

World Wide Web Search Technologies

A chapter submitted to the book:
Architectural Issues of Web-Enabled Electronic Business
edited by Shi Nansi for Idea Group Publishing

ABSTRACT

With over 800 million pages covering most areas of human endeavor, the World Wide Web is fertile ground for information retrieval. Numerous search technologies have been applied to Web searches, and the dominant search method has yet to be identified. This chapter provides an overview of existing Web search technologies and classifies them into six categories: (i) hyperlink exploration, (ii) information retrieval, (iii) metasearches, (iv) SQL approaches, (v) content-based multimedia searches, and (vi) others. A comparative study of some major commercial and experimental search services is presented, and some future research directions for Web searches are suggested.

Keywords: Survey, World Wide Web, Searches, Search Engines, and Information Retrieval.

1. INTRODUCTION

Searching for Web pages is one of the most common tasks performed on the Web. It is also one of the most frustrating. The situation is getting worse because of the Web's fast growing size and lack of structure style, as well as the inadequacy of existing Web search technologies [34]. Traditional search techniques are based on users typing in search keywords which the search services can then use to scan Web pages. However, this approach normally retrieves too many documents, of which only a small fraction are relevant to the users' need. Furthermore, the most relevant documents do not necessarily appear at the top of the query output order. A number of corporations and research organizations are taking a variety of approaches to try to solve these problems. These approaches are usually diverse and none of them dominate the field. This chapter provides a survey and classification of the available World Wide Web search techniques, with an emphasis on non-traditional approaches. Related Web search technology reviews can also be found in [25, 33, 35, 38].

Requirements of Web Searches

It is first necessary to examine what kind of features a Web search system is expected to have in order to conduct effective and efficient Web searches and what kind of challenges may be faced in the process of developing new Web search techniques. The requirements for a Web search system are listed below in order of importance:

1. Effective and efficient location and ranking of Web documents.
2. Thorough Web coverage.
3. Up-to-date Web information.
4. Unbiased access to Web pages.
5. An easy-to-use user interface which also allows users to compose any reasonable query.
6. Expressive and useful search results.
7. A system that adapts well to user queries.

Web Search Technologies

Numerous Web search technologies have been proposed and each technology employs a very different approach. This survey classifies the technologies into six categories: (i) hyperlink exploration, (ii) information retrieval, (iii) metasearches, (iv) SQL approaches, (v) content-based multimedia searches, and (vi) others. The chapter is organized as follows: Section 2 introduces the search engine structure and Sections 3 to 8 give an overview of each of the six Web search technologies in turn. A comparative study of major commercial and experimental search services is shown in Section 9 and the final section gives a summary and suggests future research directions.

2. SEARCH ENGINE STRUCTURE

Two different approaches are applied to Web search services: (i) genuine search engines and (ii) directories. The difference lies in how listings are compiled.

- Search engines, such as HotBot, create their listings automatically.
- A directory, such as Yahoo!, depends on humans for its listings.

Some search engines, known as hybrid search engines, maintain an associated directory. Search engines traditionally consist of three components: the crawler, the indexing software, and the search and ranking software [24, 56]. Figure 1 shows the system structure of a search engine.

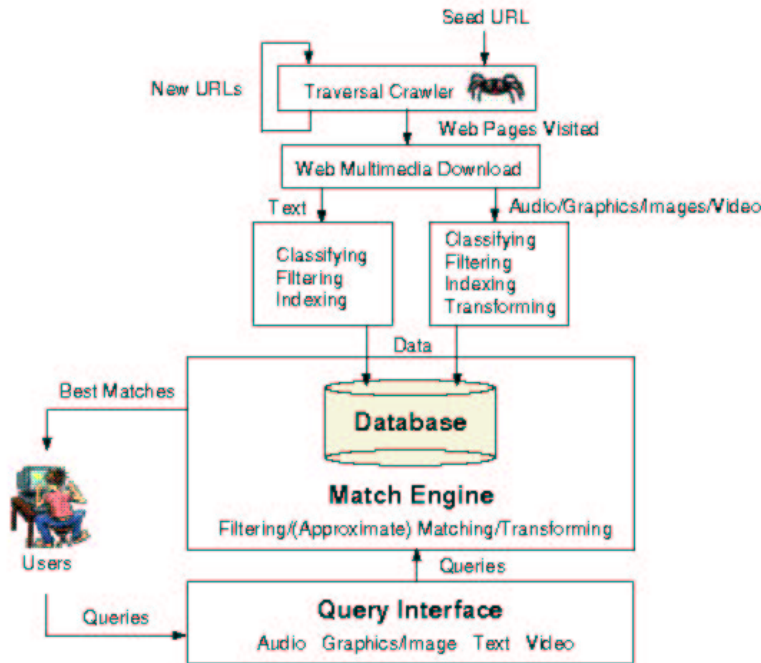


Figure 1: System structure of a search engine.

Crawler

A crawler is a program that automatically scans various Web sites and collects Web documents from them. Crawlers follow the links on a site to find other relevant pages. Two search algorithms, breadth-first searches and depth-first searches, are widely used by crawlers to traverse the Web. The crawler views the Web as a graph, with the nodes being the objects located at Uniform Resource Locators (URLs). The objects could be HTTPs (Hypertext Transfer Protocols), FTPs (File Transfer Protocols), mailto (e-mail), news, telnet, etc. They also return to sites periodically to look for changes. To speed up the collection of Web documents, several crawlers are usually sent out to traverse the Web at the same time. Three basic tools are usually used to implement an experimental crawler:

- *lynx*: Lynx is a text browser for Unix systems. For example, the command “`lynx -dump -source http://www.w3c.com/`” downloads the Web page source code at `http://www.w3c.com/`.
- *java.net*: The `java.net` package provides plenty of networking utilities. Two classes in the package, `java.net.URL` and `java.net.URLConnection`, can be used to download Web pages.
- *CPAN (Comprehensive Perl Archive Network)*: Perl has been used intensively for Web related applications. Some scripts provided by CPAN at `http://www.cpan.org/` are useful for crawler

construction.

To construct an efficient and practical crawler, some other networking tools have to be used.

Indexing Software

Automatic indexing is the process of algorithmically examining information items to build a data structure that can be quickly searched. Filtering [4] is one of the most important pre-processes for indexing. Filtering is a typical transformation in information retrieval, for example to reduce the size of a document, and/or standardize it to simplify searching. Traditional search engines utilize the following information, provided by HTML files, to locate the desired Web pages:

- *Content*: Page content provides the most accurate and full-text information. However, it is also the least-used type of information since context extraction is still far less practical.
- *Descriptions*: Page descriptions can either be constructed from the metatags or submitted by Webmasters or reviewers.
- *Hyperlinks*: Hyperlinks contain high-quality semantic clues to a page's topic. A hyperlink to a page represents an implicit endorsement of the page being pointed to. [10]
- *Hyperlink text*: Hyperlink text is normally a title or brief summary of the target page.
- *Keywords*: Keywords can be extracted from full-text documents or metatags.
- *Page title*: The title tag, which is only valid in a head section, defines the title of an HTML document.
- *Text with a different font*: Emphasized text is usually given a different font to highlight its importance.
- *The first sentence*: The first sentence of a document is also likely to give crucial information related to the document.

Search and Ranking Software

Query processing is the activity of analyzing a query and comparing it to indexes to find relevant items. A user enters a keyword (or keywords along with Boolean modifiers, such as “and,” “or,” or

“not”) into a search engine, which then scans indexed Web pages for the keywords. To determine in which order to display pages to the user, the engine uses an algorithm to rank sites that contain the keywords [58]. For example, the engine may count the number of times the keyword appears on a page. To save time and space, the engine may only look for keywords in metatags. A metatag is an HTML tag that provides information about a Web page. Unlike most HTML tags, metatags do not affect a document’s appearance. Instead, they include such information as a Web page’s contents and some relevant keywords. The following six sections give various methods of indexing, searching, and ranking.

3. HYPERLINK EXPLORATION

Hypermedia documents contain cross-references to other related documents, and these “links” function as hyperlinks, allowing the user to move easily from one to the other. Links can be tremendously important sources of information for indexers; the creation of a hyperlink by the author of a Web page represents an implicit endorsement of the page being pointed to. This approach is based on identifying two important types of Web pages for a given topic:

- *Authorities*, which provide the best source of information on the topic; and
- *Hubs*, which provide collections of links to authorities.

For the example of professional basketball information, the official National Basketball Association site <http://www.nba.com/> is considered to be an authority, while the ESPN site <http://www.espn.com/> is a hub. Authorities and hubs are either given top ranking in the search results or can be used to find related Web pages [14].

Analyzing the interconnections of a series of related pages can identify the authorities and hubs for a particular topic. A simple method to update a nonnegative authority weight x_p and a nonnegative hub weight y_p is given in [10]. If a page is pointed to by many good hubs, its authority weight is updated by the following formula:

$$x_p = \sum_{q \text{ such that } q \rightarrow p} y_q, \tag{1}$$

where the notation $q \rightarrow p$ indicates that q links to p . Similarly, if a page points to many good authorities, its hub weight is updated via

$$y_p = \sum_{q \text{ such that } p \rightarrow q} x_q. \quad (2)$$

Unfortunately, applying the above formulas to the whole Web to find authorities and hubs is impracticable. Ideally, the formulas are applied to a small collection S_σ of pages which contain plenty of relevant documents. The concept of a root set and a base set has been proposed by [30] to find S_σ . The root set is constructed by collecting the t highest-ranked pages for the query σ from a text-based search engine such as Google or Yahoo!. However, the root set may not contain most of the strongest authorities. A base set is therefore built by including any page pointed to by a page in the root set and any page that points to a page in the root set. Figure 2 shows an example of a root set and a base set. The above formulas can then be applied to a much smaller set, the base

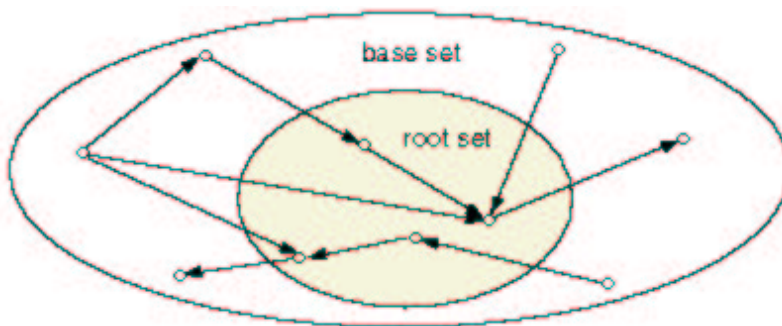


Figure 2: Expanding the root set into a base set.

set, instead of the whole Web.

In addition to the methods finding authorities and hubs, a number of search methods based on connectivity have been proposed. A comparative study of various hypertext link analysis algorithms is given in [7]. The most widely used method is a Page Rank model [8], which suggests the reputation of a page on a topic is proportional to the sum of the reputation weights of pages pointing to it on the same topic. That is, links emanating from pages with high reputations are weighted more heavily. The concepts of authorities and hubs, together with the Page Rank model, can also be used to compute the reputation rank of a page; and those topics for which the page has a good reputation are then identified [45]. Some other ad hoc methods include an HVV (Hyperlink Vector Voting) method [37] and a system WebQuery [9]. The former method uses the content of hyperlinks to a document to rank its relevance to the query terms, while the latter system studies

the structural relationships among the nodes returned in a content-based query and gives the highest ranking to the most highly connected nodes. An improved algorithm obtained by augmenting [30] with content analysis is introduced in [6].

4. INFORMATION RETRIEVAL (IR)

IR techniques are widely used in Web document searches [25]. Among them, relevance feedback and data clustering are two of the most popular techniques used by search engines.

Relevance Feedback

An initial query is usually a wild guess. Retrieved query results are then used to help construct a more precise query [12]. For example, a query is submitted to a search engine:

Which TOYOTA dealer in Atlanta has the lowest price for a Corolla 2001?

The engine may produce the following list of results:

1. Get the BEST price on a new Toyota, Lexus car or truck. <http://www.toyotaforless.com/>
2. Toyota of Glendale—Your #1 Toyota dealer. <http://www.toyota-of-glendale.com/>
3. Leith Toyota—Raleigh, North Carolina. http://www.leithtoyota.com/f_more_about_us.html
4. Atlanta rental cars & auto rentals. <http://www.bnm.com/atl2.htm>

This list includes three relevant results: 1, 2, and 3; and one irrelevant result: 4. The following two relevance feedback methods can be used to improve the next similar query:

- *Query modification:* Adjusts the initial query in an attempt to avoid unrelated or less related query results.
- *Indexing modification:* Through feedback from the users, system administrators can modify an unrelated document's terms to render it unrelated or less related to such a query.

Data Clustering

Data clustering is used to improve the search results by dividing the whole data set into data clusters. Each data cluster contains objects of high similarity, and clusters are produced that group documents relevant to the user's query separately from irrelevant ones [4]. Clustering should not be based on the whole Web resource, but on smaller separate query results. In [57], a Suffix Tree Clustering (STC) algorithm based on phrases shared between documents is used to create clusters. Beside clustering the search results, a proposed similarity function has been used to cluster similar queries according to their contents as well as user logs [54]. The resulting clusters can provide useful information for Frequently Asked Queries (FAQ) identification. Another Web document clustering algorithm is suggested in [12].

5. METASEARCHES

None of the current search engines is able to cover the Web comprehensively. Using an individual search engine may miss some critical information provided by other engines. Metasearch engines [15, 27, 47] conduct a search using several other search engines simultaneously, and present the results in some sort of integrated format. This lets users see at a glance which particular search engine returned the best results for a query without having to search each one individually. They typically do not use their own Web indexes. Figure 3 shows the system structure of a metasearch engine, which consists of three major components:

- *Dispatch*: Determines to which search engines a specific query is sent. The selection is usually based on network and local computational resources, as well as the long-term performance of search engines on specific query terms.
- *Interface*: Adapts the user's query format to match the format of a particular search engine, which varies from engine to engine.
- *Display*: Raw results from the selected search engines are integrated for display to the user. Each search engine also produces different raw results from other search engines and these must be combined and gives a uniform format for ease-of-use.

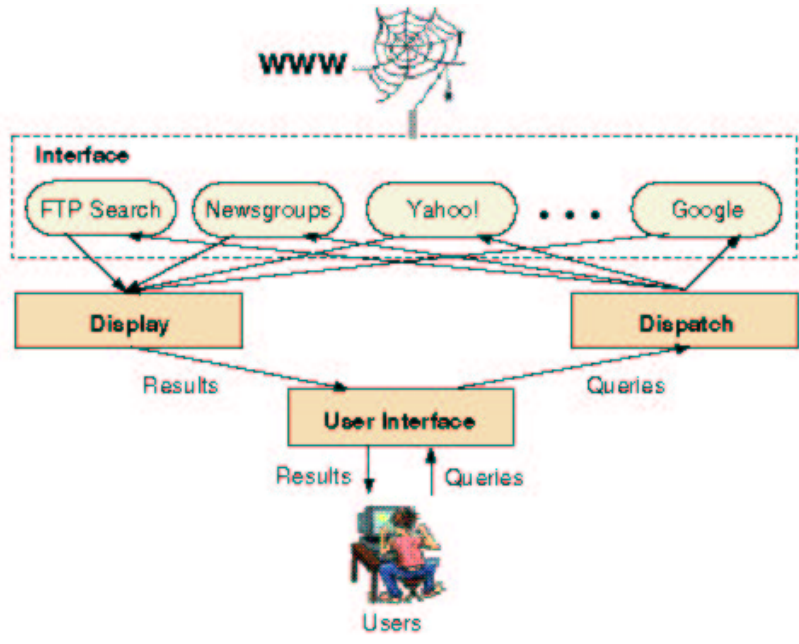


Figure 3: System structure of a metasearch engine.

Current search engines provide a multiplicity of interfaces and results which make the construction of metasearch engines a very difficult task. The STARTS protocol [23] has been proposed to standardize internet retrievals and searches. The goals are to choose the best sources (search engines) to evaluate a query, submit the query to the sources selected, and finally merge the query results obtained from the different sources. However, this protocol has received little recognition since none of the search engines most often used applies it. Another approach [29] to solving this problem is to use an adaptive model which employs a “mediator-wrapper” architecture. The mediator provides users with integrated access to multiple heterogeneous data sources, while each wrapper represents access to a specific data source. It maps a query from a general mediator format into the specific wrapper format of a specific search engine.

Metasearch engines rely on the documents and summaries returned by standard search engines. However, not all standard search engines give unbiased results and this will distort the metasearch results. The NEC Research Institute (NECI) metasearch engine [32] solves this problem by downloading and analyzing each document and then displaying results that show the query term in context. This helps users more readily determine if the document is relevant without having to download each page. The authors of Q-pilot [51] noticed that thousands of specialized, topic-specific search engines are accessible on the Web and these topic-specific engines return far better results

for “on topic” queries than standard search engines. Q-pilot dynamically routes each user query to the appropriate specialized search engines by using two methods: i) neighborhood-based topic identification, and ii) query expansion.

6. SQL APPROACHES

Learning how to use a new language is normally an arduous task for users. However, a new system which uses a familiar language is usually adopted relatively smoothly by the users. SQL (Structured Query Language) is a well-known and widely-used database language. SQL approaches [19, 40] view the World Wide Web as a huge database where each record matches a Web page, and use SQL-like languages to support effective and flexible query processing. A typical SQL-like language syntax [31, 41, 50] is

```
Query := select Attribute_List from Domain_Specifications
        [ where Search_Conditions ];
```

Three query examples are given below to explain the use of the language.

SQL Example 1 *Find pages in the World Wide Web Consortium (W3C) site and the pages have fewer than 2000 bytes.*

```
select url from http://www.w3c.org/ where bytes < 2000;
```

url is a page’s URL and each page has attributes such as *bytes*, *keywords*, and *text*.

SQL Example 2 *Find educational pages containing the keyword “database.”*

```
select url from http://%.edu/ where 'database' in keywords;
```

Regular expressions are widely used in the language, e.g., the ‘%’ is a wild card matching any string. The **in** predicate checks whether the string ‘database’ is one of the keywords.

SQL Example 3 *Find documents about “XML” in the W3C Web site and the documents have paths of length two or less from the root page.*

```

select  d.url d.title
        from Document d such that  'http://www.w3c.org/' =|→|→→ d
        where  d.text like '%XML%';

```

The ‘|’ is an alternation and the ‘→’ is a link. “=|→|→→” is a regular expression that represents the set of paths of length of one or two. The `like` predicate is used for string matching in this example.

Various SQL-like languages have been proposed for Web searches. The methods introduced previously treat the Web as a graph of discrete objects; and another object-oriented approach [2] considers the Web as a graph of structured objects. However, the latter approach has not achieved much success because of its complicated syntax.

7. CONTENT-BASED MULTIMEDIA SEARCHES

In order to allow for the wide range of new types of data which are now available on the World Wide Web, including audio, video, graphics, and images, the use of hypermedia was introduced to extend the capabilities of hypertext. The first internet search engine, Archie, was created in 1990. However, it was not until the introduction of multimedia to the browser Mosaic that the number of Internet documents began to explode. The advent of multimedia has added audio, graphics, images, video, and other types of data to the Web. Only a few multimedia search engines are available currently, most of whom use name or keyword matching where the keywords are entered by Web reviewers rather than using automatic indexing. The reason for the low number of content-based multimedia search engines is mainly due to the difficulty of automated multimedia indexing. Numerous multimedia indexing methods can be found in the literature [11, 55], yet most do not meet the efficiency requirements of Web multimedia searches, which expect both a prompt response and the search of a huge volume of Web multimedia data. A few content-based image and video search engines are available on-line [5, 22, 36, 49, 52]. Various indexing methods are applied to locate the desirable images or video. The major technologies include using camera/object motion, colors, examples, locations, positional color/texture, shapes, sketches, text, and texture as well as relevance feedback [18]. However, a de facto Web image or video search engine is still out of reach because the system’s key component—image or video collection and indexing—is not yet fully au-

tomated or is not practicable enough. Similarly, effective Web audio search engines have yet to be constructed since audio information retrieval [20] is considered to be one of the most difficult challenges for multimedia retrieval.

8. OTHERS

Apart from the above major search techniques, some ad hoc methods worth mentioning include:

- Work aimed at making the components needed for Web searches more efficient, such as better ranking algorithms and more efficient crawlers. In [58], a ranking algorithm based on a Markov model is proposed. It synthesizes the relevance, authority, integrativity, and novelty of each Web resource, and can be computed efficiently through solving a group of linear equations. A variety of other improved ranking algorithms can be found in [16, 48].
- Various enhanced crawlers can be found in the literature [1, 17, 43]. Some crawlers are extensible, personally customized, relocatable, scalable, and Web-site-specific [26, 42]. Web viewers usually consider certain Web pages more important. A crawler which collects those “important” pages first is advantageous for users [13].
- Artificial Intelligence (AI) can also be used to collect and recommend Web pages. The Webnaut system [44] learns the user’s interests and can adapt as his or her interests change over time. The learning process is driven by user feedback to an intelligent agent’s filtered selections.
- A natural language interface designed to make the system easier to use [3].

9. MAJOR SEARCH ENGINES

Some of the currently available major commercial search services are listed in Table 1 where many table entries are unable to be filled because some of the information is considered to be classified material of business [46]. In these days, most search services are backed up by or are cooperating with several other services. An independent or stand-alone service contains less information and tends to lose its users. In the table, the column *Backup* gives the major backup information provider,

Table 1: Major commercial search services. SE: Search Engine, and AS: Answering Service.

No.	Name	URL	Type	Backup	Method
1	AOL Search	http://search.aol.com/	Hybrid SE	Open Directory	
2	AltaVista	http://www.altavista.com/	SE	LookSmart	
3	Ask Jeeves	http://www.askjeeves.com/	AS		natural language
4	Direct Hit	http://www.directhit.com/	SE	HotBot	hyperlink
5	Excite	http://www.excite.com/	SE	LookSmart	
6	FAST Search	http://www.alltheweb.com/			scalability
7	Google	http://www.google.com/	SE		hyperlink
8	HotBot	http://www.hotbot.com/	Hybrid SE	Direct Hit	
9	IWon	http://www.iwon.com/	Hybrid SE	Inktomi	
10	Inktomi	http://www.inktomi.com/	SE		
11	LookSmart	http://www.looksmart.com/	Directory	Inktomi	reviewers
12	Lycos	http://www.lycos.com/	Directory	Open Directory	
13	MSN Search	http://search.msn.com/	Directory	LookSmart	
14	Netscape Search	http://search.netscape.com/	SE	Open Directory	
15	Northern Light	http://www.northernlight.com/	SE		filtering
16	Open Directory	http://dmoz.org/	Directory		volunteers
17	RealNames	http://www.realnames.com/			keywords
18	Yahoo!	http://www.yahoo.com/	Directory	Google	reviewers

and most blank methods are using keyword matching to locate the desired documents. Most search engines on the list not only provide Web search services but also act as portals, which are Web home bases from which users can access a variety of services, including searches, e-commerce, chat rooms, news, etc. Table 2 lists some major experimental search services, which use advanced search technologies not yet implemented by the commercial search services. The list in Table 2 is a snapshot of the current situation; the list is highly volatile because a successful experimental search service is usually commercialized in a short time or a prototype system is normally removed after the founders left the organizations. The above two tables list some major general-purpose search services while some of the special-purpose search services include specialty searches, regional searches, kid searches, etc. which use much small databases and therefore give more precise and limited search results.

Table 2: Major experimental search services.

No.	Name	URL	Method
1	Clever	http://www.almaden.ibm.com/cs/k53/clever.html	hyperlink
2	Grouper	http://longinus.cs.washington.edu/grouper2.html	clustering
3	HuskySearch	http://huskysearch.cs.washington.edu/	metasearch
4	ImageRover	http://www.cs.bu.edu/groups/ivc/ImageRover/Home.html	image
5	ImageScape	http://skynet.liacs.nl/	image
6	Inquirus	http://www.neci.nj.nec.com/homepages/lawrence/inquirus.html	metasearch
7	Mercator	http://www.ctr.columbia.edu/metaseek/	image
8	MetaSEEk	http://www.research.compaq.com/SRC/mercator/	crawler
9	PicToSeek	http://zomax.wins.uva.nl:5345/ret_user/	image
10	W3QS	http://www.cs.technion.ac.il/~konop/w3qs.html	SQL
11	WebOQL	http://www.cs.toronto.edu/~gus/webowl/	Object SQL
12	WebSQL	http://www.cs.toronto.edu/~websql/	SQL

10. SUMMARY

In less than a decade, the World Wide Web has become one of the three major media, with the other two being print and television. Searching for Web pages is both one of the most common tasks performed on the Web and one of the most frustrating. This chapter gave an overview of the Web search technologies currently available with an emphasis on those utilizing non-traditional approaches and classified the technologies into six categories. However, apart from the traditional keyword matching techniques, no one method dominates Web searches. The major reason for this is that none of the search techniques is able to deal effectively and efficiently with the huge volume of information posted on the World Wide Web.

Future Directions

Users of search engines often submit ambiguous queries. Ambiguous queries can be categorized into four types: (i) disorderly, (ii) incomplete, (iii) incorrect, and (iv) superfluous queries. Below are examples of perfect and ambiguous queries and the ranked search results obtained by using Infoseek at <http://www.infoseek.com/> to look for the book “Intelligent multimedia information retrieval,” edited by Mark T. Maybury [39].

- Perfect query: *Intelligent multimedia information retrieval*
 1. Intelligent multimedia information retrieval
- Disorderly query: *Multimedia information intelligent retrieval*
 1. Artificial intelligence, fuzzy Logic and neural networks

- 2. Intelligent access to information: research in natural language, information retrieval, computer vision, multimedia and databases
- 3. Multimedia color PC notebooks
- 4. Intelligent multimedia information retrieval
- Incomplete query: *Multimedia information retrieval*
 - 1. Abstract Stein Mülleler Thiel 95
 - 2. Corpora Oct 1998 to -: Corpora: TWLT 14: language technology in multimedia information
 - 3. 2.1 Introduction to the workplan
 - ...
 - 6. Intelligent multimedia information retrieval
- Incorrect query: *Intelligent multi-media information retrieval*
 - 1. Artificial intelligence research laboratory at Iowa State University
 - 2. Vasant Honavar's home in cyberspace
 - 3. CIIR multi-media indexing
 - ...
 - 31. Intelligent multimedia information retrieval
- Superfluous query: *Intelligent multimedia information retrieval systems*
 - 1. Research in multimedia and multimodal parsing and generation
 - 2. Intelligent multimedia information retrieval

This example shows that even a slight variation in the query produces significant differences among the search results. Keyword matching is used by most search engines, where users tend to submit ambiguous queries. The ambiguity creates undesired search results if keyword matching is used.

Since the introduction of eXtensible Markup Language (XML) [53], more and more Web documents are published in XML. XML document searches [28] are expected to be the next major research direction for Web searches. An XML document not only provides the same information such as keywords, hyperlinks, descriptions, etc, that a function-like HTML document supplies, but also structural information. The structural information is the most crucial feature of an XML document, and is not supplied by an HTML document.

References

- [1] Charu C. Aggarwal, Fatima Al-Garawi, and Philip S. Yu. Intelligent crawling on the World Wide Web with arbitrary predicates. In *Proceedings of the 10th International World Wide Web Conference*, Hong Kong, May 2001.

- [2] Gustavo O. Arocena and Alberto O. Mendelzon. WebOQL: Restructuring documents, databases and Webs. In *Proceedings of the 14th International Conference on Data Engineering*, Orlando, Florida, February 1998.
- [3] Ask Jeeves. <http://www.askjeeves.com/>
- [4] Ricardo A. Baeza-Yates. Introduction to data structures and algorithms related to information retrieval. In William B. Frakes and Ricardo A. Baeza-Yates, editors, *Information Retrieval Data Structures & Algorithms*, pages 13-27, Prentice-Hall, 1992.
- [5] Ana B. Benitez, Mandis Beigi, and Shih-Fu Chang. Using relevance feedback in content-based image metasearch. *IEEE Internet Computing*, 2(4):59-69, July/August 1998.
- [6] Krishna Bharat and Monika Henzinger. Improved algorithms for topic distillation in a hyperlinked environment. In *Proceedings of the 21st International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 104-111, August 1998.
- [7] Allan Borodin, Gareth O. Roberts, Jeffrey S. Rosenthal, and Panayiotis Tsaparas. Finding authorities and hubs from link structures on the World Wide Web. In *Proceedings of the 10th International World Wide Web Conference*, Hong Kong, May 2001.
- [8] Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems*, 30:107-117, 1998.
- [9] Jeromy Carriere and Rick Kazman. WebQuery: Searching and visualizing the Web through connectivity. *Computer Networks and ISDN Systems*, 29(11):1257-1267, 1997.
- [10] Soumen Chakrabarti, Byron E. Dom, S. Ravi Kumar, Prabhakar Raghavan, Sridhar Rajagopalan, Andrew Tomkins, David Gibson, and John Kleinberg. Mining the Web's link structure. *IEEE Computer*, 32(8):60-67, August 1999.
- [11] Shi-Kuo Chang and Arding Hsu. Image information systems: Where do we go from here? *IEEE Transactions on Knowledge and Data Engineering, Special Issue Celebrating the 40th Anniversary of the Computer Society*, 4(5):431-442, October 1992.
- [12] Chia-Hui Chang and Ching-Chi Hsu. Enabling concept-based relevance feedback for information retrieval on the WWW. *IEEE Transactions on Knowledge and Data Engineering*, 11(4):595-609, July/August 1999.

- [13] Junghoo Cho, Hector Garcia-Molina, and Lawrence Page. Efficient crawling through URL ordering. In *Proceedings of the 7th World Wide Web Conference*, Brisbane, Australia, April 1998.
- [14] Jeffrey Dean and Monika R. Henzinger. Finding Related Web Pages in the World Wide Web. In *Proceedings of the 8th International World Wide Web Conference*, pages 389-401, Toronto, Canada, 1999.
- [15] Daniel Dreilinger and Adele E. Howe. Experiences with selecting search engines using metasearch. *ACM Transactions on Information Systems*, 15(3):195-222, July 1997.
- [16] Cynthia Dwork, Ravi Kumar, Moni Naor, and D. Sivakumar. Rank aggregation methods for the Web. In *Proceedings of the 10th International World Wide Web Conference*, Hong Kong, May 2001.
- [17] Jenny Edwards, Kevin McCurley, and John Tomlin. An adaptive model for optimizing performance of an incremental Webcrawler. In *Proceedings of the 10th International World Wide Web Conference*, Hong Kong, May 2001.
- [18] Myron Flickner, *et. al.* Query by image and video content: The QBIC system. *IEEE Computer*, 28(9):23-32, September 1995.
- [19] Daniela Florescu, Alon Levy, and Alberto Mendelzon. Database techniques for the World Wide Web: A survey. *ACM SIGMOD Record*, 27(3):59-74, September 1998.
- [20] Jonathan Foote. An overview of audio information retrieval. *Multimedia Systems*, 7(1):2-10, 1999.
- [21] John Garofalakis, Panagiotis Kappos, and Dimitris Mouloukos. Web site optimization using page popularity. *IEEE Internet Computing*, 3(4):22-29, July/August 1999.
- [22] Theo Gevers and Arnold Smeulders. The PicToSeek WWW image search system. In *Proceedings of the IEEE International Conference on Multimedia Computing and Systems*, Florence, Italy, pages 264-269, June 1999.
- [23] Luis Gravano, Kevin Chang, Hector Garcia-Molina, Carl Lagoze, and Andreas Paepcke. STARTS: Stanford protocol proposal for internet retrieval and search. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 1997.

- [24] Ilan Greenberg and Lee Garber. Searching for new search technologies. *IEEE Computer*, 32(8):4-11, August 1999.
- [25] Kenkat N. Gudivada, Vijay V. Raghavan, William I. Grosky, and Rajesh Kasanagottu. Information retrieval on the World Wide Web. *IEEE Internet Computing*, 1(5):58-68, September/October 1997.
- [26] Allan Heydon and Marc Najork. Mercator: A scalable, extensible web crawler. *World Wide Web*, 2(4):219-229, 1999.
- [27] Adele E. Howe and Daniel Dreilinger. SavvySearch: A meta-search engine that learns which search engines to query. *AI Magazine*, 18(2), 1997.
- [28] Wen-Chen Hu, Yapin Zhong, Wei-Chuan Lin, and Jui-Fa Chen. An XML World Wide Web search engine using approximate structural matching. In *Proceedings of the 5th World Multi-Conference on Systemics, Cybernetics and Informatics*, Orlando, Florida, July 22-25, 2001.
- [29] Lieming Huang, Matthias Hemmje, and Erich J. Neuhold. ADMIRE: An adaptive data model for meta search engines. *Computer Networks (The International Journal of Computer and Telecommunications Networking)*, 33(1-6):431-448, 2000.
- [30] Jon M. Kleinberg. Authoritative sources in a hyperlinked environment. *JACM*, 46(5):604-632, September 1999.
- [31] David Konopnicki and Oded Shmueli. Information gathering in the World Wide Web: The W3QL query language and the W3QS system. *ACM Transactions on Database Systems*, 23(4):369-410, December 1998.
- [32] Steve Lawrence and C. Lee Giles. Context and page analysis for improved Web search. *IEEE Internet Computing*, 2(4):38-46, July/August 1998.
- [33] Steve Lawrence and C. Lee Giles. Searching the World Wide Web. *Science*, 280:98-100, 1998.
- [34] Steve Lawrence and C. Lee Giles. Accessibility of information on the Web. *Nature*, 400:107-109, 1999.
- [35] Steve Lawrence and C. Lee Giles. Searching the Web: General and scientific information access. *IEEE Communications*, 37(1):116-122, 1999.

- [36] Michael S. Lew. Next generation Web searches for visual content. *IEEE Computer*, 33(11):46-53, November, 2000.
- [37] Yanhong Li. Toward a qualitative search engine. *IEEE Internet Computing*, 2(4):24-29, July/August 1998.
- [38] Hongjun Lu and Ling Feng. Integrating database and World Wide Web technologies. *World Wide Web*, 1(2):73-86, 1998.
- [39] Mark T. Maybury. Intelligent multimedia information retrieval. MIT Press, 1997.
- [40] Alberto O. Mendelzon and Tova Milo. Formal models of Web queries. *Information Systems*, 23(8):615-637, 1998.
- [41] Alberto O. Mendelzon, George Mihaila, and Tova Milo. Querying the World Wide Web. *International Journal on Digital Libraries*, 1(1):54-67, 1997.
- [42] Robert C. Miller and Krishna Bharat. SPHINX: A framework for creating personal, site-specific Web crawlers. In *Proceedings of the 7th International World Wide Web Conference*, Brisbane, Australia, April 1998.
- [43] Marc A. Najork and Janet Wiener. Breadth-first search crawling yields high-quality pages. In *Proceedings of the 10th International World Wide Web Conference*, pages 114-118, Hong Kong, May 2001.
- [44] Zacharis Z. Nick and Panayiotopoulos Themis. Web search using a genetic algorithm. *IEEE Internet Computing*, 5(2):18-26, March/April 2001.
- [45] Davood Rafiei and Alberto O. Mendelzon. What is this page known for? Computing Web page reputations. In *Proceedings of the 9th International World Wide Web Conference*, Amsterdam, Netherlands, May 2000.
- [46] Search Engine Watch. <http://www.searchenginewatch.com/>
- [47] Erik Selberg and Oren Etzioni. The MetaCrawler architecture for resource aggregation on the Web. *IEEE Expert*, 12(1):8-14, January/February 1997.
- [48] Amit Singhal and Marcin Kaszkiel. A case study in Web search using TREC algorithms. In *Proceedings of the 10th International World Wide Web Conference*, pages 708-716, Hong Kong, May 2001.

- [49] John R. Smith and Shih-Fu Chang. An image and video search engine for the World-Wide Web. In *Proceedings of the Symposium on Electronic Imaging: Science and Technology - Storage and Retrieval for Image and Video Databases V, IS&T/SPIE*, San Jose, California, February 1997.
- [50] Ellen Spertus and Lynn Andrea Stein. Squeal: A structured query language for the Web. In *Proceedings of the 9th International World Wide Web Conference*, Amsterdam, Netherlands, May 2000.
- [51] Atsushi Sugiura and Oren Etzioni. Query routing for Web search engines: Architecture and experiments. In *Proceedings of the 9th International World Wide Web Conference*, Amsterdam, Netherlands, May 2000.
- [52] Leonid Taycher, Marco La Cascia, and Stan Sclaroff. Image digestion and relevance feedback in the ImageRover WWW search engine. In *Proceedings of the International Conference on Visual Information*, San Diego, December, 1997.
- [53] eXtensible Markup Language (XML). <http://www.w3.org/XML/>
- [54] Ji-Rong Wen, Jian-Yun Nie, and Hong-Jiang Zhang. Clustering user queries of a search engine. In *Proceedings of the 10th International World Wide Web Conference*, Hong Kong, May 2001.
- [55] Atsuo Yoshitaka and Tadao Ichikawa. A survey on content-based retrieval for multimedia databases. *IEEE Transactions on Knowledge and Data Engineering*, 11(1):81-93, January/February 1999.
- [56] Budi Yuwono and Dik Lun Lee. WISE: A World Wide Web resource database system. *IEEE Transactions on Knowledge and Data Engineering*, 8(4):548-554, August 1996.
- [57] Oren Zamir and Oren Etzioni. Web document clustering: A feasibility demonstration. In *Proceedings of the 19th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 46-54, Melbourne, Australia, 1998.
- [58] Dell Zhang and Yisheng Dong. An efficient algorithm to rank Web resources. In *Proceedings of the 9th International World Wide Web Conference*, Amsterdam, Netherlands, May 2000.