# Corpus collection and analysis for the linguistic layman: The Gromoteur

## Kim Gerdes

LPP, Université Sorbonne Nouvelle & CNRS

## Abstract

This paper presents a tool for corpus collection, handling, and statistical analysis: the Gromoteur. The paper explains the scientific needs that lead to the development of the tool and its main characteristics. Different usage schemes are explored in which the tool can contribute to simple NLP tasks that are often beyond the reach of the common linguist. In particular corpus collection from online resources or heterogeneous file structures becomes an easy task with the Gromoteur. But also tokenization, lemmatization, or part of speech annotation of various languages are a matter of a few clicks in this tool. The Gromoteur also allows for simple statistical analysis of the collected language data, including a section map, specificity computation for any set of sections, and collocation computation for n-grams.

**Keywords :** corpus collection, web crawler, online corpus, the Web as corpus, tagger, lemmatizer, statistical analysis, specificity computation, collocations
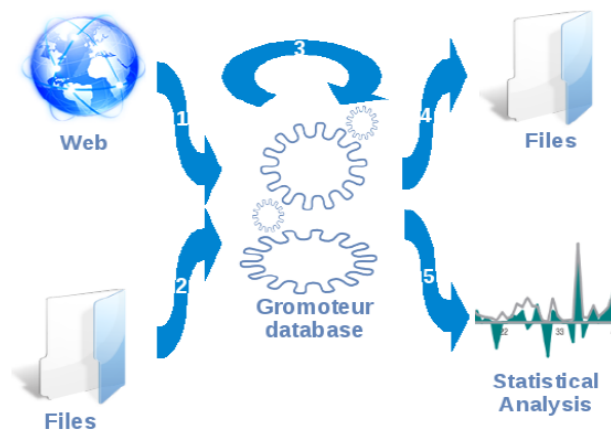
## 1. Introduction

The sudden abundance of linguistic data, written and spoken, through the World Wide Web, has brought linguistics from data scarcity and the obligation to rely on intuition  in a short time to the possibility to verify and even quantify many of the formally purely theoretical postulates. In spite of this paradigm change, common linguists run up against a multitude of technical roadblocks to use "Big Data" for their research. If they want to use large corpora, they mainly end up either with simple results provided by all purpose search engines or with specific tools that have predefined corpora like Google's Ngram viewer or the WebCorp Linguist's Search Engine. A simple task like downloading all available articles of an online newspaper or a blog, coding them uniformly, cleaning them and making simple measures of the distribution of words is beyond the reach of a linguist not trained on technical aspects of encoding and scripting.

The Gromoteur attempts to remedy this need. It is a tool for linguists that gives easy access to textual corpora. It allows to get pages from the Web or from local files, treat them, analyze them, and output results. It is tailored for common needs of linguists that want to work on data they either have in various file formats on the hard drive, or that they have access to on the Web. All configurations and usage is controlled from a graphical interface, and in common usage scenarios no scripting or tweaking with configuration files is necessary.

## 2. Usage Schemes

The following figure illustrates the different paths linguistic data can take in the Gromoteur: The Gromoteur can take data from the internet (1) and from local files (2), extract parts of pages, sort, and clean the data (3), export the data to different text files and concordancer formats, and a simple integrated tool allows for simple statistical analyses.

## 2.1. The crawler

The central part of the Gromoteur is the crawler with a wide range of configuration options.

After creating a database (or using an existing one), the user can fill it with data from the internet. A dialog wizard guides the user through the configuration. The dialog distinguishes between *normal* and *expert* configuration options, the latter being accessible by clicking on the "expert" button of the dialog. The normal user is thus exposed only to a few essential options among the *expert* options presented below, the simple options being only the website to start crawling, the restriction on the websites to visit, and the number of pages to retrieve.

The first question is only relevant if the user has filled the database with the crawler before, and wants to complete the database now:

- new data is appended to the tables
- only equal URLs are overwritten
- the database is completely erased and started from zero
- the search starts with URLs that were collected during the last crawl and remained to be explored

The user is then presented with the choice between the two basic ways of finding data on the internet:

1. by giving seed URLs as a starting point
2. by providing search terms for a search using the Bing search engine

The seed URLs can either be pasted directly into the in the configuration dialog or provided in a file containing a URL per line The number of results for any Bing search are directly given in the interface and the result pages can also be opened and examined in a browser by a simple click. The Bing query syntax allows for all the common restrictions on language, URL, and file type (see http://onlinehelp.microsoft.com/en-us/bing/ff808421.aspx for details). Moreover the Bing locale can be set automatically or manually. The user can decide to just download the search results or to use the results as seeds from which to keep crawling. To use the Bing search API, Gromoteur is preconfigured with a Bing user ID, but in case the users of Gromoteur want to make extensive use of this feature they are urged to get their own Bing user ID, in order not to attain the daily limit of searches per user ID.

The next screen in the configuration wizard asks to put restrictions on the way to walk the web. As shown in the screenshot below, this includes the order in which the pages are taken,

restrictions on the URL and the terms on the page, the depth of the crawl, and the file types. The Gromoteur can also check the language of each page, using the letter trigram based Language Detector (http://elizia.net/languageDetector/), thus staying on a unilingual search path. This is particularly interesting for the corpora creation of rare languages where the examples pages can be found by a keyword search, which is then controlled by the language testing.



The restrictions on the URLs and the web content can be a simple word or part of the URL, but they are in fact Perl/Python style regular expressions. Although regular expressions are a quite powerful tool allowing to go beyond simple keyword or URL matching, they turn out to be the biggest obstacle for the non-technical oriented users of the Gromoteur in the crawler configuration. The following dialog page controls the quantity and speed of the crawl.

The final configuration screen of the crawler controls the number of threads of the crawl, the usage of proxies, the identification of Gromoteur when crawling, and the handling of the robots.txt file.[1]



The number of threads can significantly alter the speed with which the pages are taken from the Web. However, it is difficult to find the optimal configuration to enhance the speed. The main factors for the download speed are of course the user's internet connection, the distance to the webserver and the webserver's speed. So in theory, a higher number of parallel threads can only improve the speed, but in reality, a very high number of threads can slow down the hole process because of a slower performance of the crawler, but also because of the reaction of the server, in particular if the download is taken place from one unique server (and not from a variety of servers), as the multiple downloads can be interpreted as a denial of service attack by the webserver.
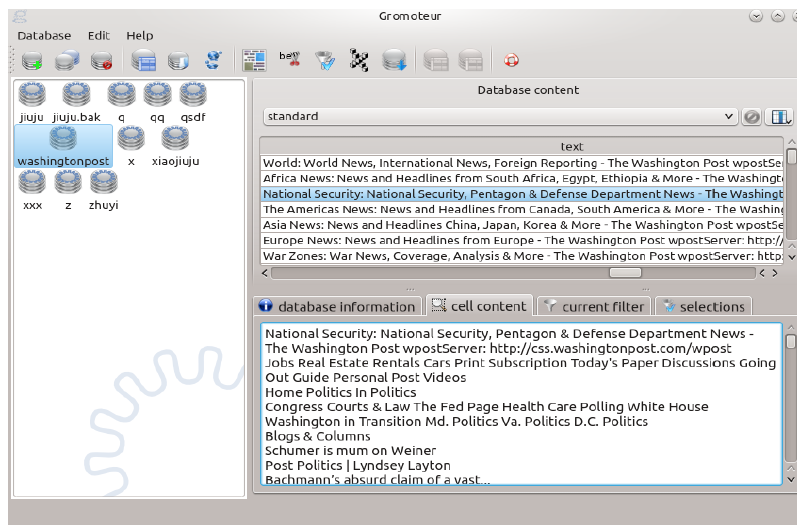
## 2.2. The file importer

Files can be imported by pointing to folder from which all text, html, or pdf files are imported. This can be very helpful for users that have their data in form of pdf ebooks are similar formats that are difficult to handle without some specialized tools. Alternatively textual data can be imported from a simple tab-separated (csv) format, as can easily be produced from a spreadsheet.

Both in the crawling and file importing process, a wide range of heuristics is applied for the cleaning of the import files and the transformation into uniform UTF-8 forms stored in a Sqlite data base. The database stores the source files as well as the text files, the text files are stored using the full text search module (FTS-4) that provides a full-text index allowing for very fast retrieval also in very large databases (of many gigabytes of text).

## 2.3. Cleaning, sorting, analyzing

Gromoteur's main window shows the different databases in the left panel and the the content of the selected database in the right panels. Each database has different columns whose complete content is shown when clicked upon. Different selection filters can be applied to any of the columns.

---

[1]The robots.txt file is a Web convention allowing web site owners to control the behavior of Web crawlers when visiting their websites.
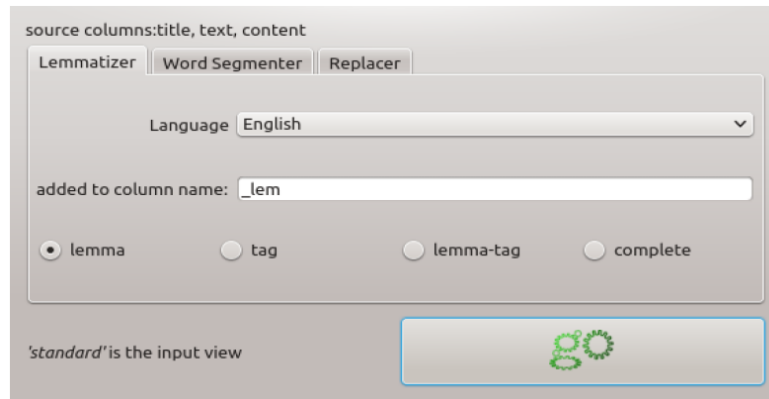
A very important feature of the Gromoteur is the possibility to extract parts of webpages into separate columns. A graphical tool shows the downloaded webpage, and clicking on the parts of the page that are interesting for the linguistic analysis selects and highlights all similar sections of the page.



The tool automatically show the specific HTML tags of the selected parts of the webpage. Those tags can also be adjusted manually. After testing on a number of pages, the linguist can decide to apply this extraction of segments on the whole corpus, resulting in a new column that only contains the central part of the page, for example the content of the article and not
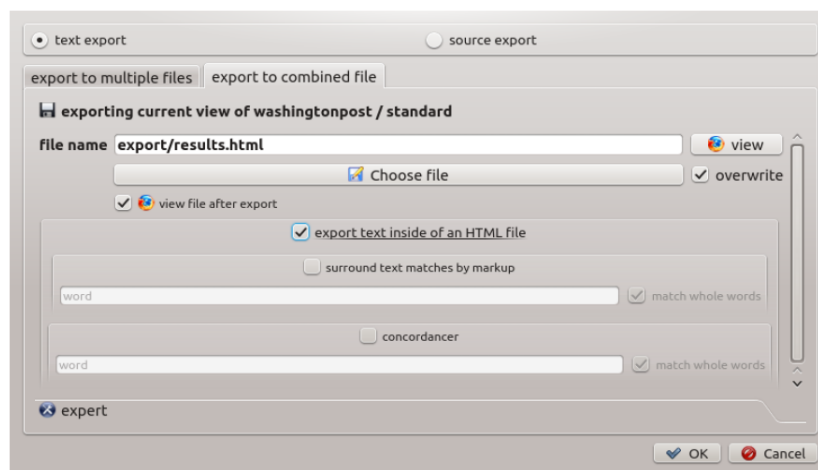
the surrounding information or advertisement. Moreover, a selection of tools allows to word-segment Chinese texts using Urheen[2] or Mmseg[3], lemmatize and tag English[4] and Chinese texts, and do global corrections using regular expressions. Gromoteur comes with the taggers of the Pattern module (De Smedt & Daelemans, 2012) which includes simple statistical taggers of varying quality for English, Spanish, German, French, Italian, and Dutch.



A simple replacement option relying on regular expressions allows to make changes and corrections across the whole corpus.

## 2.4. Exporting the texts

Once the collection, cleaning, analysis, and correction process is finished and at least one of the column contains desirable linguistic data, there are different ways of exporting the column(s) for further inspection or use in other tools: After having selected the desired columns of the table in the main window, the export dialog offers two basic options: export to multiple files or to a unique combined file. The multiple file export provides the choice of including various separations between columns and naming options. The combined file export can just combine the data in a unique text file, but it can also produce html files with highlighted words as well as tables in a concordancer style.



---

---
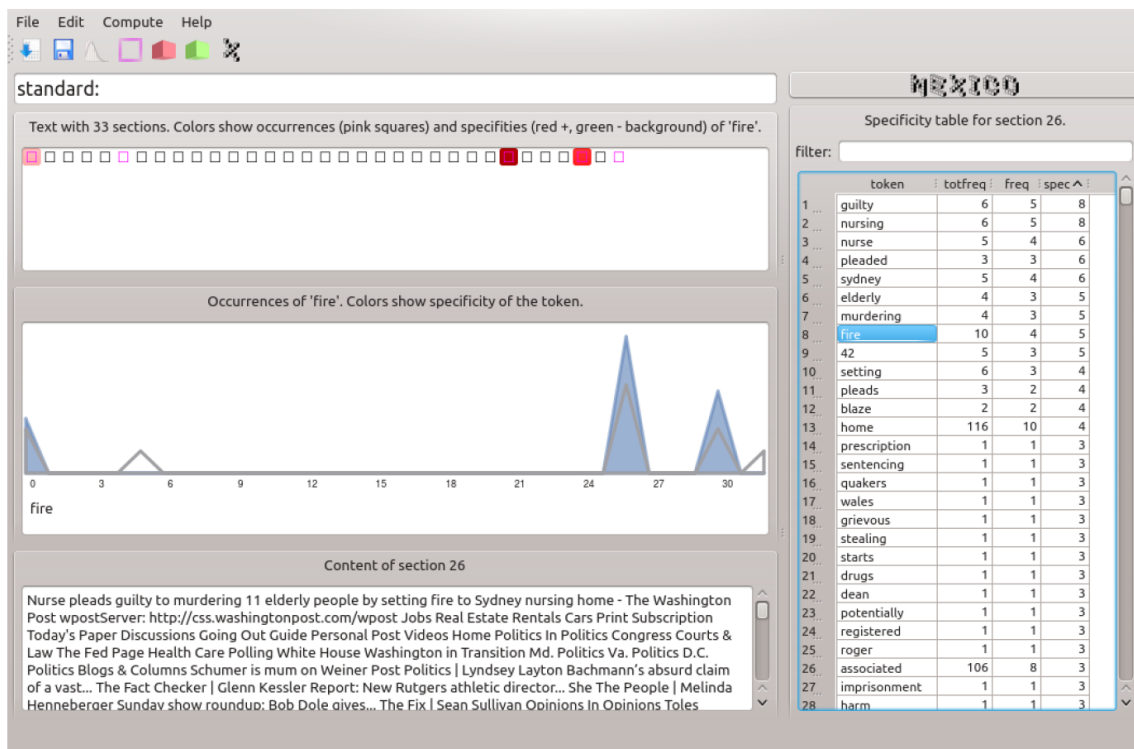
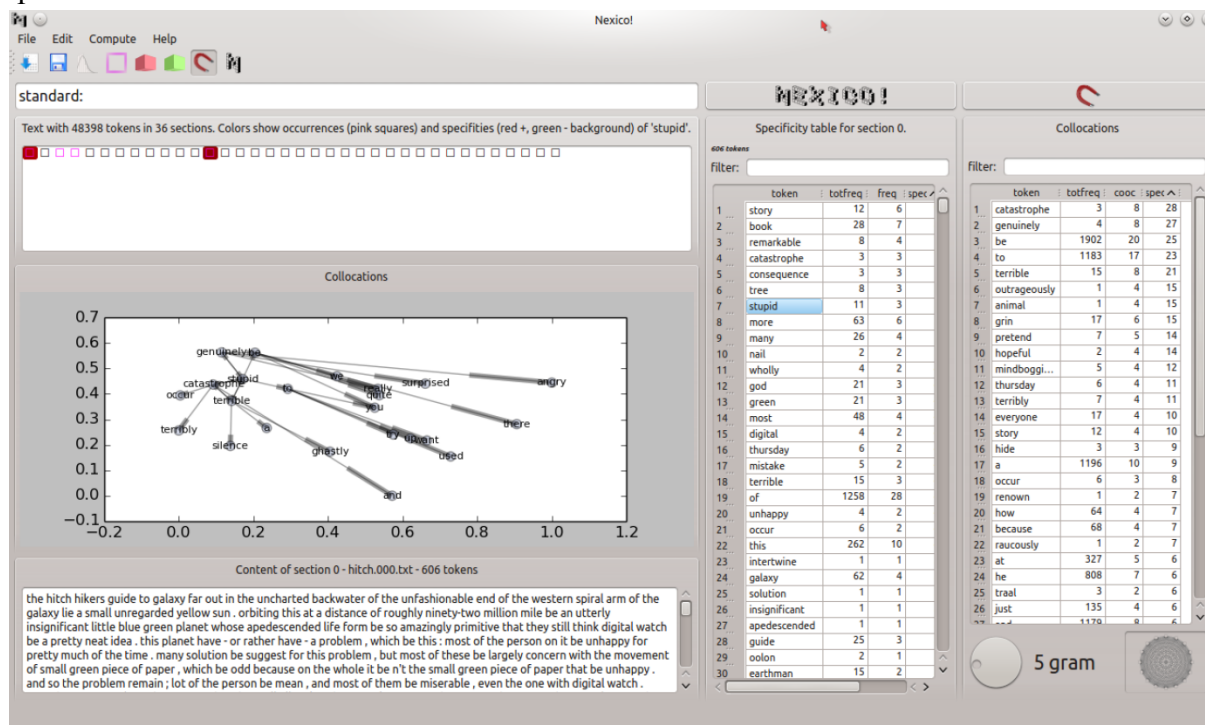| science 73784 What Modern Humans | Can | Learn From The Neanderthals' Extinction It |
| that modern humans inherited from them. How | can | we avoid meeting the Neanderthals' fate? T |
| bate over what happened to Neanderthals | can | be boiled down to two dominant theories: E |
| ca, nicknamed Mitochondrial Eve. If all of us | can | trace our roots back to one African woman, |

## 2.5. Analyzing the textual data

The Gromoteur also includes a simple statistical tool named Nexico as it resembles the Lexico 3 tool (Fleury and Salem, 2002) in its use of a section map, which allows for an easy overview of the collected data. For each section (i.e. each web page or local file), the specific words are computed, using the cumulative hypergeometric distribution of the words of the section compared to the complete corpus, i.e. an exact Fisher test is applied for each word in each section. As (Pedersen, 1996) remarks, the computation of the hypergeometric distribution is prohibitively heavy and, at the time, was only applicable to small corpora. He shows that the results (at least for the dependent bigrams he's working on) are more reliable than mutual information on sparse phenomena like collocations. The use of the hypergeometric distribution itself instead of some approximating function allows us to avoid the term "stop words" altogether: if a functional token (like a determiner or an auxiliary verb) appears in the results it is not due to its hight global frequency but rather because it really has a highly improbable frequency in the analyzed section if the words were distributed arbitrarily.



For our implementation, we developed a very effective computation of the cumulative hypergeometric distribution, with configurable thresholds. Contrary to the rest of Gromoteur, this function had to be written in C for performance reasons.

The specificity table shown on the right of the screenshot above shows the specificity of all the words in the selected sections compared to the rest of the corpus. Specificities are the

logarithmic order of the p-value of the exact Fisher test. The positive specificity of "8" for "guilty", for example, indicates that the surplus of the form "guilty" has a probability of approximately $10^8$. Inversely, negative specificities show an under-representation of a form in the selected parts. Moreover, the statistical window shows a graphical representation of the occurrence and specificity of forms over the whole corpus. All tables and graphics can easily be exported in order to be used in word processing software or for further analysis in a spreadsheet.



The window of collocation analysis by default is the whole text (webpage, import document...), but the window can be reduced to any kind of n-gram for the retrieval of collocations of a more syntactic nature. The same approach based on the hypergeometric distribution avoids again the usage of stopwords, contrarily to, for example, the computation of collocations based on mutual information. In the above screenshot the right most table shows the statistical collocates of the word "stupid", which contains "be" and "to" because they genuinely are collocates of "stupid" also in a linguistic sense (the function verb and the governed preposition in construction such as "Is is stupid to ..." or "He was so stupid to ..."). The collocation graph shows the strength of the collocative relations and allows for an easy exploration of the collocative relations by "clicking" through the graph.

## 3. General aspects, comparison, limitations

The Gromoteur is programmed in python and QT using some modules written in C. All data is stored in Sqlite databases. This allows for easy scaling on tens of thousands of webpage but the Gromoteur is essentially geared towards small size corpora consisting of a few hundred pages. The tool has already been used in various projects even before its recent completion including various newspaper analyses and the construction of a German-Chinese corpus from a bilingual website.

Most other crawlers are either command line tools built around the *wget* command that need hand-crafted scripting and tweaking or are web-based tools that have to be installed on a server, like for example the Babouk crawler (de Groc, 2011). To other knowledge, no other

desktop tool exists for crawling mid-size corpora that is geared toward linguists and sufficiently easy to install and to use. Webcorp Live (Kehoe & Renouf, 2002) is a Web-based tool that allows to use the web as a concordancer as the results taken from Google are presented in a more linguist-friendly manner. However, Webcorp is not a crawler as the whole web pages cannot be downloaded for further use. This also excludes statistical analysis of the results or further linguistic transformations like tagging or lemmatization. The Webcorp Linguist's Search Engine gives access to a set of preconfigured and preanalyzed English language webdata. In this sense, Gromoteur is a very different tool, as it collects the data locally and gives complete freedom to the user which data to collect and how to transform it.

Gromoteur is distributed under a GPL license in a Linux and a Windows version on http://gromoteur.ilpga.fr. The source code is available on https://launchpad.net/grosmoteur

However, the desktop architecture imposes certain limitations: The crawling cannot be distributed on different machines with various IP addresses, which makes it more vulnerable to server-side exclusions from downloading. Work in progress attempts to make the heuristic tools of Gromoteur accessible independently of the graphical interface and thus scriptable.

## References

Fleury S. and A. Salem (2002). *Lexico 3. Outil de statistique textuelle.* Paris: Université Sorbonne Nouvelle-Paris 3

Gerdes K. & Samvelian, P. (2008). A Statistical Approach to Persian Light Verb Constructions., *Proceedings of the Colloque Lexique et grammaire*, L'Aquila.

de Groc C. (2011). Babouk: Focused web crawling for corpus compilation and automatic terminology extraction. In *Proceedings of the IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology*. Lyon, France.

Kehoe A. & A. Renouf (2002). "WebCorp: Applying the Web to Linguistics and Linguistics to the Web". *WWW2002 Conference.*

Pedersen T. (1996). « Fishing for exactness », in *Proceedings of the South-Central SAS Users Group Conference*, Austin, TX. www.d.umn.edu/~tpederse/Pubs/scsug96.pdf

De Smedt T. & Daelemans W. (2012). "Pattern for Python". *Journal of Machine Learning Research*, 13: 2031–2035.