

## CARTOGRAPHIC DATA-MINING

*RUHSTORFER M., ECKRICH N.*

*Kartographie Huber, MüNCHEN, GERMANY*

This article deals with cartographic data mining – especially the challenge of finding cartographic material in the Internet. Each map production starts with a research for the best available base material. This is the linchpin for quality in terms of positional and temporal accuracy as well as the richness of the content.

In pre-internet times it was usually sufficient to contact the responsible Land Surveying Office and license their data. Another possibility was to acquire data from the local city building authority or - for the task of producing touristic maps – to cooperate with a guidebook author. Today – in the age of Internet and Web 2.0 - there are several other possibilities. Nevertheless, the challenge of finding the best fitting data remains.

Main problems for this task are: Data discovery, Multilingualism, Spatial reference, the variety of formats

### EXISTING SEARCH ENGINES AND THEIR USABILITY FOR CARTOGRAPHIC DATA MINING

The well-known search engines like Google, Bing or Yahoo! perform their indexing mainly based on different types of keyword-search. For images this approach would only work if all of them would be provided with a descriptive metadata set. This isn't the case, therefore keyword indexing can neither provide information about image content nor spatial extent. Also Spatial Web Services, which – pushed through the implementation of INSPIRE - are increasingly important for the cartography branch, can't be queried by ordinary search engines. Cartographic data mining will – and has to - be able to consider such services, and other “deep web” data sources as possible input.

### IMPROVED MULTILINGUALISM WITH A SPECIALIZED THESAURUS

The latter of the problems can be easily addressed by the use of a thesaurus specialized on cartographic terms. A thesaurus consists of a collection of terms and their translation into different languages.

### IMPROVED THEMATIC SEARCH CAPABILOTIES WITH FOCUSED CRAWLERS

Universal search engines have to deliver results for wide range of different topics. This has to be reflected in their ranking hits. Search engines which are optimized for a certain area of interest make use of so called Focused Crawlers. These software agents utilize content based analysis to rate the relevance of a website for a certain topic. Also the crawling is focused on URLs which are more likely relevant for the search queries without initial knowledge of the target website (based on URL- & Hyperlink-analysis).

Thereby Focused Crawlers are better suited to create high-quality and relevant document repositories for a certain domain compared to universal search engines, while saving computing power and disk usage.

The following search engines have been evaluated for their potential in finding cartographic data

	<i>Nutch</i>	<i>Sphinx</i>	<i>Bingo</i>	<i>Combine</i>
<i>Link-Topology</i>	<i>Opic-Implementation</i>	<i>Yes</i>	<i>Hits, Page-Order</i>	<i>---</i>
<i>URL-Order</i>	<i>URL-Filter with rules</i>	<i>URL-Filter with rules</i>	<i>Link following for positive classification</i>	<i>URL-Filter with rules</i>
<i>Additional Information</i>		<i>SQL</i>	<i>Taxonomy positive/negative</i>	<i>Thesaurus Integration</i>
<i>Classification</i>	<i>Lucene Engine</i>	<i>Classifier</i>	<i>SVM</i>	<i>Automatic Classifier</i>
<i>Protocols</i>	<i>http, ftp</i>	<i>http</i>	<i>http</i>	<i>http, ftp</i>
<i>Multilingualism</i>	<i>Yes, Language identifier</i>	<i>No, Language identifier</i>	<i>Yes</i>	<i>No</i>
<i>Open Source</i>	<i>Apache Software License</i>	<i>Yes</i>	<i>Yes</i>	<i>GPL</i>
<i>Platform</i>	<i>Java, and other portings</i>	<i>PHP, Perl, Python, Ruby, Java etc.</i>	<i>Java VM</i>	<i>Linux</i>
<i>Code</i>	<i>Open Source</i>	<i>Open Source</i>	<i>Open Source</i>	<i>Open Source</i>
<i>Development Activity</i>	<i>High</i>	<i>High</i>	<i>Low</i>	<i>Low</i>
<i>Documentation</i>	<i>Good</i>	<i>Good</i>	<i>Good</i>	<i>Good</i>
<i>File types</i>	<i>HTML, TXT, PDF, and others</i>	<i>HTML, TXT,</i>	<i>HTML, TXT, DOC, PDF</i>	<i>HTML, TXT, PDF, DOC, PS, TeX</i>
<i>Homepage</i>	<i>http://nutch.apache.org/</i>	<i>http://www.sphinxsearch.com/</i>	<i>http://www.mpiinf.mpg.de/departments/d5/software/bingo/index.html</i>	<i>http://combine.it.lth.se</i>

The evaluation of the Focused Crawlers led to the decision for Kartographie Huber to rely on Nutch. Reasons are the active developer community, which hopefully will result in constant improvements, as well as the good documentation.

### IMPROVED RESULTS WITH CONTENT-BASED IMAGE RETRIEVAL

The last issue which has to be addressed is the relatively limited image recognition capabilities of existing search engines.

- Text based methods: Text based retrieval methods determine the graphical primitives of images. These primitives are the distribution of colors within an image, the color histogram, the edge histogram, the dominant color, the image size, the color spectrum and the file format.
- Image based methods: So called QBVE-systems are based on an ancient approach of IBM. They make use of the fact that an example image will greatly improve the possibility to find similar images. This explains the acronym “query by visual example”. It’s distinguished between systems which use example drawings of the user and systems which provide given example images.
- Combined Methods: There exist approaches which combine the afore mentioned methods. This eliminates many of their shortcomings. A functionality which isn’t provided till now – but should be implemented within the development of a cartographic data mining search engine – is the possibility to create a spatial reference of the image extent.

Within the research for cartographic data mining the systems shown in table 3 have been evaluated .

	<i>Caliph&amp;Emir</i>	<i>VizIR</i>	<i>IMGseek</i>	<i>WISE</i>	<i>Rummager</i>	<i>BrisC</i>
<i>Open Source Software</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Yes</i>
<i>Multilingualism</i>	<i>No</i>	<i>No</i>	<i>No</i>	<i>No</i>	<i>No</i>	<i>No</i>
<i>Documentation</i>	<i>Good</i>	<i>None</i>	<i>Available</i>	<i>None</i>	<i>Available</i>	<i>Minimalistic</i>
<i>Colors</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>
<i>Structures</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>
<i>Textures</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>
<i>Metadata retrieval</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>
<i>Actuality</i>	<i>11-2010</i>	<i>6-2003</i>	<i>1-2009</i>	<i>10-2010</i>	<i>05-2010</i>	<i>Not anymore</i>
<i>Platform</i>	<i>Java</i>	<i>Java</i>	<i>C++, Python</i>	<i>-</i>	<i>Java</i>	<i>C#, .net 2.0</i>

The evaluation showed that the Caliph and Emir software show promising results for their usability for our purpose.

### RESULTS

The technologies described in this article will only improve the possibilities for cartographic data mining to a limited amount. It’s the combination of these technologies which we think make our approach interesting. A focused crawler will limit the results which have to be checked by the CBIR-System to a reasonable amount. Thereby the data repository will be “double-checked” to contain only the most relevant results. And the thesaurus for cartographic terms will enable the editor to retrieve information not only in languages he can master but also in additional languages relevant for the map production.

### CONCLUSION AND FUTURE PLANS

Cartographic data mining is necessary to keep maps up-to-date and competitive to navigational data and other data providers. The classical maps have to be most precise to fit the markets needs and therefore cartographers need the best information available. To enable cartographers and editors to make best use of the internet they have to be provided with the right tools. In case of the Internet in our opinion this should be a search engine which combines focused crawling with Content-Based Image Retrieval functionality as well as multilingualism. Such a tool will help to restrict the information flood to a manageable amount. The problem isn’t that the information isn’t available. The problem is to know how to find it.