

Search Result Visualisation with xFIND

Keith Andrews, Christian Gütl, Josef Moser, Vedran Sabol, Wilfried Lackner

Institute for Information Processing and Computer Supported New Media (IICM),
Graz University of Technology
A-8010 Graz, Austria

Abstract

The xFIND gatherer broker architecture provides a wealth of metadata, which can be used to provide sophisticated search functionality. Local or remote documents are indexed and summaries and metadata stored on an xFIND broker (server). An xFIND client can search a particular broker and access rich metadata for search result presentation, without having to fetch the original documents themselves. Search result sets are not only presented as a traditional ranked list, but also in an interactive scatterplot (Search Result Explorer) and using dynamic thematic clustering (VisIslands).

1 Introduction and Status of Research

The amount of information readily accessible to individuals at their workplace and at home is rapidly increasing. There are now more than one billion unique, indexable web pages [Inktomi and NEC Research Institute, 2000]. Monolithic, centralised search engines are increasingly unable to cope with the exponential growth of the web. Since individual pages are visited perhaps only once or twice a year, the centralised index is inherently out of date. More general search queries often return many hundreds or thousands of matching documents. Hence the motivation for both a scalable resource discovery framework and for visualisation tools to help end users explore search result sets.

WAIS [Kahle et al., 1992] was a client-server indexing and retrieval system, using a protocol derived from Z39.50 [ANSI/NISO, 1995]. The Harvest system [Bowman et al., 1995] is a distributed, gatherer-broker information discovery and access system for the web, an integrated set of tools, written largely in C, for gathering information from diverse repositories, building topic-specific content indices, and searching the indices.

Envision [Nowell et al., 1996, 1997] plots the result of bibliographic searches in a two-dimensional scatterplot. The mapping of particular attributes to visual representations such as the x-axis, y-axis, icon size, and icon shape, is controlled by drop-down menus. Bead [Chalmers, 1993, 1996b,a] uses a force-directed placement technique to lay out relationships between documents in a corpus as a landscape. More similar documents lie closer together in the landscape. Searches can be made and the results highlighted. ThemeScape (part of the SPIRE text visualisation suite) [Wise, 1999] examines a corpus of text documents and extracts a set of discriminating terms (words), usually nouns, which characterise topics in the corpus. A combination of stoplists, synonym substitution, and statistical analysis is used to select 200 to 300 “good” discriminating terms. A landscape of topical terms is generated and documents contribute to the height field at each point, based on their contribution to each topic. Self-organising maps (SOMs) [Kohonen, 2000; WEBSOM, 2000] use neural networks to organise a set of text documents. The neural network is trained

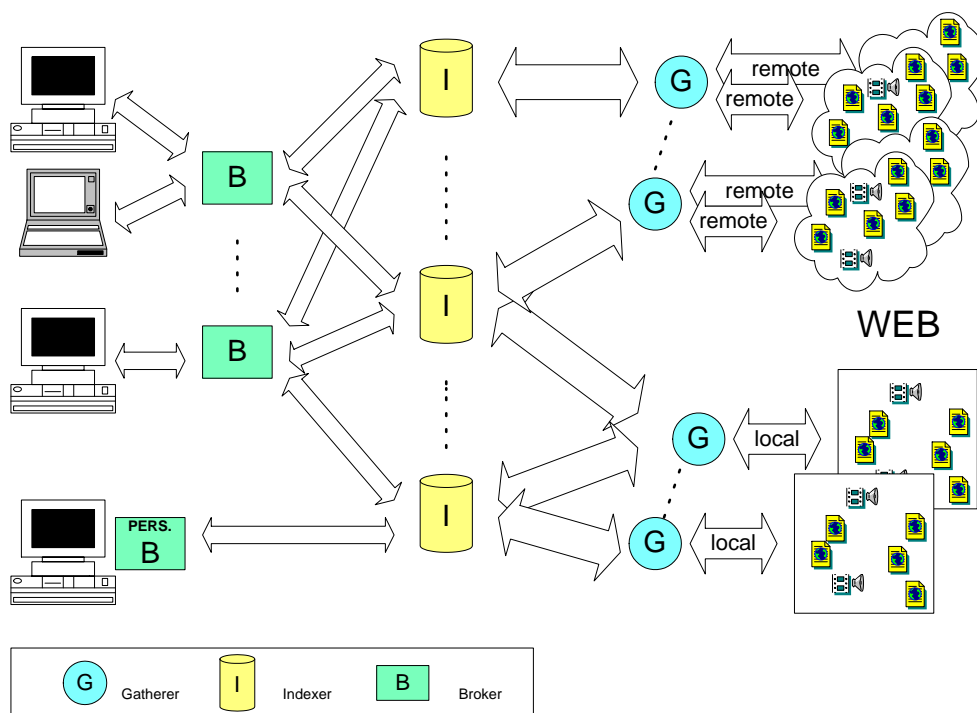


Figure 1: The distributed architecture of the xFIND system.

initially on a sample set of documents and thereafter can assign further documents to their corresponding location on a regular (hexagonal) grid in unsupervised mode.

2 The xFIND System

xFIND[Gütl, 2000; Gütl et al., 1998], the extended Framework for Information Discovery, is a framework for distributed information discovery and knowledge management. For scalability, xFIND uses a gatherer-indexer-broker architecture, similar to that pioneered by the Harvest[Bowman et al., 1995] system. As shown in Figure 1, xFIND consists of *gatherers*, *indexers* and *brokers*. For reasons of portability and platform independence, xFIND is implemented entirely in Java.

In addition to traditional metadata like title, keywords, and description, xFIND also gathers and indexes document headings, and information about embedded links and images. Metadata regarding aspects of information quality, such as authority, diction, and target audience can be manually or semi-manually added.

2.1 The xFIND Gatherer

The gatherer gathers information about documents and resources, both local and remote, and pre-processes this information. The gathering process is adaptive. At configurable intervals, entire servers, particular sub-sites, or individual documents can be gathered. To reduce network load, a local gatherer can be configured to gather information from a locally mounted file system. The gatherer pre-processes HTML and plain text documents (soon also PDF) and generates a pre-defined set of metadata (a document descriptor or summary) for each document. Thumbnails of embedded images are also generated and stored in the document descriptor. Only the document descriptors are passed on to one or more xFIND indexers.

2.2 The xFIND Indexer

The document descriptors harvested by a gatherer can be fetched by one or more authorised indexers. The indexer indexes a set of document descriptors and renders them searchable. An indexer may specialise, for example, in a particular topic or geographical location. Statistical information such as term frequency and discriminating terms (discriminators) are generated. Furthermore, trusted external systems (rating systems, ACF systems, etc.) are allowed to contribute additional metadata fields. The xFIND architecture also provides for the contents of a particular index to be replicated in whole or in part, in order to minimise network loads.

2.3 The xFIND Broker

An xFIND broker is the starting point for user interactions. A broker can distribute its search queries to a particular set of indexers. The broker is also able to expand queries using a thesaurus. The results of a distributed search are collated and compiled into a uniform search result set. Brokers can be individually tailored for a division, a department, a group of employees or even for a single user as well as supporting particular topics.

The xFIND broker provides standard search functionality (simple, extended, and expert search) through HTML forms. Search queries can combine both full-text and descriptive and evaluative metadata. The standard result set is a linear list ranked by relevance.

2.4 Rich Metadata

The enhanced metadata set can be divided into two main parts: extracted document information and quality metadata. Automatically extracted metadata are generated by the gatherer, which identifies the URL, mime type, file size, and creation or modification time of each document. Further information extraction depends on the type of the object, as well as the proper filter for processing the object. At the time of writing, plain text and HTML filter are available for textual information, and an image filter (supporting gif and jpg file format) for multimedia documents. Support for PDF and audio and video formats is planned. For HTML documents, common meta-attributes such as title, keywords, description, and language are parsed and processed by the gatherer. The full-text content is retained for indexing. Headings, links, and (thumbnails of) embedded images are extracted and retained as metadata.

The gatherer also creates an electronic fingerprint of each information object. This fingerprint suffices to determine the trustworthiness of information in case of replication and allows detecting the origin of every piece of information.

The xFIND system supports the integration of external metadata, for example for non-textual objects. Since authors rarely enrich individual documents with additional metadata, xFIND supports the definition of metadata for an entire document structure, a directory, or a particular document by inserting additional meta data files. More specific metadata overrides more general metadata.

3 Search Result Visualisation

The richness of metadata provided by xFIND can greatly aid users during the search process.

For the visualisation examples presented here, an xFIND broker specialising in the topic of Knowledge Management will be used. At the time of writing, the broker has access to an index of some 44,878 documents, gathered from the sites shown in Table 1, as well as a number of other sites less frequently. The query used in each case is the single word “agents”, for which there are 314 matching documents.

Daily http://agents.www.media.mit.edu/groups/agents http://bots.internet.com http://agents.umbc.edu
Once a Week http://www.uibk.ac.at/sci-org/voeb http://www.ai.mit.edu http://xfind.iicm.edu

Table 1: These sites are indexed for the broker on Knowledge Management, as well as a number of other sites at less frequent intervals.

3.1 Ranked List Search Results

The default presentation of search results by an xFIND broker takes the form of a traditional ranked list. Figure 2 shows the first two matching documents. The most relevant document to the query “agents” is entitled “Coordination as Distributed Search”. Note the context of the query term “agents” is shown for each document, the most important terms (discriminators) contained in each document are listed, and thumbnail images are shown for documents containing embedded images. Figure 3 shows 13 thumbnails for the ninth document in the result list.

3.2 Interactive Scatterplots with the Search Result Explorer

The Search Result Explorer uses a scatterplot (starfield display) to allow interactive exploration of the search result set based on the rich metadata associated with each object, in a manner similar to Envision[Nowell et al., 1996]. Documents are plotted according to two of their metadata attributes (corresponding to the x and y axes). Further metadata attributes can be mapped to icon size and icon colour, allowing four dimensions of metadata to be visualised and explored simultaneously. If too many documents would be mapped to the same proximity, a group icon is used to represent that subset of documents. For group icons, the size and colour of the group icon is determined (under user control) by the maximum, minimum, median, or average value of the group’s members. Since it is possible to zoom in on specific areas of the display, an overview window is provided in the lower left corner to help maintain context and orientation.

Figure 4 shows the first 210 (a user-configurable limit) of the 314 matching documents plotted by relevance on the y axis and document size on the x axis. The most relevant document is shown at the top of the plot. At the moment, the colour of each document icon is determined by the document’s age, from yellow older documents to white recent documents. Relevance is mapped to icon size, providing a redundant encoding. More relevant documents are both larger and towards the top of the plot.

In Figure 5, this document has been selected and its metadata displayed. It is the same document from 10th October 1998 entitled “Coordination as Distributed Search”.

Figure 6 illustrates the interactive nature of the plot. Most of the rich metadata attributes provided by xFIND can be mapped to either axis or to icon size or colour. Figure 7 shows the result of the change. The y axis now corresponds to the modification date of the document, and document relevance is mapped to both icon size and icon colour (more relevant are orange, less relevant are white). It can be seen at a glance, that the most relevant documents are about a year old and reasonably small.

3.3 Dynamic Clustering with VisIslands

The VisIslands interface supports dynamic thematic clustering of search result sets, in a manner similar to SPIRE’s Themescape[Wise, 1999] and its commercial successor Cartia[Cartia, 2000].

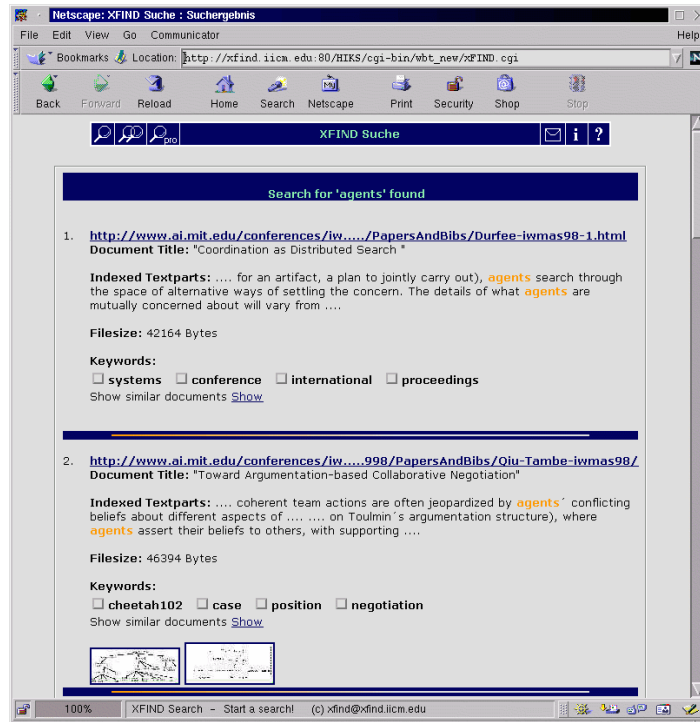


Figure 2: The standard ranked list returned by an xFIND broker. Note the context of the query term “agents”, the most important terms (discriminators) contained in each document, and the thumbnail images.

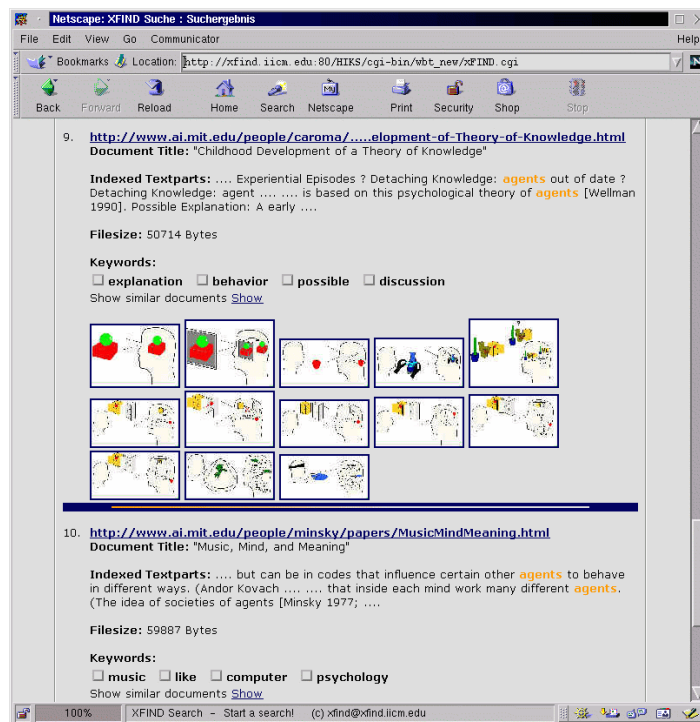


Figure 3: Image thumbnails of embedded images are generated and associated as metadata with a HTML document object.

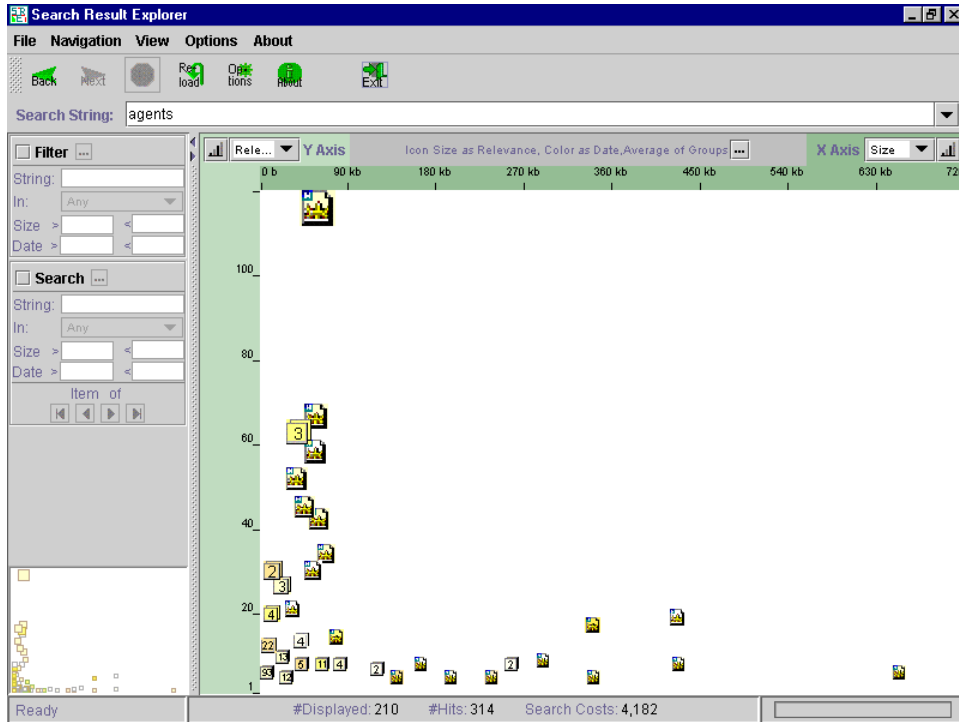


Figure 4: The Search Result Explorer plots search results along two axes. Here, document relevance is mapped to the y axis and document size to the x axis. More relevant documents also have larger icons. Older documents are yellowish.

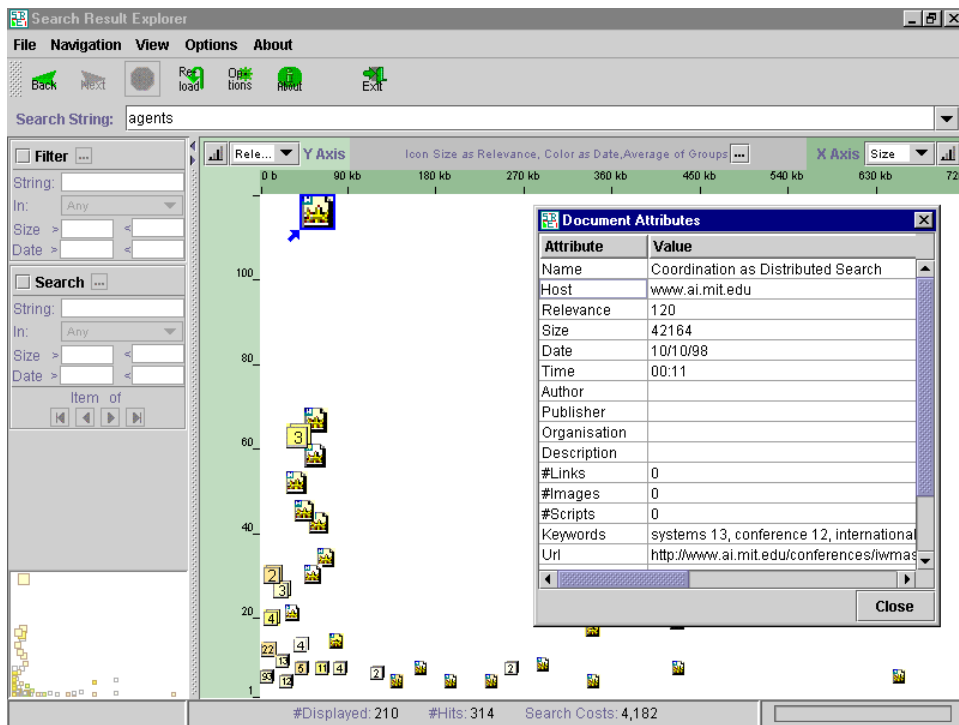


Figure 5: Search Result Explorer: The top matching document has been selected and its metadata displayed.

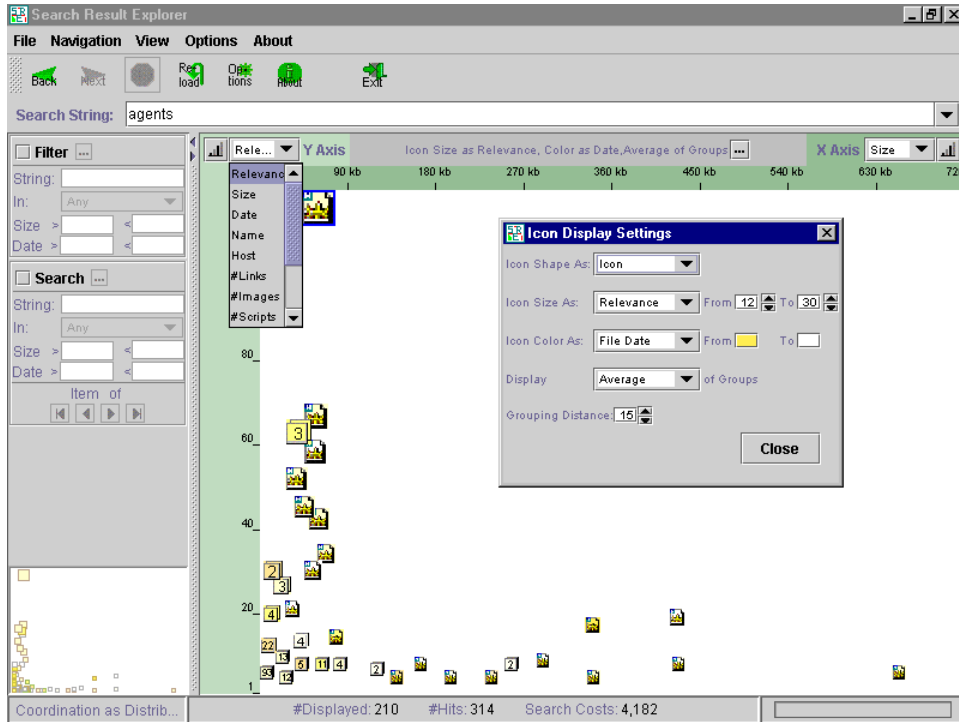


Figure 6: Search Result Explorer: Most of the rich metadata attributes provided by xFIND can be mapped to either axis, or to icon size or icon colour.

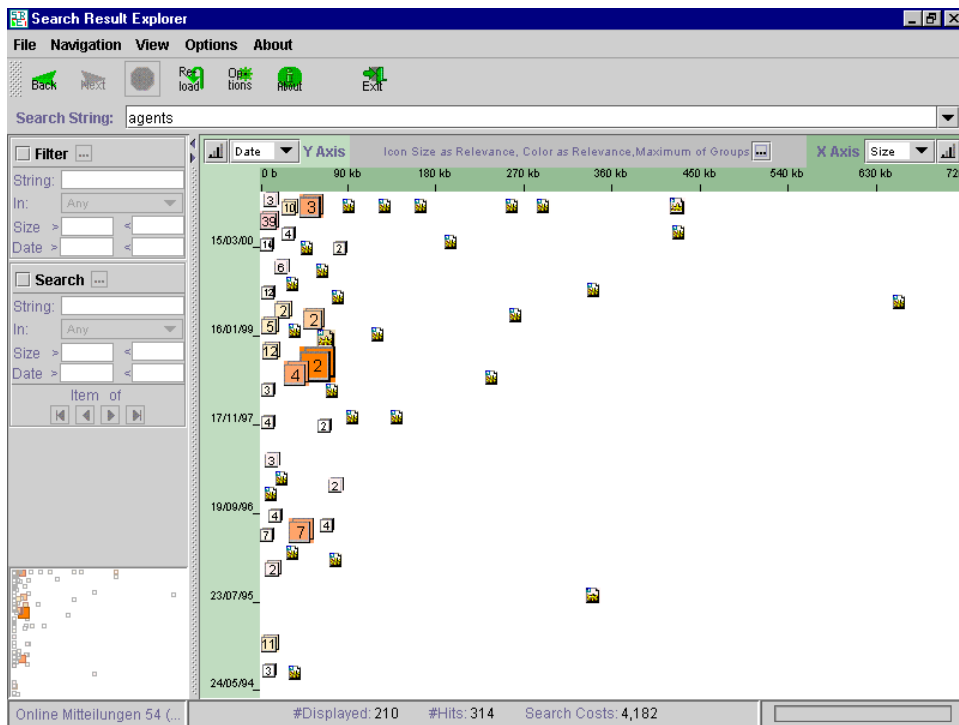


Figure 7: Search Result Explorer: An alternative view. Modification date has been mapped to the y axis. More relevant documents are now both larger and more orange.

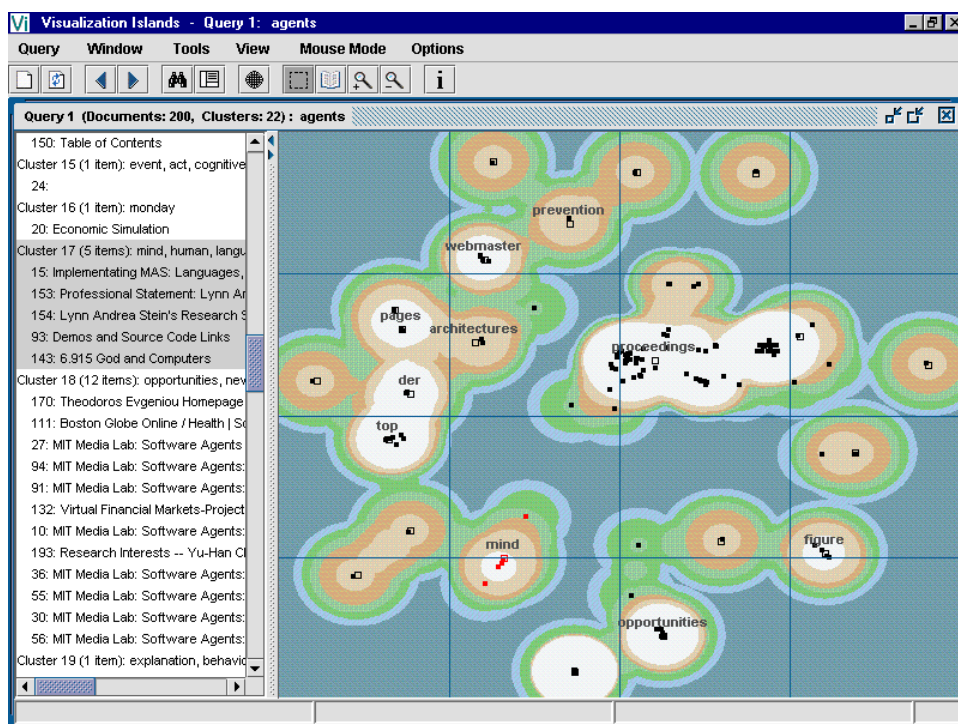


Figure 8: VisIslands: Pre-clustering using hierarchical agglomerative clustering has identified the 22 clusters shown in the left-hand pane. Cluster 17 has been selected.

The search result set is first pre-clustered using hierarchical agglomerative clustering (or optionally k-means clustering)[Anil K. Jain and Flynn, 1999]. The cluster centroids are then distributed randomly in the viewing rectangle. The documents belonging to each cluster, as determined by the initial pre-clustering, are then placed in a ring around each centroid. This arrangement is fine-tuned using a linear iteration force-directed placement algorithm derived from Chalmers [1996b]. Documents similar to one another are attracted towards each other. After a certain cut-off point, the arrangement has stabilised, and each document contributes its weight to the height field of the grids within which it lies. Dense areas of many documents have corresponding peaks. The overall result is like a contour map of islands. A more three-dimensional visualisation of the islands would also be possible, but has not yet been implemented.

Figure 8 shows the islands visualisation for the first 200 documents matching our example query. Pre-clustering has identified 22 clusters. Cluster 17 concerns “mind, human, language” and has been selected. Note the corresponding visual cluster of red documents in the islands display. Figure 9 shows the metadata associated with Cluster 17.

Figure 10 focuses on Cluster 22, containing 108 documents. Note that fine tuning with force-directed placement has attracted one document which pre-clustering assigned to Cluster 22 over towards the “webmaster” and “architectures” clusters. Zooming in on Cluster 22, Figure 11 shows that, in fact, many of the documents assigned to Cluster 22 on pre-clustering, should perhaps have been assigned to the cluster called “erl”auterungen”.

4 Future Work

The xFIND project is ongoing. Current work includes development of a utility to automatically detect the language of a document, so that stop list and stemming algorithms can be applied automatically. Also, a filter for indexing PDF documents is planned.

On the visualisation side, the Search Result Explorer and VisIslands are currently independent xFIND

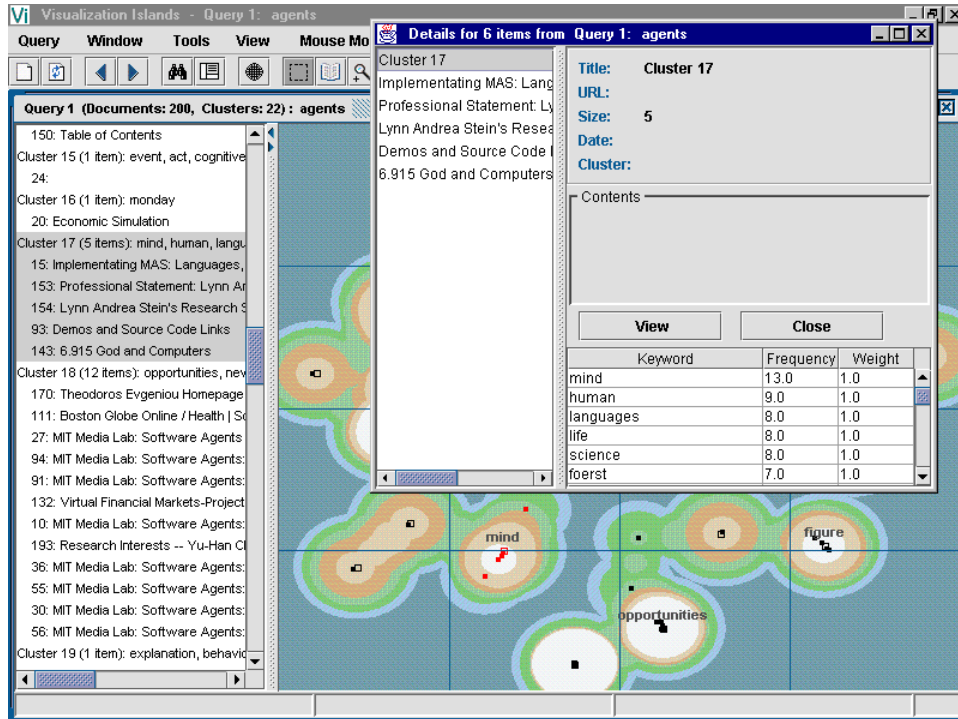


Figure 9: VisIslands: The metadata associated with Cluster 17 is displayed. Its most frequent terms include “mind”, “human”, and “language”.

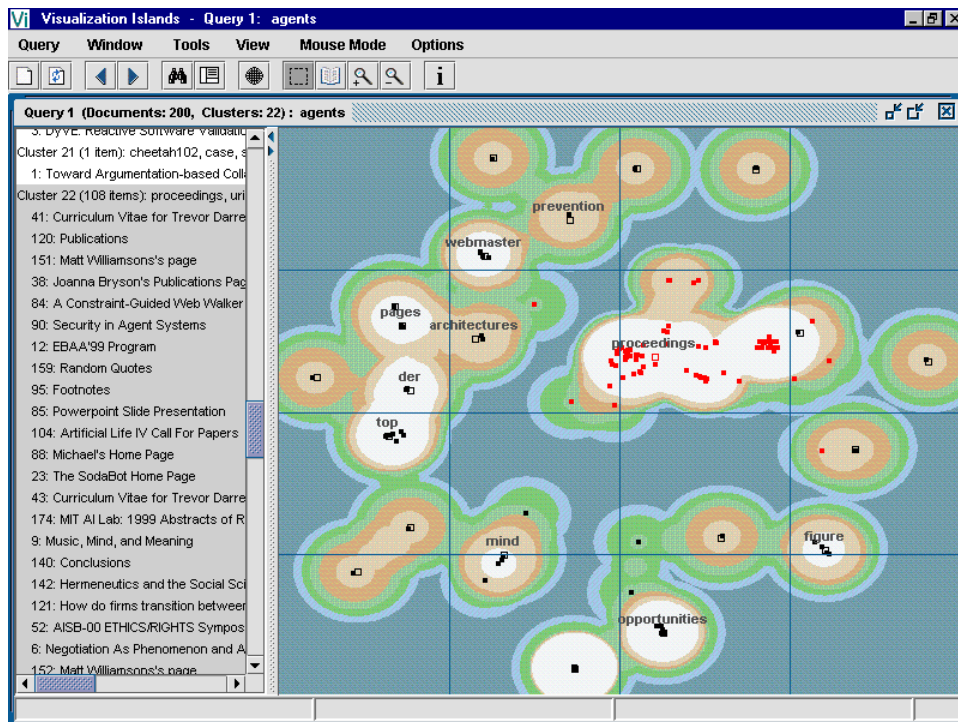


Figure 10: VisIslands: Cluster 22 deals with a variety of topics including “proceedings” and “conference”.

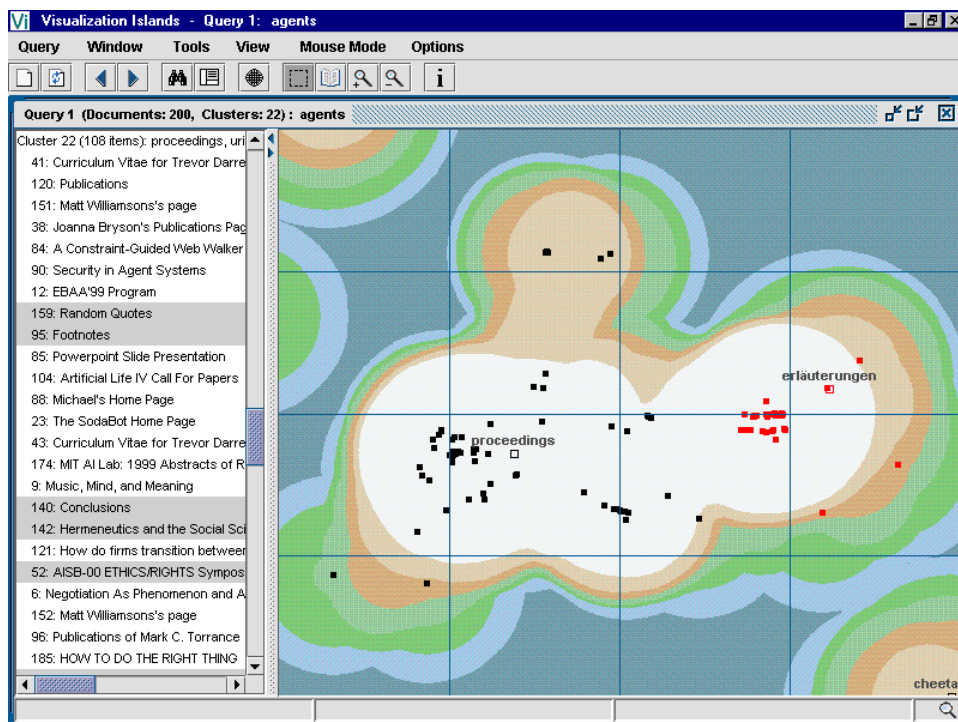


Figure 11: VisIslands: After zooming in on Cluster 22, the group of documents on the right-hand peak has been manually selected.

clients. A Master's thesis just commencing to build an integrated search client incorporating query formulation and history, ranked result lists, and a framework for plugging in one or more synchronised visualisations, embracing both the Search Result Explorer and VisIslands.

5 Concluding Remarks

The xFIND infrastructure builds on the original work of Harvest and provides a rich, highly configurable knowledge management framework. The Search Result Explorer and VisIslands visualisations demonstrate the potential of information visualisation techniques applied to the exploration of search result sets. It is intended to accompany this paper with a live demonstration of the system.

6 Acknowledgements

We would like to acknowledge the support of the IICM, Graz University of Technology, partial funding from the Austrian Ministry of Science, and the contributions of former colleagues and students: Jürgen Heber, Axel Jurak, Bernhard Knögler, Herbert Legenstein, Susanne Mayr, and Erwin Weitlaner.

References

- Anil K. Jain, M. N. M. and Flynn, P. J. (1999). *Data Clustering: A Review*. ACM Computing Surveys, 31(3):264–323. <http://www.acm.org/pubs/citations/journals/surveys/1999-31-3/p264-jain/>.
- ANSI/NISO (1995). *Z39.50-1995*, Library of Congress. <http://lcweb.loc.gov/z3950/agency/document.html>.

- Bowman, C. M., Danzig, P. B., Hardy, D. R., Manber, U., Schwartz, M. F. and Wessels, D. P. (1995). Harvest: A Scalable, Customizable Discovery and Access System. Technical Report CU-CS-732-94, University of Colorado. <ftp://ftp.cs.colorado.edu/pub/cs/techreports/schwartz/Harvest.ps>.
- Cartia (2000). Mapping the Information Landscape. <http://www.cartia.com/>.
- Chalmers, M. (1993). *Using a Landscape Metaphor to Represent a Corpus of Documents*. In Spatial Information Theory, Proc. COSIT'93, pages 377–390, Boston, Massachusetts (1993). Springer LNCS 716. <http://www.dcs.gla.ac.uk/~matthew/papers/ecsit93.pdf>.
- Chalmers, M. (1996a). *Adding Imageability Features to Information Displays*. In Proc. UIST'96, Seattle, Washington (1996a). ACM. <http://www.dcs.gla.ac.uk/~matthew/papers/uist96.pdf>.
- Chalmers, M. (1996b). *A Linear Iteration Time Layout Algorithm for Visualising High-Dimensional Data*. In Proc. Visualization'96, pages 127–132, San Francisco, California (1996b). IEEE Computer Society. <http://www.dcs.gla.ac.uk/~matthew/papers/vis96.pdf>.
- Gütl, C. (2000). xFIND: Extended Framework for Information Discovery. IICM, Graz University of Technology. <http://xfind.iicm.edu/>.
- Gütl, C., Andrews, K. and Maurer, H. (1998). *Future Information Harvesting and Processing on the Web*. In Proc. European Telematics: Advancing the Information Society, Barcelona, Spain (1998). http://www2.iicm.edu/~cguetl/papers/fihap/fihap_en.html.
- Inktomi and NEC Research Institute (2000). Web Surpasses One Billion Documents. Press Release. <http://www.inktomi.com/new/press/billion.html>.
- Kahle, B., Morris, H., Davis, F., Tiene, K., Hart, C. and Palmer, R. (1992). *Wide Area Information Servers: An Executive Information System for Unstructured Files*. Electronic Networking: Research, Applications and Policy, 2(1):59–68.
- Kohonen, T. (2000). *Self-Organizing Maps*. Springer, third edition. ISBN 3540679219.
- Nowell, L. T., France, R. K. and Hix, D. (1997). *Exploring Search Results with Envision*. In CHI'97 Demonstration (Extended Abstracts), pages 14–15, Atlanta, Georgia (1997). ACM. <http://www.acm.org/sigchi/chi97/proceedings/demo/1tn1.htm>.
- Nowell, L. T., France, R. K., Hix, D., Heath, L. S. and Fox, E. A. (1996). *Visualizing Search Results: Some Alternatives to Query-Document Similarity*. In Proc. SIGIR'96, pages 67–75, Zurich, Switzerland (1996). ACM.
- WEBSOM (2000). WEBSOM - Self-Organizing Maps for Internet Exploration. Helsinki University of Technology. <http://websom.hut.fi/websom/>.
- Wise, J. A. (1999). *The Ecological Approach to Text Visualization*. Journal of the American Society for Information Science, 50(9):814–835. http://www.vistg.net/hat/Wise_draft/Ch5.Wise.html.