.31	ровный	205	4.27	излишний	660	5.86	смочить
квартира	1366	3.95	усилие	4131	6.25	гладкий	146	4.12	надлежащий	563	5.67	мысль
бой	586	3.93	малое	1968	6.12	земельный	290	3.92	немыслимый	397	5.53	повзросление
жизнь	5441	3.92	потеря	3707	6.11	апелляционный	83	3.91	хирургический	475	5.48	мочить
столбик	187	3.91	ущерб	2715	6.09	длительный	397	3.85	единый	2288	5.27	спрашивать
заработок	387	3.84	вмешательство	2359	6.03	Послеоперационный	39	3.82	специальный	3111	5.13	выжить

Av X/X Av	100363	-0.2	ZX	231029	-0.2	XZ	226728	-0.3	Pp X	70963	-0.1	X Pp	79547	-0.1
невозможно	7192	8.32	куда	1227	5.54	всякий	21391	8.49	ко	271	3.49	со	2113	3.69
немыслимо	439	6.99	тут	1299	4.65	какой-либо	10149	8.07	сверх	23	2.68	над	724	3.57
нельзя	4436	6.65	даже	4526	4.49	никуда	887	6.39	за	2747	2.45	для	6323	3.02
извне	309	5.98	ж	324	4.4	никак	1459	6.18	на	12506	2.44	на	16431	2.83
практически	3089	5.85	никак	412	4.35	таковой	1174	6.03	обо	39	2.41	к	4975	2.75
трудно	1261	5.83	поэтому	1540	4.28	оное	252	5.04	до	1499	2.34	с	8158	2.4
сложно	797	5.79	ведь	1216	4.15	онный	281	4.87	во	994	2.33	от	3131	2.4
вообще	3670	5.78	теперь	1167	4.12	чей-либо	190	4.68	посреди	19	2.29	о	2583	2.28
скучно	252	5.77	здесь	1347	4.05	какой	2738	3.89	в	23004	2.25	из-за	206	2.24
вовсе	983	5.53	не	40489	3.96	ничто	1527	3.82	через	483	2.21	об	527	2.21
бесплатно	507	5.37	бы	4480	3.82	то	13929	3.79	при	1481	2.2	свыше	40	2.21
можно	13945	5.32	весь	8436	3.81	они	13906	3.63	к	3292	2.16	около	223	2.2
почти	2600	5.25	тоже	1295	3.73	она	9102	3.56	из	2797	2.11	до	1265	2.09
желательно	325	4.96	туда	277	3.73	он	15472	3.55	об	486	2.09	пред	27	2.07
возможно	826	4.93	как-то	277	3.62	пять	539	3.45	сквозь	29	2.09	между	375	2.04
тяжело	302	4.91	просто	1910	3.61	ваш	2980	3.43	под	609	2.06	в	19486	2.01