

# Using a Graph-Based Data Mining System to Perform Web Search

**Diane J. Cook, Nitish Manocha, and Lawrence B. Holder**  
**Department of Computer Science and Engineering**  
**University of Texas at Arlington**  
**{cook, manocha, holder}@cse.uta.edu**  
**<http://ailab.uta.edu/subdue>**

## Abstract

The World Wide Web provides an immense source of information. Accessing information of interest presents a challenge to scientists and analysts, particularly if the desired information is structural in nature. Our goal is to design a structural search engine that uses the hyperlink structure of the Web, in addition to textual information, to search for sites of interest. Our structural search engine, called WebSUBDUE, searches not only for particular words or topics but also for a desired hyperlink structure. Enhanced by WordNet text functions, our search engine retrieves sites corresponding to structures formed by graph-based user queries. We hypothesize that this system can form the heart of a structural query engine, and demonstrate the approach on a number of structural web queries.

## 1 Introduction

The World Wide Web (WWW) provides an immense source of information. When responding to a query, web search engines sift through keywords collected from millions of pages. However, these search engines ignore the importance of searching the structure of

hyperlink information and cannot resolve complicated queries like *'find all pages that link to another page containing term x'*. A hyperlink from one page to another represents an internal organization that a web page designer wants to establish, or relationships that exist between web sites. A structural search engine can search not only for particular words or topics but also for a desired hyperlink structure.

Web information can be represented as a labeled graph. SUBDUE (Cook, 2000) is a data mining tool that discovers repetitive substructures in graph-based data. We therefore hypothesize that SUBDUE can form the heart of a structural search engine, called WebSUBDUE. To verify our hypothesis and demonstrate the use of WebSUBDUE as a structural search engine, we perform various experimental queries on the Computer Science and Engineering Department of the University of Texas at Arlington web site and compare the results with those of posting similar queries to the search engine, Altavista (Altavista).

## **2 Background**

Structural web search is the process of searching the web for a specific hyperlink structure combined with textual content. Sometimes, it is not sufficient to apply purely text-based methods to find a large number of potentially relevant pages. People are likely to surf the web using its graph structure. Current web search engines can be used in order to search for some keywords or some combination of keywords without forcing any hyperlink structure between web pages. Consequently, the results of a particular search engine would be a number of hits, each containing one web page.

In response to a user query, a structural search engine searches not only the text but also the structure formed by the set of pages to find matches to the query. The web consists of hyperlinks that connect one page to another, and this hyperlink structure contains an enormous amount of information. Specifically, the creation of a hyperlink by the author of a web page represents a relationship between the source and destination pages. Such relationships may include a hierarchical relation between a top-level page and a child page containing more detailed information, relationships between the current, previous and next sites in a sequence of pages, or an implicit endorsement of a page which represents an authority on a particular topic. A structural search engine utilizes this hyperlink structure along with the textual content to find sites of interest. Structural web queries can locate online presentations, web communities, or nepotistic sites. By searching the organization of a set of pages as well as the content, a search engine can provide richer information on the content of web information.

Even a purely textual search can face challenges due to the fact that human language is very expressive and full of synonymy (different words having the same meaning). A query for 'automobile', for example, may miss a deluge of pages that focus on the term 'car'. One strategy to overcome this limitation is to augment search techniques with stored information about semantic relations between words. Such compilations, typically constructed by a team of linguists, are sometimes known as semantic networks, following the seminal work on the WordNet project (Miller et al., 1991). A structural search engine could make use of this information to improve the retrieval of web sites that reflect the structure and the semantics of the user query.

Much research has focused on using hyperlink information in some method to enhance web search. The Clever system (Clever) scans the most central, authoritative sites on broad search topics by making use of hyperlinks. A respected authority is a page that is referenced by many good hubs; a useful hub is a location that points to many valuable authorities. Clever uses an iterative process to identify the best hub and authority pages for a particular topic based on the strength of the hub/authority pages to which they are connected. Google (Brin and Page, 1999) is another search engine that uses the power of hyperlinks to search the web. Google evaluates a site by summing the scores of other locations pointing to the page. When presented with a specific query, Google ranks relevant pages by their evaluation measure.

Although these systems use certain types of hyperlink structures to rank retrieved web pages, they do not perform searches for a range of structural queries. In contrast, WebSUBDUE searches for any type of query describing structure embedded with textual content. WebSUBDUE's functionality is enhanced using tools provided by the lexicon database WordNet.

### **3 Data Preparation**

Viewed as a graph, the World Wide Web structure exhibits properties that are of interest to a number of researchers (Kleinberg et al., 1999). We transform web data to a labeled graph for input to the WebSUBDUE system. Data collection is performed using a web robot written in Perl. The web robot only follows links to pages residing on specified servers. For example, if `http://cygnus.uta.edu` is a defined server then the web robot will explore the URL `http://cygnus.uta.edu/projects.html`, but a URL outside the domain, such as `http://ranger.uta.edu`,

will not be explored. As it traverses a web site, the web robot generates a graph file representing the specified site, described as a set of labeled vertices with labeled edges connecting the vertices.

The web robot scans each page for URL references contained in that page. A depth-first search through the space of connected web pages on the specified site is executed. The Perl script uses the built-in hash data type / associative array to store, for each web page, a set of URLs that contain links to the web page. For example, if URL 'A' and URL 'B' link to URL 'C', then the hash table entry for key 'C' is defined as *Hash\_Variable*{'C'} = 'A B'. In this case, URLs 'A' and 'B' are considered to be parent pages, and URL 'C' is considered to be the child page. If the child URL has already been defined, then the parent URL is appended to its values. If the child URL has not been defined, then a new entry is created in the associative array with the child URL as the key and the parent URL as the value.

All the URLs in the current page are checked to see if they have been visited. If a URL has not been visited, it is pushed onto the stack and marked as visited. When all the URLs in the page have been pushed onto the stack, the topmost URL is popped off of the stack, and this URL becomes the new parent.

For each new parent URL on the specified site, the robot fetches the page header using the 'request' routine in the Perl HTTP module and checks the page fetch response code. If the response code is RC\_OK, then the page currently exists and is added to the graph. The page is then scanned for URLs and all reference entries are added to the associative array. This process continues until the stack is empty or the number of visited URLs exceeds the user-defined limit.

Once the URLs have been scanned and the associative array has been created, a labeled graph is generated representing the website. In this graph, each URL is represented as a vertex labeled ‘*\_page\_*’, and an edge labeled ‘*\_hyperlink\_*’ points from parent URLs to child URLs (from each value of a *Hash\_Variable* entry to the key itself).

After adding an edge to the graph representing a hyperlink relation, the script determines the document type by examining the ‘Content-Type’ tag in the page header. If the document is either an HTML file or a text file, the file is processed to obtain keyword information. All of the HTML tags, numbers, punctuation, single characters, articles, prepositions, pronouns and conjunctions are removed from the file, and all the remaining unique words are added to the graph. A vertex labeled with the corresponding word is added to the graph, and an edge labeled ‘*\_word\_*’ connects the page in which the word is found to the word vertex.

The current implementation of this search engine allows the user to create a graph for a new domain or search an existing graph. An example graph representing a collection of web pages for domain cygnus.uta.edu is shown in Figure 1. WebSUBDUE can process structural queries on this graph. The input file containing the graph information is also shown in Figure 1. All of the structural queries follow the input representation shown here. In the input file, a graph is represented as a set of vertices (denoted by “v”, a vertex id, and a label) and a set of edges (denoted by “d” for directed edge or “u” for undirected edge, the source and destination vertex ids, and a label). Hits are reported by the vertex id and corresponding URL.

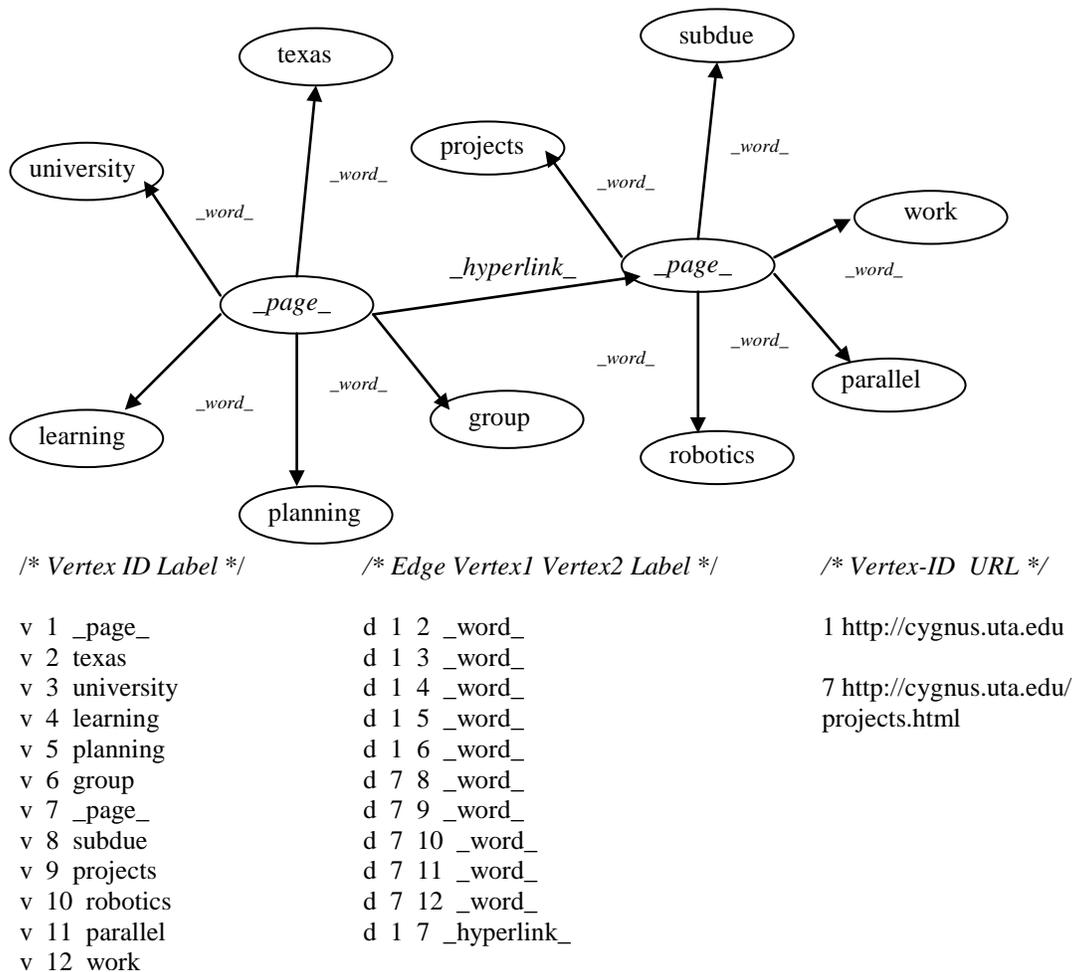


Figure 1. Graph representation of a web site.

Although the web robot creates a single graph through this process, new web sites can be added to an existing graph, allowing web pages to be incrementally added. In the current implementation, only files which are reachable from a root page are considered, which may exclude some pages. The types of files which are currently handled by the robot contain extensions .html (.htm), .txt, .exe, .gif, .pdf, .ps, .tar, .gz, .jpg, or .mpg.

## 4 The SUBDUE Knowledge Discovery System

SUBDUE is a knowledge discovery system that discovers patterns in structural data. SUBDUE accepts data in the form of a labeled graph and performs various types of data mining on the graph. As an unsupervised discovery algorithm, SUBDUE discovers repetitive patterns, or subgraphs, in the graph. As a concept learner, the system identifies concepts that describe structural data in one class and exclude data in other classes. The system can also discover structural hierarchical conceptual clusters. In this paper, we introduce the application of SUBDUE as a structural search engine, called WebSUBDUE, in which the search query and the WWW are represented as labeled graphs. The discovered instances are reported as the results of the query. SUBDUE accepts input in the form of a labeled graph, where objects and attribute values map to vertices, and relationships between objects map to edges.

### 4.1 Discovery Search Algorithm

When performing unsupervised knowledge discovery, SUBDUE uses a polynomial-time beam search for its main discovery algorithm. The goal of SUBDUE's search is to find the subgraph which yields the best compression of the input graph. A substructure in SUBDUE consists of a subgraph defining the substructure and all of the substructure instances throughout the graph. The initial state of the search is the set of subgraphs consisting of all uniquely labeled vertices. The only search operator is the *Extend Subgraph* operator, which extends a subgraph in all possible ways by a single edge and a vertex or by a single edge only if both vertices are already in the subgraph. Substructures are kept on an ordered queue of length determined by the beam width.

To evaluate subgraphs, the Minimum Description Length (MDL) principle is used, which states that the best theory is one that minimizes the description length of a database (Rissanen, 1989). Using this principle, a substructure is evaluated by its ability to minimize the description length of the graph when compressed using the substructure.

The search terminates upon reaching a user-specified limit on the number of substructures extended, or upon exhaustion of the search space. Once the search terminates and returns the list of best subgraphs, the graph can be compressed using the best subgraph. The compression procedure replaces all instances of the subgraph in the input graph by a single representative vertex. Incoming and outgoing edges to and from the replaced subgraph will connect to the new vertex that represents the subgraph. The SUBDUE algorithm can iterate on this compressed graph to generate a hierarchical description of discovered substructures.

## **4.2 Search for Predefined Substructures**

SUBDUE supports biasing the discovery process toward specified types of substructures. Predefined substructures can be provided as input by creating a predefined substructure graph that is stored in a separate file. SUBDUE will try to locate and expand these substructures, in this way “jump-starting” the discovery process. The inclusion of background knowledge has been shown to be of great benefit. This allows the user to provide SUBDUE with background knowledge about the kinds of substructures for which the database is being mined.

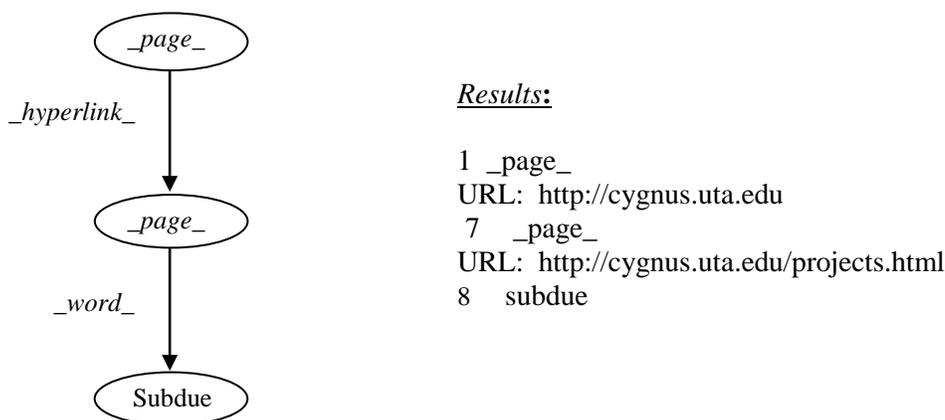


Figure 2. Structural search query.

For our structural search engine, WebSUBDUE invokes SUBDUE in predefined substructure mode, where the search query is represented in the form of a substructure. SUBDUE discovers the instances of the predefined substructure in the graph, representing the results of the query. WebSUBDUE reports the graph vertices, edges, and corresponding URLs for each discovered instance. For example, Figure 2 shows the representation of a structural search query to ‘find all pages which link to a page containing the term *Subdue*’, along with the results of applying the search to the graph shown in Figure 1.

### 4.3 Inexact Graph Match

In many databases, patterns exist in the database with small variations between instances. SUBDUE applies a polynomial-time inexact match algorithm that allows for small differences between a substructure definition and a substructure instance. The dissimilarity of two graphs is determined by the number of transformations needed to make one graph isomorphic to the other. The allowed transformations include adding or deleting an edge or vertex, changing an edge or vertex label, and reversing the direction of an edge. Two graphs match if the number of required

transformations is no more than a user-defined threshold value multiplied by the size of the larger graph. If the threshold value is defined as 0, an exact match is required.

Figure 3 shows an example pair of graphs (one could represent a substructure definition, and the other a possible instance of the substructure appearing somewhere in the input graph). The least-cost transformation to make graph G1 isomorphic to G2 would be to map the node labeled “A” in G1 to the node labeled “A” in G2, and to map the node labeled “B” in G2 to the top left node labeled “B” in G2. The corresponding cost is calculated as the cost for adding a node (the bottom node labeled “B” in G2) and two corresponding edges. Because the size of the larger graph is 3 vertices + 4 edges = 7, graph G1 would be considered an instance of G2 if the threshold value is  $3/7$  or greater.

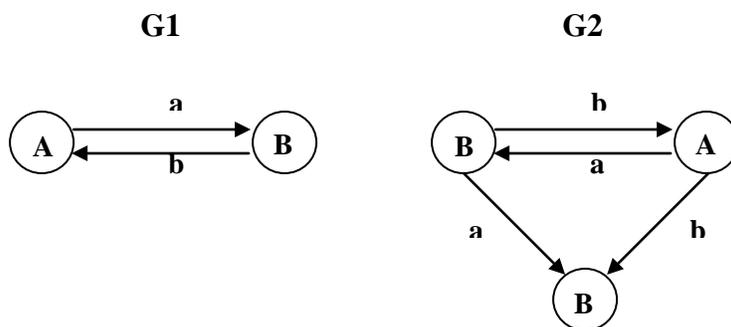


Figure 3. Example inexact graph match.

To find the least-cost match, SUBDUE searches in a branch-and-bound fashion through the space of possible matches. The first complete match found is thus the least-cost match and is returned. If the number of partial matches considered exceeds a user-defined function of the size of the larger graph, the search switches to a hill climbing search. As a result, the inexact graph

match is constrained to run in time polynomial in the size of the larger graph. SUBDUE's inexact graph match algorithm can be used by WebSUBDUE to efficiently find web sites that closely, but not exactly, match the user query.

## 5 WordNet

WordNet is an electronic lexicon database which attempts to organize information according to the meanings of the words instead of the forms of the words (Miller et al., 1991). The WordNet database and related tools can be downloaded from [www.cogsci.princeton.edu/~wn](http://www.cogsci.princeton.edu/~wn). In this project, our goal is to augment the structural search technique with stored information about relations between words. In particular, WebSUBDUE uses WordNet to increase its search capacity by searching for words similar to the query term rather than searching only for the query term itself.

WordNet contains standard information found in dictionaries and thesauri, but also stores relationships between words including hypernyms, hyponyms and synonyms. A hypernym of a term is a more general term that fits the statement “\_\_\_ is a kind of \_\_\_.” For example, a dog is a kind of canine, so canine is a hypernym of dog. Hyponyms are the opposite of hypernyms. Different words having the same meaning are called synonyms.

WordNet provides several ways to identify related terms. Synsets provide one method of grouping synonym terms. A word can be part of several different synsets, each representing a different context. WordNet also has morphological features. Although word base forms are usually stored in WordNet, searches may be performed on inflected forms. We apply a set of morphology functions to the input search string, generating a form that is present in WordNet.

For example, if a user searches for the word *teaching*, WordNet will return the base form *teach*. Similarly, if the query term *Oct* is not present in WordNet, WordNet will return the inflected form *October* which does reside in its dictionary.

We enhance the capabilities of WebSUBDUE's structural search engine by integrating some of WordNet's functionality. WebSUBDUE uses WordNet to enhance its search capabilities by finding words similar to the search terms. In particular, vertices representing a word within a page are matched with vertices containing a different label if they hold one of the valid relationships defined by WordNet. Valid considerations include matching the base form of the word (e.g., match *base* or *basis* with query *bases*), derived forms of the word (e.g., match *Jan* with query *January*), and synonyms (e.g., match *employment*, *position*, or *work* with query *job*).

## **6 Experimental Results**

We conducted a number of experiments to evaluate the capabilities of the WebSUBDUE structural search engine. Where possible, we compare query results of WebSUBDUE with search results generated using a popular search engine, Altavista. We select Altavista because it offers features such as a simple keyword search, searches with advanced options, image searches, audio searches and video searches. Altavista's advanced search features include the use of Boolean expressions, the limitation of a search to a particular host or domain, and other search criteria that provide a valuable point of comparison for the results discovered by WebSUBDUE.

## 6.1 Test Domain

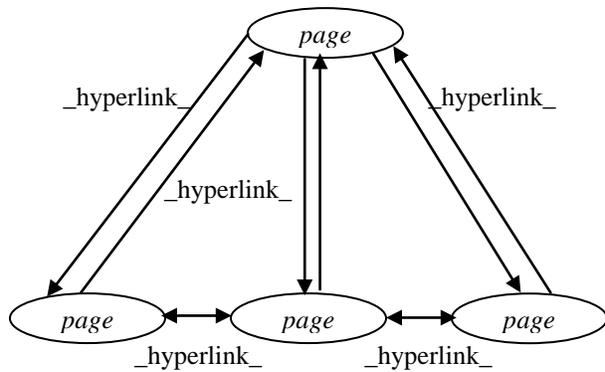
To demonstrate the use of WebSUBDUE as a structural search engine, we restrict its search space to a single sample domain. The sample domain chosen is the web domain for the Computer Science and Engineering department (CSE) of the University of Texas at Arlington (UTA). A graph file of the CSE-UTA domain was created using the web robot described earlier. The test domain includes all web pages starting with the root URL <http://www-cse.uta.edu>. The graph representation of the CSE-UTA web site contains 113,933 vertices and 125,657 edges covering 5,825 URLs.

## 6.2 Structural Searching

To demonstrate the structural searching capability of WebSUBDUE, we posed several queries to the system and compared the results to similar searches using Altavista. The search range of Altavista was also restricted to only the CSE-UTA domain by using the define-host advanced search option in Altavista. All results were found within the [www-cse.uta.edu](http://www-cse.uta.edu) domain, thus we eliminate the domain name in our summaries.

### 6.2.1 Searching for Presentation Pages

One type of structural query might be to find online lecture materials or HTML papers discussing a particular topic. A query to search for all presentation-style pages was posed to WebSUBDUE. A presentation page is defined here as series of pages focusing on a topic, where each page contains links to the *next* page, the *previous* page and the *top* page, the top page being a menu page which has links to all the other pages. This query structure is shown in Figure 4.



1. ~holder/courses/cse5311/lectures/126/126.html
2. ~cook/ai1/syllabus/syllabus.html
3. ~holder/courses/cse2320/lectures/13/13.html
4. ~holder/courses/cse5311/lectures/112/112.html
5. ~holder/courses/cse2320/lectures/116/116.html
6. ~holder/courses/cse2320/lectures/110/110.html
7. ~holder/courses/cse5311/lectures/119/119.html
8. ~holder/courses/cse5311/lectures/19/19.html
9. ~cook/ai1/lectures/113/113.html
10. ~cook/ai1/lectures/11/11.html
11. ~reyes/node48.html
12. ~reyes/node1.html
13. ~cook/ai1/lectures/16/node6.html
14. ~cook/ai1/lectures/17/17.html
15. ~holder/courses/cse5311/lectures/123/123.html
16. ~holder/courses/cse2320/lectures/12/12.html
17. ~cook/ai1/lectures/110/110.html
18. ~cook/ai1/lectures/112/112.html
19. ~holder/courses/cse5311/lectures/110/110.html
20. ~holder/courses/cse2320/lectures/115/115.html
21. ~holder/courses/cse5311/lectures/14/14.html
22. ~holder/courses/cse5311/lectures/12/12.html

Figure 4. Presentation page structural query.

In response to the presentation page query, WebSUBDUE discovered 22 unique instances or 22 presentation pages on the CSE-UTA site. The discovered URLs are listed in Figure 4. These discovered web pages all presentation pages covering various topics with a top menu page containing a link to sub-topic pages and each sub topic page pointing to the next sub topic page, the previous page and the top menu page.

The Altavista search engine has no method of responding to structural search queries. To determine whether Altavista could find the sites discovered by WebSUBDUE, we provide Altavista with additional information. All the discovered pages represent documents that were generated using the LaTeX2HTML conversion program and contain links to the icon files *next\_motif.gif*, *up\_motif.gif* and *previous\_motif.gif*, representing contiguous sections of the document. To assist Altavista, we provide information about icon file links using the advanced

search option: “*host:www-cse.uta.edu AND image:next\_motif.gif AND image:up\_motif.gif AND image:previous\_motif.gif.*” This search option asks Altavista to search in the CSE-UTA domain for pages containing links to next icon, previous icon, and top icon files.

Altavista discovered 12 instances, and WebSUBDUE discovered 22 instances. Because the actual names of the icons varied between presentations, not all presentation pages were discovered using Altavista. This shows the difficulty of posing these types of structural queries using a text search engine, whereas WebSUBDUE just requires the structure information.

### **6.2.2 Searching for Pages With At Least $N$ Inlinks**

Another type of structural query is to find pages with at least  $n$  links, where the value of  $n$  is given in the query structure. This type of page could be labeled as a reference site for a particular topic, or a site that is recommended by many other sites. WebSUBDUE can use its structural search capabilities to find all the pages with at least a specified number of inlinks, which are pages with at least a specified number of other pages containing a hyperlink to the reference page. We posed this query on the CSE-UTA domain to find pages having at least 35 inlinks. Figure 5 shows the structure used as an input query to WebSUBDUE and the results of the query. Most of the retrieved sites represent icon files used by many web pages. In contrast, Altavista cannot perform this type of structural query to find pages with a minimum number of inlinks.

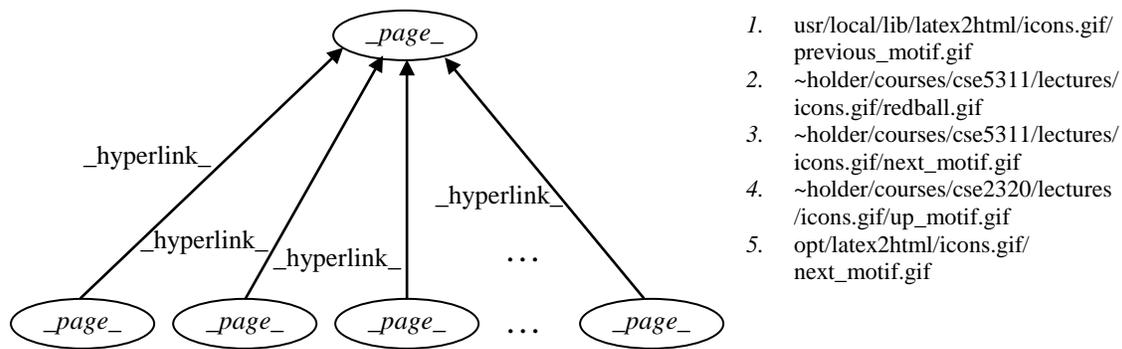


Figure 5. Structural query to find pages with at least 35 inlinks.

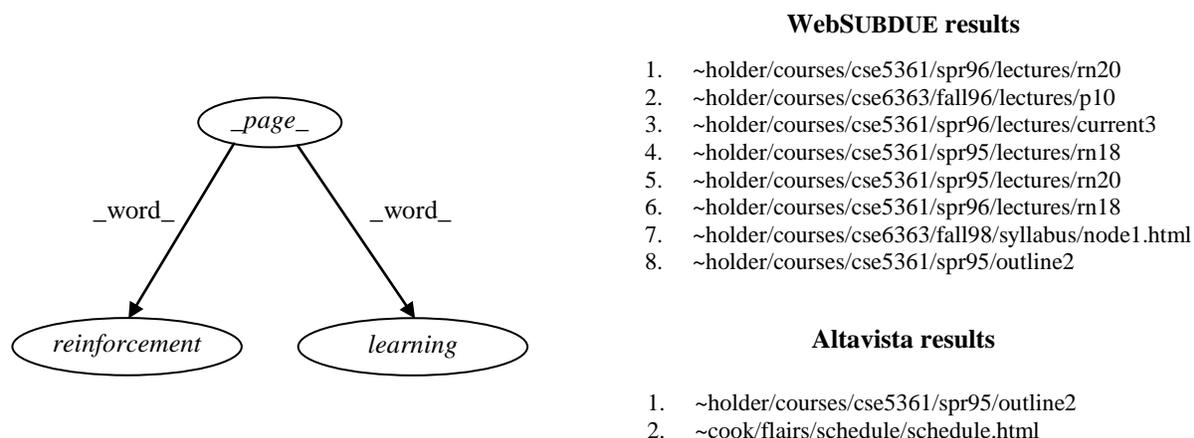
### 6.3 Text Searching

WebSUBDUE can perform text searching as well as structure searching. To demonstrate the text searching capabilities of WebSUBDUE, various search queries were run on the CSE-UTA domain and the results compared with the results of similar queries posted to Altavista.

#### 6.3.1 Searching for pages on ‘Reinforcement Learning’

A query to search for all pages on ‘*reinforcement learning*’ was posed to WebSUBDUE. WebSUBDUE searched for pages containing the words “*reinforcement*” and “*learning*”. The search query substructure is shown in Figure 6, along with results for WebSUBDUE and Altavista.

WebSUBDUE found 8 hits for this query, representing pages that focus on reinforcement learning or contain the words “reinforcement” and “learning” in the text. A similar query to find pages on reinforcement learning was posed to Altavista with the following option: “*host:www-cse.uta.edu AND text:reinforcement AND text:learning*”.



1. ~holder/courses/cse5361/spr96/lectures/rn20
2. ~holder/courses/cse6363/fall96/lectures/p10
3. ~holder/courses/cse5361/spr96/lectures/current3
4. ~holder/courses/cse5361/spr95/lectures/rn18
5. ~holder/courses/cse5361/spr95/lectures/rn20
6. ~holder/courses/cse5361/spr96/lectures/rn18
7. ~holder/courses/cse6363/fall98/syllabus/node1.html
8. ~holder/courses/cse5361/spr95/outline2

1. ~holder/courses/cse5361/spr95/outline2
2. ~cook/flairs/schedule/schedule.html

Figure 6. Structural query to find pages on ‘Reinforcement Learning’.

Altavista discovered 2 instances of the search pattern, as shown in Figure 6. Instance (1) discovered by Altavista was also discovered by WebSUBDUE. Instance (2) is not discovered by WebSUBDUE, because the page cannot be accessed from any CSE-UTA page and is therefore not included by the web robot. The other pages discovered by WebSUBDUE were overlooked by Altavista.

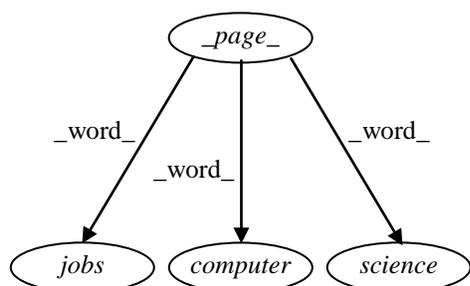
### 6.3.2 Searching for Pages on ‘jobs in computer science’ Using Synonyms

WebSUBDUE can be used with the synonym searching capability to find pages containing text on ‘jobs in computer science’. The structural query to perform this search is shown in Figure 7. WebSUBDUE discovered 33 instances of the search query, as summarized in Figure 7 with the matched terms found in the corresponding page. WebSUBDUE finds ‘*employment*’, ‘*work*’, ‘*job*’, ‘*problem*’, ‘*task*’ as synonyms and allowable forms of the search terms. Only exact matches were found for ‘*computer*’ and ‘*science*’ in these pages.

A similar query was posed to Altavista using the syntax “*host:www-cse.uta.edu AND text:jobs AND text:computer AND text:science*”. Altavista discovered 2 instances of the search query. WebSUBDUE can thus cover a wider search range using synonyms along with the keywords themselves.

#### WebSUBDUE results

1. about\_uta\_engr.html (employment)
2. Undergrad/undergrad.html (work)
3. ~cook/ai1/hw/h16 (job)
4. ~cook/grt/index.html (job)
5. Undergrad/bscse.html (work)
6. Graduate/mse.html (work)
7. ~cook/pai/pai.html (problem)
8. jobs.html (jobs)
9. ~reyes/node11.html (jobs)
10. ~al-khaiy/cse3320\_summer99/summer99.htm (task)
11. ~al-khaiy/cse3322\_fall99/ann/syllabus.htm (work)
12. ~umbaugh/Course-Syllabus-Spring-2000.html (work)
13. ~holder/research/ugv.html (work)
14. Undergrad/97/97program\_info\_1.html (work)
15. ~piotr/cse5361/spr2000/syllabus/syllabus.html (work)
16. ~holder/courses/cse2320/spr99/syllabus/syllabus.html (work)
17. ~holder/courses/cse2320/fall98/syllabus/syllabus.html (work)
18. Staff/gta\_fall99.html (work)
19. Graduate/mse.html (work)
20. ~holder/courses/cse5361/spr99/syllabus/syllabus.html (work)
21. Graduate/graduate\_overview.html (work)
22. Undergrad/97/97program\_info\_2.html (job)
23. ~holder/courses/cse6363/fall96/lectures/final (problem)
24. ~holder/courses/cse5361/spr96/README (problem)
25. ~holder/courses/cse2320.html (problem)
26. ~ramirez/pubs/confs/FLAIRS98.htm (problem)
27. ~holder/research/ml.html (problem)
28. ~holder/courses/cse5311.html (problem)
29. ~holder/courses/cse5361/spr95/README (problem)
30. ~holder/courses/cse5311/summ98/syllabus/syllabus.html (problem)
31. ~cook/aa/syllabus (problem)
32. ~cook/ai1/lectures/applets/pd/ga-axelr.htm (problem)
33. Undergrad/97/97program\_info\_4.html (problem)



#### Altavista results

1. Graduate/mse.html
2. jobs.html

Figure 7. Structural query to find pages on ‘jobs in computer science’.

## 6.4 Combined Searching

This section demonstrates the combined power of structural and text searching using WebSUBDUE. In these experiments, WebSUBDUE uses a combination of text/synonym match and structure match to satisfy user queries.

### 6.4.1 Searching for Hub and Authority Pages on Algorithms

Web search engines typically index several thousand relevant pages on a topic, but the user will only be willing to look at a small number of these sites. Hyperlink information can be used to filter this search by identifying hub and authority pages. A hub page on a topic is a page that has links to many other pages on that topic, or that links to many authorities on a topic. A good authority is a page that is pointed to by many good hubs, while a good hub is a page pointed to by many authorities (Kleinberg, 1998). The HITS algorithm (Chakrabarti, 1999) uses a series of matrix calculations to find good hub and authority pages. WebSUBDUE can discover these hub and authority pages by specifying the hub and authority page structure.

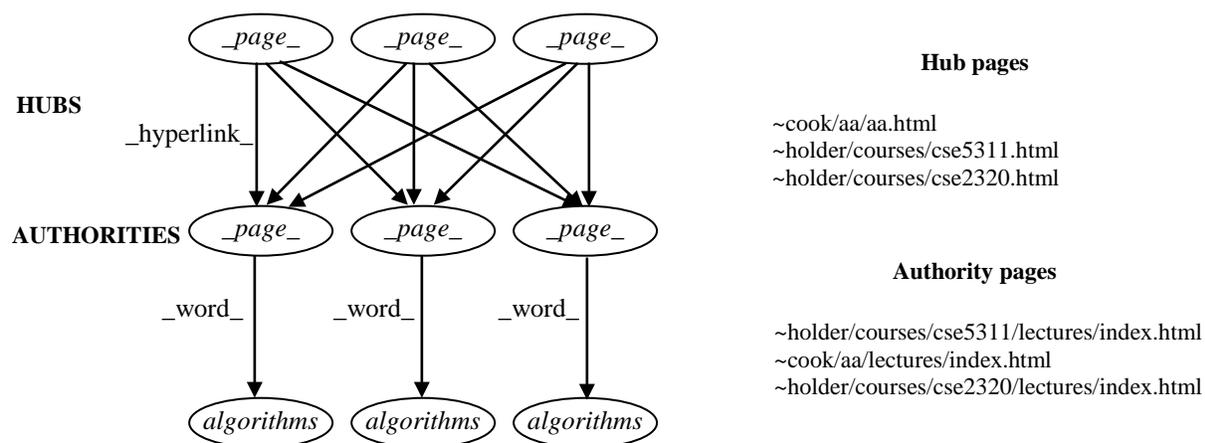


Figure 8. Structural query to find hub and authority pages on 'algorithms'.

A query to search for all hub and authority pages on ‘*algorithms*’ was posed to WebSUBDUE (Altavista cannot process this type of structural query). The substructure input to WebSUBDUE is shown in Figure 8. The query structure represents both an authority page as a page that is pointed to by many (in this case, at least three) hubs, and a hub as a page pointed to by many (in this case, at least three) authorities. WebSUBDUE discovered one instance of the search query, and the corresponding three hub sites and three authority sites are summarized in Figure 8.

#### **6.4.2 Inexact match for hub and authority pages on ‘*algorithms*’**

SUBDUE’s inexact graph match capability can be used to enhance the WebSUBDUE’s web search results. Using an inexact match, web sites can be retrieved which approximate the query structure. The amount of allowable difference is controlled by the user-defined match threshold.

To demonstrate the impact of using an inexact graph match, we repeat the search for hub and authority pages that focus on ‘*algorithms*’, but this time allow an inexact threshold of 0.2. Because all transformations are assigned a unit cost and the size of the query structure is 9 vertices + 12 edges = 21, up to  $0.2 * 21 = 4$  (4.2) transformations are allowed. A total of 13 instances are identified with this query, as opposed to the single instance found using an exact match. Two examples of retrieved instances are shown in Figure 9. The first instance is missing two edges that are provided in the query structure. The second instance is missing three edges and includes an extra node with an associated edge (this combination is treated as a single

transformation). These matches provide indications that the related sites are fairly good hubs and authorities, and the match cost provides a useful ranking measure among the retrieved sites.

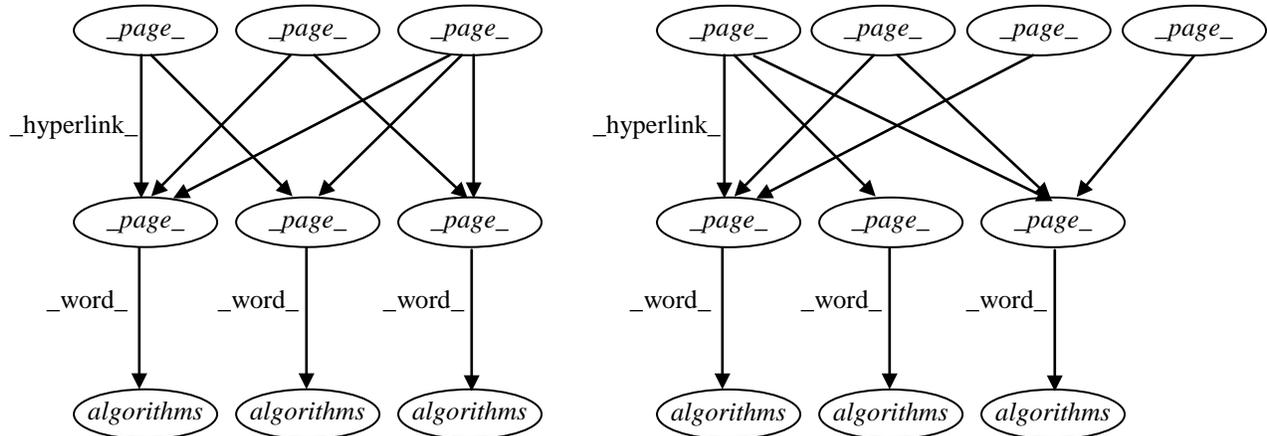


Figure 9. Inexact match results for hub / authority query.

Results of these experiments indicate that WebSUBDUE can successfully perform structural search, text search, synonym search, and combinations of these searches. These experiments also demonstrate the advantage of structural search engines and the inability of the Altavista search engine to perform these functions.

## 7 Conclusions

In this paper, we introduce the concept of a structural search engine and demonstrate that web search can be enhanced by searching not only for text but also for structure that has been created to reflect relationships between individual pages. We have demonstrated how WebSUBDUE can be used successfully as a structural search engine, and how the search can be enhanced using WordNet functions and an inexact graph match algorithm.

To improve the textual search capabilities of WebSUBDUE, additional features of WordNet, such as hyponyms and hypernyms, can be integrated. Distances between words can be calculated using techniques such as marker passing (Lyтинен et al., 2000) or Latent Semantic Analysis (Landauer et al., 1998) and calculated as part of the graph transformation cost. A user interface that facilitates generation of structural queries is also being considered.

Currently, WebSUBDUE must traverse the domain being searched, convert the data to a graph, and use the graph representation as a basis for answering posed queries. We are currently refining this process to couple WebSUBDUE with an existing search engine such as Google. Using the key words, Google will return pages that will structurally searched by WebSUBDUE to satisfy the structural query. This will allow us to test the process on a much larger set of domains. Consideration of additional performance criteria will also be measured.

In addition to using WebSUBDUE as a structural search engine, we are also exploring other methods by which SUBDUE can perform data mining on web data. The unsupervised discovery, concept learning, and hierarchical structural clustering capabilities may provide useful insight into the structure and content of web data.

## REFERENCES

1. The Altavista search engine. <http://www.altavista.com>.
2. S. Brin and L. Page, "The Anatomy of a Large-Scale Hypertextual Web Search Engine", *Computer Networks*, 30(1-7), pages 107-117, 1998.

3. S. Chakrabarti, B.E. Dom, D. Gibson, J. Kleinberg, R. Kumar, P. Raghavan, S. Rajapalan, and A. Tompkins, "Mining the Link Structure of the World Wide Web", *IEEE Computer*, 32(8), pages 60-67, 1999.
4. The Clever Project by IBM. URL : <http://www.almaden.ibm.com/cs/k53/clever.html>.
5. D. J. Cook and L. B. Holder, "Graph-Based Data Mining", *IEEE Intelligent Systems*, Vol. 15, No. 2, pages 32-41, 2000.
6. I. Jonyer, D. J. Cook, and L. B. Holder, "Discovery and Evaluation of Graph-Based Hierarchical Conceptual Clusters", *Journal of Machine Learning Research*, 2, pages 19-43, 2001.
7. J.M. Kleinberg, S.R. Kumar, P. Raghavan, S. Rajagopalan, and A. Tompkins, "The Web as a Graph: Measurements, Models, and Methods", *Proc. International Conference on Combinatorics and Computing*, 1999.
8. J.M. Kleinberg. "Authoritative Sources in a Hyperlinked Environment", *Proc. Ninth ACM-SIAM Symposium on Discrete Algorithms*, 1998.
9. T.K. Landauer, P.W. Foltz, and D. Laham, "Introduction to Latent Semantic Analysis", *Discourse Processes*, 25, pages 259-284, 1998.
10. S.L. Lytinen, N. Tomuro, and T. Repede, "The Use of WordNet Sense Tagging in FAQFinder", *Proc. AAAI Workshop on Artificial Intelligence for Web Search*, 2000.
11. G.A. Miller, R. Beckwith, C. Fellbaum, D. Gross and K.J. Miller, "Introduction to WordNet: An On-line Lexical Database", *International Journal of Lexicography*, 3(4), pages 235-244, 1991.

12. J. Rissanen, *Stochastic Complexity in Statistical Inquiry*, World Scientific Publishing Company, 1989.