Université de Montréal

# Domain Specific Web Search

par

Zheng Zhen

Département d'informatique et de recherche opérationnelle

Faculté des arts et des sciences

Thèse présentée à la Faculté des études supérieures
en vue de l'obtention du grade de Maîtrise ès sciences (M.Sc.)
en informatique

Décembre 2004

Université de Montréal

Faculté des Études Supérieures

Ce mémoire intitulé :

**Domain Specific Web Search**

présenté par:

Zheng Zhen

a été évalué par un jury composé des personnes suivantes :

Guy Lapalme

Président Rapporteur

Jean Vaucher

Directeur de recherche

Jian Yun Nie

Membre du Jury

Mémoire accepté :  11 janvier 2005

# Résumé

Un nouveau modèle de recherche d'information développé dans ce mémoire traite des problèmes de précision et complétude pour les systèmes de recherche courants. De nos jours, la recherche d'information sur le Web est devenue de plus en plus importante dans toutes les applications informatiques. Il y a beaucoup de moteurs de recherche commerciaux puissants tels que Google, Askjeeve, Yahoo et de MSN etc... Tous ces systèmes de recherche semblent pouvoir accepter des requêtes arbitraires, avec un secteur de recherche large. Cependant, leurs résultats manquent de détails dans la plupart des situations. Ce mémoire se concentre sur une situation différente: des recherches récurrentes dans quelques domaines ciblés. Ceci exige des résultats de recherche plus précis et plus complets.

Pour créer un tel système, nous avons utilisé plusieurs progiciels du domaine public et nous présentons un survol des systèmes disponibles.

Comme expérimentation finale, nous avons utilisé notre application de recherche pour trouver des appels d'offres de 5 villes canadiennes choisies aléatoirement. Les résultats prouvent que nos techniques fonctionnent mieux que les moteurs de recherche conventionnels.

**Mot-clé :**

domaine spécifique, recherche, recherche documentaire

# Abstract

A new information-searching model is developed in this thesis to deal with the precision and completeness problems for current searching systems. Nowadays, web searching has become more and more important in all computer applications. There are many powerful commercial search portals such as Google, Askjeeve, Yahoo and MSN; these searching systems support arbitrary queries, with a wide search area, and their results lack in depth in most situations. This thesis focuses on a different situation: recurrent searches in some specified domain. This requires the search results to be more focused.

To create such a system, we utilized several public domain packages and we present a survey of available systems.

As a final experiment, we tested our search application to find call for tenders (CFT) from 5 randomly selected Canadian cities. Results show that our techniques work much better than conventional search engines.

**Keyword:**

domain specific, search, information retrieval

# List of Tables

# List   of   Figures

# Acknowledgement

Special thanks are due to my supervisor, professor Jean Vaucher. His help and advice enabled me to improve and complete this work. With his knowledge and straightforward way he has inspired me to do my best.

I am also gratful to all the professors and staff members in the department of computer science and operations research, for their efficiency, professionalism and sympathy.

Finally, very special thanks goes to my wife for her support and encouragement in making this dream coming true.

# Chapter 1

# Introduction

For the procurement of expensive goods and services, governments and corporations rely on the call-for-tender process: publishing descriptions and conditions for what they need and selecting from submitted tenders. Currently, the trend is to publish Call For Tenders (CFTs) on the Web and many government bodies have set-up major tendering sites which list thousands of CFTs from their various departments. In Canada, MERX (http://www.merx.com) is an example; in the US, FedBizOpps.gov is the single government point-of-entry (GPE) for Federal government procurement opportunities over $25,000; TED (Tenders Electronically Daily) is the equivalent for the European community. Other organizations, such as SourceCan (http://www.sourcecan.com), provide an aggregation service by collecting and indexing documents from other sites.

A quick look an SourceCan (October 2004) shows clearly the importance of this phenomenon: it shows 3400 open bids from Canadian sources, 14,000 from US governments sites, 33,000 European Union public tenders, 2,100 for the World Bank and over 20,000 more opportunities from other sources. However, this abundance of listings hides very real problems when it comes to exploit the data. First, there is no unique standard format for these documents. Most are loosely structured and meant for human interpretation; some documents are quite detailed whereas others are limited to a one-line title. Secondly, product codes – when they are used – may follow different standards. Indexing is therefore partial and approximates. The sheer volume means that effective monitoring of these sites must be done automatically; but the automatic classification of CFTs and identification of pertinent opportunities is difficult.

This situation motivated the launch of a joint university / industry research initiative with researchers from the Universite de Montreal, Nstein and CIRANO. The project, "Matching Business Opportunities on the Internet", is aimed at developing systems which can find and classify CFTs accurately. It has several facets: one deals with the design of collector agents that systematically visit known sites and download new documents. Another,

the biggest, concentrates on the natural language analysis and classification of documents [F. Paradis et al. 2004]. The part which is the topic of this thesis is the design of discovery agents responsible for finding and mapping new CFT sites so that they can be exploited by the collector agents.

In spite of the volume of data concentrated in the major sites listed above, it is important to continuously scour the Web for others. One reason is that many CFTs linked to local governments, large companies, utilities (like Hydro-Quebec) or particular industrial sectors (like Construction, Transportation or Wood products) are published on local sites and not listed on the major ones. This was confirmed by our initial studies. During the last year and a half, the MBOI collector robots visited about 35 sites daily in addition to the major sites. These are especially important because for many companies, CFTs closer to home are easier to take advantage of. A final justification for efficient continuous discovery is that sites are often reorganized and the scripts used by the collector agents no longer work. Over a year, we found that about half the sites changed and we had to rediscover how to extract CFTs from them.

Our discovery agents are basically web spiders, but they target specific kinds of documents and the techniques they require are quite different from those used for commercial search engines (CSEs) like Google or Yahoo. Those systems cater to a heterogeneous clientele and their coverage of the Web is broad but shallow. They need massively parallel search to handle the estimated million pages added to the web every day and must use techniques such as Page Rank [Page et al. 1998] to rate Web pages objectively.

In contrast, we are interested in targeted search with a narrow but deep focus. In particular, we would like to search the invisible web accessible only by complex web navigation involving entering data into forms and which is often not indexed by CSEs. Another difference is that our agents are not primarily interested in documents; they are interested in sites which contain documents. This is because the sites are dynamic and the actual job of document retrieval is left to collector agents. A collector agent doesn't need a list of documents to retrieve; rather it needs a script detailing how to navigate a site and identify new CFTs.

Early in the project, it became evident that intelligent discovery entailed the combination of a multitude of techniques that would have to be adapted to the reality of the

web: what do CFT pages look like? How are CFT sites organized? What data formats are used? Therefore it was decided to develop a prototype to evaluate the aptness of an initial architecture based on a combination of published techniques. Only once each module achieved a certain level of performance could we start more systematic adaptation and evaluation. The prototype was tested on the task of finding the CFTs of a city government given only the city name as starting point.

This thesis presents what we learned from the prototype. It is organized as follows. In Chapter 2, we review some general principles and processes of information retrieval. In Chapter 3, we briefly review current research and list the techniques that should be incorporated in a focused spider. In Chapter 4 and Chapter 5, we describe the architecture of our system and comment on the public domain software tools that were used in its implementation. In Chapter 6, we present some experimental results that show the performance of the system and led to the heuristics used in various modules. Finally, in Chapter 7 gives some conclusions.

# Chapter 2

# Information Retrieval

With the explosion of new sciences and technologies, the problem of information storage and retrieval has been receiving increasing attention. To face the vast amount of information, an accurate and speedy access to information is eagerly awaited. However, difficulties come from either missing relevant information or duplicating useless information. With the rise of computers, there have been a great number of ideas about how to provide a rapid and intelligent retrieval system by computers. In this chapter, we look at information retrieval in classical context. Retrieval from the web is considered in the next chapter.

## 2.1 Information retrieval and Data Base

Database systems are designed to provide efficiently access to large quantity of information. They are not suited for information retrieval (IR) needs.

If we have a database-generated site, such as an e-commerce catalog or a stock exchange for on-line quotations, we might wonder why we need another type of retrieval -- IR. There are some obvious reasons for pursuing a new type of retrieval:

1) Database information is stored in separate fields, but searchers do not want to do searching by choosing fields one by one manually.

2) The response time of database searching can be extremely long, especially for the multiple words searching.

3) The results coming from the database searching cannot be sorted by relevance automatically.

Then, what is so-called IR?

A straightforward definition is given by Lancaster,F.W. [Lancaster,F.W 1968]: "An information retrieval system does not inform (i.e. change the knowledge of) the user on the

subject of his inquiry. It merely informs on the existence (or non-existence) and whereabouts of documents relating to his request."

In IR, we are more interested in finding those items that partially match the requests, and then we will select the best matching ones. The inference and model used in IR are inductive inference and probabilistic modeling. The preferred query language of IR is natural language and the query may be incomplete and vague. Moreover, in IR we search for relevant documents that sometimes don't have exactly matching items. Some small errors in matching sometimes do not affect the performance of the system significantly.

We will discuss and explain how to achieve the above IR features in chapter 5.We'll also introduce an IR system construction tool -- Lucene there. Lucene is a tool that provides some convincing evidence of its efficiency and effectiveness of those IR features [Apache Jakarta Project 2003].

## 2.2   Information retrieval process

The typical information retrieval (IR) system is illustrated in figure 2.1. Just like a library's retrieval system, this system includes a storage of documents and books, an index of this huge storage, and a database of the indexed documents. When users, such as a student or a teacher, require some information, they first need to formulate questions, and then send a query to this system. After receiving a query and matching each items stored in the database to this query, the system retrieves the information and replies to the user. When this system is on-line, it is possible for the user to change his requests frequently for different searching purposes. The whole procedure may be repeated and refined for several times, so as to improve the subsequent retrieval results. This repeated procedure is called relevance feedback. From this procedure, we can see that the basic elements of information retrieval are: documents storage, indexing, index storage, query, matching (searching) and feedback.

Figure 2.1   *Information retrieval (IR) system*

We can abstract the system in a further step, circling three components: input, processor and output.

First, to start with the input, it can link the document storage and users' problems. It should show the best link or the best relationship between the document storage and users' problems, but it usually does not.

Then second, the processor is the core part of the retrieval system for the retrieval processes. The main problem here is how to obtain representations of documents and queries respectively. A lot of computer-based library retrieval systems don't store original documents

or queries, they only store some representations of them. It means that the text of a document is lost, once it has been processed for the purpose of generating its representation. The process may involve structuring the information in some appropriate ways, such as classifying the information. It will also involve performing the actual retrieval function, that is, executing the search strategy in response to a query. In the above diagram, the representation of each document and query has been placed in processor side to emphasize such a fact -- to form these representations is actually an important task during the retrieval process. So the structure of documents is often considered as a part of this process. Finally, the output will come out with a set of references or document numbers.

### 2.2.1 Indexing process

Indexer transfers those documents into one index file which provides for the next searching.



Figure 2.2    *Indexer*

The core of indexer is the document processor, which normally includes 9 steps:

❶ Normalize the document stream to a predefined format.

❷ Break the document stream into desired retrievable units.

❸ Isolate and "metatag" subdocument pieces.

❹ Identify some potential elements in documents for indexing.

❺ Delete "stop words" like: a, an, the, on, at etc.

❻ Stem from terms.

❼ Extract the index entries.

❽ Compute the weights.

❾ Create and update the main inverted file against which the search engine searches to match queries to documents.

Steps ❶,❷, and ❸ are Preprocessing. These first three steps simply standardize the multiple formats encountered when deriving documents from various providers. These steps serve to merge all the data into a single consistent data structure that all the downstream processes can be handled. Step two is important because the pointers stored in the inverted file will enable a system to retrieve units of various sizes — site, page, document, section, paragraph, or sentence.

The fourth Step identifies elements to index. The identification of the potential indexable elements in documents will seriously affect the nature and the quality of the document representation, which the search engine will depend on. To design the system, we must define the word "term". It usually includes the alpha-numeric characters between blank spaces and punctuations. In most way, it should also consider the non-compositional phrases (for example "hot dog"), the multi-word names, and inter-word symbols, such as hyphens and apostrophes, for example, "small business men" versus "small-business men". Each search engine has a series of rules that are used to define a term suitable for indexing. The document processor will tokenize the input streams according to these rules.

Next Step is to delete "stop words". This step mainly contributes to save system resources by eliminating "stop words" from the above process. It may comprise up to 40 percent of text words in a document. A "stop word" list typically contains those words known to convey no useful meanings, such as Articles (*a, the*), Conjunctions (*and, but*), Interjections (*oh, but*), Prepositions (*in, over*), Pronouns (*he, it*), and formats of the "be" verb (*is, are*). To delete stop words, an algorithm compares index term candidates in the documents with a "stop word" list, and eliminates them from the documents for searching.

The sixth step is term stemming. Stemming is a process of removing word suffixes, which achieve two goals: 1). To improve efficiency, stemming reduces a number of unique words in the index, reduces the storage space required for the index, and speeds up the search process. 2). For effectiveness, stemming improves recall by reducing all forms of the word to a basic or stemmed form. For example, if users ask for *analyze*, they probably also want documents which contain *analysis, analyzing, analyzer, analyzes*, and *analyzed*. Therefore, the document processor stems document terms to *analy-* so that documents that include various forms of *analy-* will have equal possibilities to be retrieved. Of course, this stemming

may negatively affect precision: when users query an exact word format, all forms of a stem will be matched, instead of the formatted word in the query form.

Systems may implement either a strong stemming algorithm or a weak stemming algorithm. A strong stemming algorithm will strip off both inflectional suffixes (*-s, -es, -ed*) and derivational suffixes (*-able, -aciousness, -ability*), while a weak stemming algorithm will strip off only the inflectional suffixes (*-s, -es, -ed*).

The seventh Step is to extract the index entries. Having completed steps ❶ through ❻, the document processor extracts the remaining entries from the original document.

The output of step ❼ is to insert and store in an inverted file that lists the index entries and an indication of their positions and frequencies of occurrences. The specific nature of the index entries, however, will vary based on the decision in Step ❹ , which will concern what constitutes an "indexable term." Some document processors will even have phrase recognizers, for example, Named Entity recognizers and Categorizers.

The eighth Step is Term weight assignment. Weights are assigned to each term in the index file. The simplest search engines only assign a binary weight: 1 for presence and 0 for absence. The more sophisticated the search engine is, the more complex the weighting scheme is. To measure the frequency of occurrence of a term in the document creates more sophisticated weighting. Generally, "tf/idf." is an optimal weighting in IR. This algorithm measures the frequency of occurrence of each term within a document. Then it compares that frequency in one document with the frequency of occurrence in the entire database.

Not all terms are good as "discriminators" — which means, not all terms can single out one document from another very well. A simple example is the word "the". This word appears in too many documents to help distinguish one from another. A less obvious example would be the word "antibiotic". For a sport database, when we compare each document to the whole database, the term "antibiotic" could be a good discriminator, and therefore, it should be assigned at a high weight. Conversely, in a health or medicine database, "antibiotic" could be a poor discriminator because it occurs very common. Practically, the TF/IDF weighting scheme assigns higher weights to those terms that really distinguish one document from the others.

The last step is to create index. The index or inverted file is the internal data structure that stores the index information and will be searched by each query. Inverted files can be a simple listing of alpha-numeric sequence in a group of documents/pages indexed, with an overall identification numbers of the documents in which the sequence occurs to a complex list of entries, the tf/idf weights, and pointers. The more complete the information in the index is, the better the search results are.

We'll give a concrete example in chapter 4, in which Lucene will be introduced as a very good indexing tool for constructing indexer and search engine.

## 2.2.2 Query process

Query processing has six steps at different level. It shares many steps with document processing. The whole processing is shown in the figure 2.3 below.

**Level I:**

  Query term(s) ---> tokenizing----->parsing--->   ——————————>

**Level II:**

    deleting stop words ---> stemming words --->   ——————————>

**Level III:**

    expanding   term(s) ---> computing weight ---->   ——————>

Figure 2.3   *Query Processor*

Level I:

**Tokenizing:** As soon as a user inputs a query, the search engine will tokenize the query stream, i.e., break it down into understandable segments. Usually a token is defined as an alpha-numeric string that occurs between white spaces and/or punctuations.

**Parsing:** Because users may utilize some special operators in their query, including Boolean, adjacency, or proximity operators, the system needs first to parse the query into query terms and operators. These operators may occur in the form of reserved punctuations (e.g., quotation marks) or reserved terms in specialized formats (e.g., AND, OR). At this point, a search engine may take the list of query terms and search them with the inverted file. In fact, this is the point at which most of the available public search engines perform the search.

Level II:

**Stop list and stemming:** Some search engines will go the stop-list and stem the query, similar to the processes described above in the Document Processor section. The stop list might also contain those words from frequently occurring query/phrases. However, because most of the available public search engines encourage short queries, the engines may drop these two steps. Some particular search engines can create a query representation combining with its matching function and perform the search against the inverted file at this point.

Level III:

**Query expansion**: Because most of the users usually search some single statement of information with a query, it is highly probable that we can use synonyms as a substitute, not be limited on the exact query terms, to search in the documents for meet users' needs better. Therefore, more sophisticated systems may develop the query into all possible synonymous terms --- even broader and narrower terms.

**Query term weighting** (assuming more than one query term): The final step in query processing is to compute weights for each term in the query. Sometimes those users control this step by simply indicating either pre-fixed/required each term's weights or which term in

the query must be considered in each retrieved document to ensure relevance. After this final step, the expanded, weighted query is searched against the inverted file of documents.

Although a system can choose any one of the 3 kinds of query processors to match the query to the inverted file, numerous steps and documents make this process more expensive for processing with the computational resources and responsiveness. However, the longer the waiting time is, the higher the quality of results will be. Thus, there is a tradeoff between time and quality. Available public search engines usually choose speed, not quality, having too many documents for searching. Sometimes for some special case, users may occasionally obtain better quality, for example, in domain specific web search.

## 2.3    Areas of search

Effort to improve information retrieval (IR) can be subdivided in content analysis, information structure and evaluation. Briefly, the content analysis is concerned with describing the contents of documents in a form suitable for computer processing. The information structure is concerned with exploiting relationships between documents to improve the efficiency and effectiveness of retrieval strategies, and the third one -- evaluation is concerned with the measurement of the effectiveness of the retrieval.

Usually we use the vocabulary-usage frequency in the document text to determine which word is sufficiently significant to represent or characterize this document in the computer. Thus, a list of what might be called 'keywords' derive form each document. In addition, the frequency of the occurrence of these words in the text body could also be used to indicate a degree of significance. This provides a simple weighting scheme for the 'keywords' in each list and makes a representation in the form of a 'keyword description'. The use of statistical skill about distributions of words in documents reveals statistical associations between keywords. These associations can provide a base for the construction of a thesaurus – an aid to information retrieval (IR). The frequency that any two keywords occur together in the same document is called co-occurrence frequency. It can be used to measure of the association between keywords [Baroni et al. 2002].

Specifically, the information structure is a logical organization of information. The development of information structure should empasize on the logical structure of the information (documents, files). Some earlier experiments for those document retrieval

systems usually adopted a series of file organizations. Up to now, it is more often to use the inverted file or clustered files for on-line retrieval.

Evaluation of retrieval systems historically has been proved extremely difficult. As early as 1966, Cleverdon listed six main measurable quantities[Cleverdon et al. 1966]:

(1) The *coverage* of the collection, that is, the extent of the system included relevant matters;

(2) The *time la*g, that is, the average interval between the spending time of the search request and finding out an answer;

(3) The form of output *presentation*;

(4) The *effort* involved on the part of the user in obtaining answers to his search requests;

(5) The *recall* of the system, that is, the proportion of relevant material actually retrieved in answer to a search request;

(6) The *precision* of the system, that is, the proportion of retrieved material that is actually relevant.

There are two factors: the recall and the precision that are used to measure the effectiveness of the retrieval system. In other words, they are the rulers of the system abilities to retrieve relevant documents and to hold back the non-relevant ones at the same time. The more effective the system is, the more it will satisfy the user. It is assumed that the precision and the recall are sufficient for the measurement of effectiveness.

Nowadays, all the IR sub-categories --- content analysis, information structures and evaluation, have been already applied in computerized IR systems, in on-line IR systems, and even in the web searching systems. For example, we can give a keyword to calculate those documents' or web pages' tf / idf, so as to determine the most relevant files or pages. In order to make the information structure for those documents or web pages, computer should get the logical structure of the documents/web pages first [Carchiolo et al. 2003].

Although, web pages have a lot of differences with traditional documents, the basic principles of IR can also apply in the web. Chapter 3 will mainly discuss a situation of web search --- IR on the web.

# Chapter 3

# Survey of Web Search

The web search problem actually is a problem of IR on the web. The whole web search process basically includes three parts: search engine, indexer and crawler [Sherman et al. 2003]. Different users may pay attention to different parts. One can complete its own page-ranking algorithm, while others may add more function to their crawler parts [Brin et al. 1998]. Referring to the figure 3.1, following the order from top to bottom, we'll go into details for each part of web search one by one.

Generally, as an important concept, a crawler-based searching system has three basic elements [Crimmins 2002]. One is that the crawler downloads the new web pages from the web continuously. The second is that the indexer gathers the words from those documents, whether they are HTML pages, local files or database records, and the indexer picks up those words into one index for a fast information retrieval. The third element is the search engine itself, which accepts the queries, locates the relevant pages in the index and formats the results in an HTML page (see figure 3.1). Indexer and Search engines are heavy-duty server programs, thus it is required that processors be fast for the processing, hard-disk space be significant for the index and a great deal of bandwidth must be available for responding to many simultaneous searching requests. An exact configuration depends on the numbers of pages, but most of the search engines require Intel Pentium or Sun Solaris processors, Microsoft Windows NT/2000 or Unix, and at least a T1 line.

For those specific domain search applications [Oyama et al. 2004], an optimal solution can make a crawler more powerful and smarter by making it able to follow extended links intelligently while it is crawling on the fly.

Figure 3.1 W*eb search system*

The term "specific domain", it does not only mean a domain specified before searching, but also mean a domain specified with certain arbitrary. For a domain specific searching, we can easily apply heuristic search. Thus, the crawler is able to crawl with a focused crawling to reduce the search space and the storage of downloaded pages and indexed files. With the help of the basic crawler-based search engine technologies, we develop a domain specific information searching system on the base of the huge web resources.

## 3.1 Crawler

A web crawler is a program which automatically traverses the web by downloading documents and following links from page to page [Koster 1999]. In the input side, the features of the downloaded web pages can highly influence the behavior of crawler. The workflow of crawler looks like the following chart.



Figure 3.2    *Crawler*

After the crawler gets an initial URL, it will visit this URL and download this web page. It extracts URLs from downloaded pages, sends extracted URLs to the scheduler, and continuously begins the next download loop. The scheduler will be responsible for choosing a URL from those extracted URLs for guiding the crawler on the web in its further step. This is the basic continuous process. In reality, we often complement its functions with intelligent crawling.

Among all the web search systems, some crawlers are indiscriminate crawlers, which means that they will download web pages without any choosing, and any new pages are possible to be downloaded. But there are still some others that can follow some rules when crawling [Mladenic 1999]. Some interesting methods proposed are those of *Fish Search* [Bra 1994] and *focused crawling* [Chakrabarti et al. 1999]. The essential idea in focused crawling is that there is a short range topical locality on the web. This locality may be used in order to

design effective techniques for resource discovery by starting at a few well chosen points and maintaining the crawler within the ranges of these known topics. These crawlers only download new pages with an interesting topic. Hub pages (web pages containing links to a large number of pages on the same topic) and authority pages (documents whose content corresponds to a given topic) may be used for the purpose of such crawling. Generally, the hubs on a subject are likely to point to authorities and vice-versa. In addition, structural properties of world wide web like linkage locality (web pages on a given topic are more likely to link to those on the same topic) and sibling locality (if a web page points to certain web pages on a given topic then it is likely to point to other pages on the same topic) can be applied in crawling [Kleinberg 1998].

Besides, to implement a focus crawling, a crawler must create some learning modes, in other words, it should have some intelligence [Aggarwal et al. 2001]. There are 6 learning modes proposed by Aggarwal:

1) Probabilistic mode for priorities

2) Content based learning

3) URL token based learning

4) Link based learning

5) Sibling based learning

6) Combining the preferences

Dr. J.H. Cho also gave some more complex methods through URL analysis in 1998[Cho et al. 1998]. The advantage of URL analysis is that it is fast and light-weight (don't rely on much storage) so that a crawler can apply it on the fly.

Since some learning modes are based on much calculation, only part of them are proper for creating a light-weight crawler, such as mode 1), 3), 4) and mode 5).

When wandering on the web, crawler can explore new URL with different algorithm: breadth-first, depth-first or best-first. Whatever breadth-first or depth-first, it must go through all the web pages (nodes) to finish a complete solution. However, it is obviously impossible to crawl all the 500 billions pages on the web for each searching, even for Google, which can only crawl and store 5 billions web pages. To try a heuristic function, the best-first search algorithm is an optimal solution. If a crawler can keep users' keyword in mind, explore new web pages along a route close to the users' interest, it will come out with what we want. Web searching means a focusing crawling that is able to crawl specific topics on portions of the web quickly without having to explore all web pages. To realize this selective crawling, a crawler should be able to select the next node based on some criteria.

## 3.2 Search engine

Search engine working process looks like the following chart. It mainly includes query processing, search and match process, and page ranking[ Elizabeth Liddy 2001].
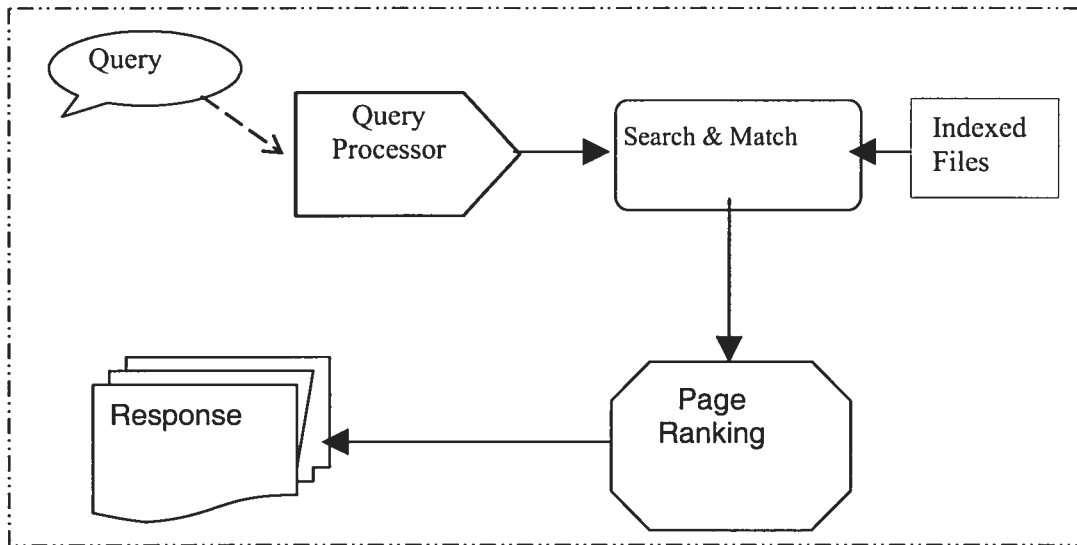


Figure 3.3   *Search engine*

### 3.2.1 Search and match function

Search and match function depends on the information retrieval mode of the searching system. About the searching the inverted file for documents to meet the query requirements, please refer to the above "matching" process. While the required computational processing for simple, un-weighted, non-Boolean query matching is far simpler than a weighted, Boolean model. The search and match function is related to each other clearly and easily, as the search system information retrieval model, the query processor model and the document processor model are determined above.

### 3.2.2 Page ranking

To determine the degree to which those subsets of documents or pages match the query requirements well or not, a similarity score can be computed between the query and each document / page, based on the scoring algorithm used by the system. Scoring algorithms rankings are based on the presence of query terms, term frequency, tf/idf, Boolean logic, and query term weights. Some search engines use scoring algorithms not only based on document contents, but also based on the relationship among documents or past retrieval history of documents /pages. After computing the similarity of each document in the subset of documents, the system provides an arranged list to the user. Some systems' orders also depend on query weighting mechanisms. However, after the search engine determines ranks, the user can simply click and follow hyperlink to the selected document /page. More sophisticated systems, as figure 2.1 shows, will allow the user to provide some relevance feedback or to modify the query based on the results received.

### 3.3 Performances of a web search system

General speaking, web search has these basic characters:

1) web search has to face heterogeneous storage with various format, language etc;

2) web pages are never indexed on the web;

3) web search face an open, dynamic, and variable situation. The content of a specific web page may be changed periodically, or even totally disappear.

4) all the pages on the web are connected only by hyperlink within pages. They are absolutely in a huge network.

In order to discuss the web search, we modify the Figure 3.1 to Figure 3.4. To improve performances of a web search system, a crawler of a search system can extend its function for different designs. We will discuss this more in the following chapters. The processor is a pure IR application, basically the same as the traditional IR system. For the output layer, there is an additional page ranking procedure. For any web search, matching is partial match or best match; given a proper query, it is quite easy to find numerous responses to meet users' requests on the web; any web search is an open, dynamic process, i.e., users can browse webs following the search results, and any bias or error may lead users to browse outside of his concentration. Therefore, an important step of the page ranking is to ensure that users get partial but best match corresponding his/her query. This is also an example of evaluation – the most difficult one of the three IR subcategories.
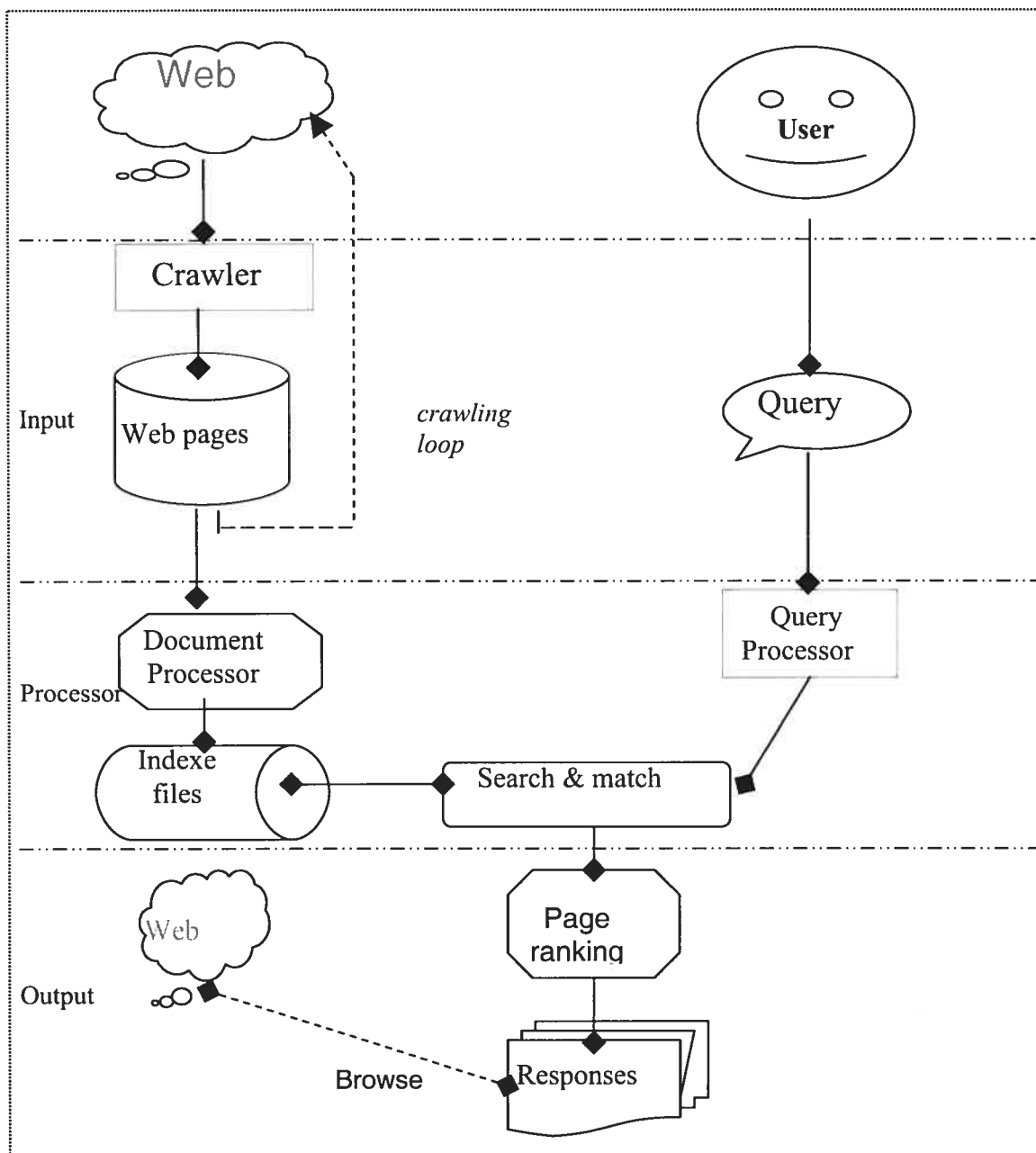
Figure 3.4 *Layers of a web search system*

# Chapter 4
# **Domain specific web search**

Web search systems are characterized by their strategies. Our crawler combines several techniques. It starts with a small set of interesting pages obtained by interrogating commercial search engines with some judiciously chosen keywords. The crawler then proceeds to explore the neighbouring web by following links from the initial page set. However, the real problem a targeted crawler must address is the analysis of page content to determine the pertinence of each page and its links. Pages on commercial sites generally follow a standard layout with the interesting data in the center surrounded by generic header and footer information combined with navigation sidebars. In order to intelligently classify pages, the crawler must first do a *structural analysis* to identify the various functional blocks on a page (header, footer, etc...) and extract its information core; then it can go on to a *lexical analysis* of this core to determine whether the page is relevant to the area of interest.

## **4.1    Proposed architecture**

Figure 4.1 shows the general architecture of the system. The initial query will include a city name where from a menu of city plus user selected terms and the results will be used as seeds (initial URLs), the number of seeds is set as a parameter. We start with a small set of web page URLs obtained from commercial search engines such as Google or Yahoo. Then we proceed to explore links from those pages. In order to explore efficiently, we must evaluate the pertinence both of pages and links, and discard irrelevant data. This is underlined by the several places where data is discarded.

The diagram shows that various techniques come into play. After a quick initial evaluation, such as file type etc., the crawler begins to do a real job: the classification of pages. In what follows we consider the classification of pages, tunneling, structure analysis and content analysis.
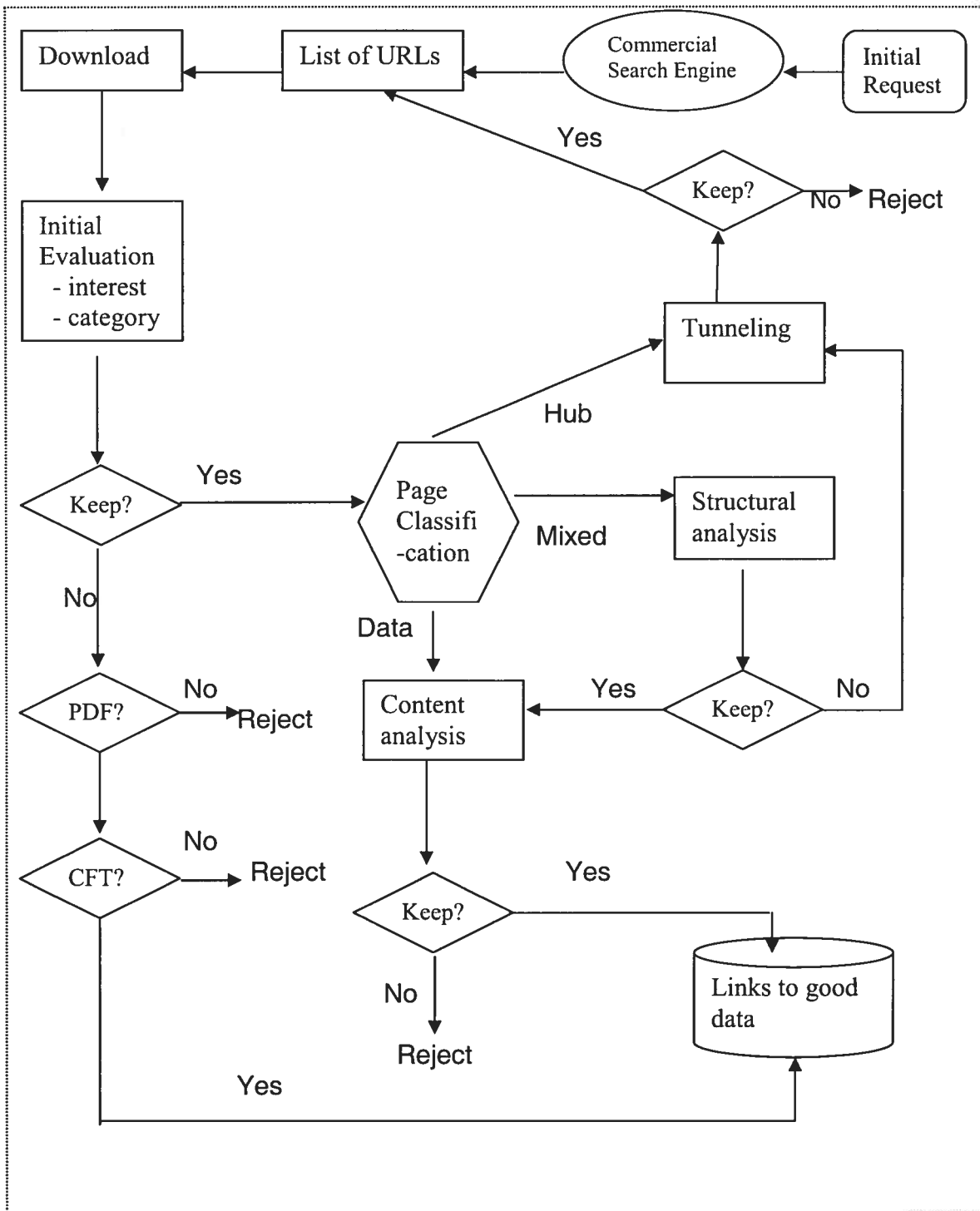
Figure 4.1    Domain specific searching system

## 4.2    Classification of web pages

To focus on a specific domain, the crawler should pick out links from downloaded pages judiciously. Besides, the particular treatment varies with the web page being crawled. For the different type of web pages, the selection process shouldn't be the same. For some web pages, the crawler may pick out links and crawls further; for some others, it may store or discard them so as to terminate the current crawling process. Normally speaking there are three types of web pages: hub, data and mixed pages. A hub page is a web page on which most of its content are links (see Figure 4.2). A data page is a page on which most of the content is pure text (see Figure 4.3). A mixed page is a web page with a mixture of links and text (see Figure 4.4 type ❶ & Figure 4.5 type❷).



Figure 4.2    *Typical "Hub" page*

Figure 4.3    *Example of a data page*

Figure 4.4    *A "Mixed" web page (type ❶)*

Figure 4.5    *A "Mixed" web page (type ❷)*

Type ❶ pages can be treated as a data page, whereas type ❷ pages are rather like hub pages. This is because each piece of text in type ❷ page is a brief introduction for the nearest

link. Although it is not easily detected by a crawler, the ratio of links/strings is useful in differentiating between type ❶ and type ❷ pages.

## 4.3    Tunneling

When a crawler explores a hub page, it will analyze anchor words of all the links within the page so as to explore some links further with close meaning of the query term. Sometimes, however, the anchor words in relevant links do not match any in the query.

For example, if we search "tender" on the following hub page (see the figure 4.6),



Figure 4.6    *Tunneling*

we can't directly find any anchor word matching the term "tender", but we should explore some pertinent links like "Shopping" or "Business Finder" first rather than exploring some other irrelevant links or even discard this hub page.

The phenomenon that a crawler has to pass through irrelevant pages whose content does not match exactly on query in order to reach useful pages is called tunneling (Figure 4.6). Tunneling often occurs in hub pages. In this case, the crawler needs to adjust whatever widening the matching criteria or extending the query term to tackle the tunneling [Bergmark et al. 2002].

Another example of tunneling is when the hub pages provide "search" function. The crawler can auto-fill "search" form with the query term, send a requirement, and get a response page to avoid tunneling.

## 4.4    Web page structural analysis

Text data, hyperlinks, images are organically placed inside web pages. An explicit visual structure of a web page can help us to understand it semantically. Such structure can be partly used in web pages' classification since an analysis of the structure may reveal logical relevancy between sections in the web page.

Web page structural analysis combines both structural and semantics information of a page in order to get to a logical page schema.

### 4.4.1    Web page structure

Semantic information of a web page often needs to be organized with various sections so as to represent author's view or linkage relation. This structural requirement reflects in both coding html tags and arranging logical content [Carchiolo et al. 2003].

#### (1) HTML structure

Tags are used to organize blocks of text representing the same information. They also control graphical appearance of a web page. Tags like <table>, <p>, <hr>, <map>, <a>, <frame>, <td>, <tr>, <option>, <area>, <li> often present a logical section or a subsection in a web page. Besides some properties play an important role in a semantic partition of a web page, such as <td bgcolor =>, <front size=>, <td width=>.

As mentioned above, although a lot of html tags can represent relative independent logical meaning in a web page, the <table> tag is better for analyzing. In a correctly

formatted web page, various logical sections are located between opening tags (i.e. <table>) and closing tags (i.e. </table>). All the <table> tags of a web page can be constructed in a tree structure. Analyzing this tree can helps in recognizing logical structure of a web page.

### (2) Logical structure

Often the HTML structure reflects the logical organization of a web page. Just as written text follows certain conventions for example starting with a title page, then an introduction, a body and a conclusion followed by references; the constituent parts of a web page often fall in certain typical categories.

Generally a web page can be partitioned as several logical sections:

1) a *document information section:* metadata concerning the document, as document type specified in <!doctype>, referring DTD, author etc.

2) a *logical heading:* the name, title or the logo of the site, usually placed at the beginning of the document. Denoted by the tag <head>.

3) a *logical footer:* generally placed at the bottom of page and whose content is often a mail link to webmasters or text with copyright, privacy etc.

4) *logical links:* a set of HTML tags representing a link towards another page. Denoted by the tag <a href>.

5) an *index:* a group of logical links having similar characteristics and graphical appearance. Denoted by the tag <a href>.

6) *logical data:* containing raw data --- information core of the page

7) *interactive sections:* forms where users interact with the page; this section generally includes codes such as cgi, php, applets or script, whose meaning is hard to understand. This section often contains <form>, <input> or <select>. It is easily to be identified or automatic filled by a piece of program.

### 4.4.2   Analyzing structural information

Not all the web pages require structural analysis; only those pages with a relative complex structure, similar to Figure 4.4 or Figure 4.7. These web pages normally have the logical section 1) to 6) or 1) to 7) mentioned above. Since most logical data sections mainly contain raw data, as a simplified example of using structural analysis, we extract the

logical data section from these web pages for lexical analysis. To do the structural analysis, we decompose a web page owning a complex structure and being composed with a lot of tables to a set of sub-tables, analyze each of them, and then pick out the proper one --- the logical data section of the web page. When a crawler explores a mixed page, it will analyze the structure and content of that page. The crawler will extract the main content of that page and then do content analysis. For example, see Figure 4.7, Figure 4.8.



Figure 4.7    *"Mixed" page before data extraction*

The analyzing process decomposes a mixed page with a lot of imbricated tables, then counts links and strings of each table, calculates a ratio (number of links/ number of strings) and a value for each table.

The value degree = "tender" occurring frequency / ratio, if ratio not equal zero; otherwise indicate an experimental value.

A table with the biggest value degree is extracted out.

Figure 4.8    *"Mixed" page after extraction*

We call the main content part extract from the "mixed" web page an information core. If the information core matching a primary requirement of querying, for example, it contains the term "tender", then the page where the information core locates will be sent to do content analysis for further evaluation. Otherwise the page will be classified as a hub page (see figure 4.5).

## 4.5    Web page content lexical analysis

Lexical analysis is to classify web pages according to their content. Lexical analysis will transform the web pages into a representation suitable for classification task. It includes the following steps: collecting a training set, extracting features of it, calculating similarity, classifying the web page automatically.

### 4.5.1 Collecting a training set

For the domain specific search, a training set relate to a specific domain. To collect the training set, a set of related web pages (HTML files) should be downloaded firstly. After removing HTML tags and stop-words from those HTML files, (if necessary, perform word stemming) the set of web pages is transformed to a collection of documents. In the "call for tender" domain, different sites have various written formats, hence collecting a set of tenders from several cities is necessary to extract representation of tenders. The examples of various tenders are shown as following (Figure 4.9, 4.10, 4.11 & 4.12).



Figure 4.9    *City of Victoria, BC*

Purchasing & Tenders

Address http://www.city.red-deer.ab.ca/Connecting+with+Your+City/City+Services+and+Departments/Treasury/Purchasing+and+Tenders/default.htm

## Purchasing & Tenders

The Purchasing Department is responsible for the acquisition of all materials, equipment and services required by The City, and the disposal of all surplus materials and equipment. Purchasing also provides information on products, suppliers/vendors, specifications, etc. to operating departments. The City of Red Deer requires that purchasing be done on a competitive basis.

The services provided are:

- Acquisition of materials, equipment and services required by The City
- Disposal of all surplus materials and equipment
- Control and maintenance of adequate levels of inventory

Please contact the Purchasing Agent for more information.

### City Tendering

The City of Red Deer invites prospective suppliers of various goods and/or services wishing to Tender on City requirements, to request a Bidder's Application form from The City of Red Deer Purchasing Department.

If you require more information on the tendering process, please refer to Purchasing Tendering Procedures.

| | |
|---|---|
| Office Hours: | Monday to Friday, 8:00 a.m. – 4:30 p.m. |
| E-mail: | purchasing@reddeer.ca |
| Phone: | (403) 342-8273 |
| Fax: | (403) 341-6960 |
| Street Address: | Fourth Floor, City Hall |
| | 4914 - 48 Avenue |
| Mailing Address: | Purchasing |
| | P.O. Box 5008 |

Figure 4.10   *City of Red deer, Alberta*

Address http://www.city.london.on.ca/purchasing/manpg.htm

Constitution and Terms of Reference - January 20, 2003

### Active Tenders/Quotes/RFPs

| File # | File Name | Plan Takers List | Registration Form | Closing Date | Status |
|---|---|---|---|---|---|
| T04-30  Addendum #1  Addendum #2  Addendum #3 | Tender 04-30 Automatic Irrigation Upgrades at Fanshawe Golf Course | List | add me to this mailing list | Wednesday, September 1, 2004 | Open |
| RFP04-15 | Request for Proposal 04-15 Two Megawatt Landfill Gas Power Generation | List | add me to this mailing list | Friday, August 13, 2004 | Open |
| | Request for Proposal | | add me to | Friday, | Open |

Figure 4.11 *City of London, Ontario*

Figure 4.12 *Town of Orangeville, Ontario*

### 4.5.2 Feature extraction

From the Figure 4.9 to 4.12, we can see various formats of tenders. After collecting a set of cities' tenders, we need to extract their common features. The most commonly used document representation is the so-called vector space model. In this model, documents are represented by vectors of words. A collection of documents which is represented by a word-by-document matrix $\mathbf{A}$, where each entry represents the occurrences of a word in a document, i.e., $\mathbf{A} = (\mathbf{a}_{ik})$, where $\mathbf{a}_{ik}$ is the weight of word $i$ in document $k$. These techniques are well known in the IR fields, for a further overview, see [Aas et al. 1999].

There are several ways of determining the weight $\mathbf{a}_{ik}$, most of the approaches are based on two empirical observations regarding text:

1) the more times a word occurs in a document, the more relevant it is to the topic of the document.

2) the more times the word occurs throughout all documents in the collection, the more poorly it discriminates between documents.

Suppose we have the statistic matrix A0 derive from the collection:

$$
A0 = \begin{array}{l} \text{word}_1 \\ \text{word}_2 \\ \dots \\ \text{word}_M \end{array} \begin{pmatrix} f_{11}, f_{12}, \dots f_{1N} \\ f_{21}, f_{22}, \dots f_{2N} \\ \dots \\ f_{M1}, f_{M2}, \dots f_{MN} \end{pmatrix}
$$

let $f_{ik}$ be the frequency of word $i$ in document $k$, M number of words in collection N number of documents in collection, $n_i$ is the total number of times word $i$ occurs in the whole collection.

If adopting the tf/idf weighting, $a_{ik} = f_{ik} * log(N / n_i)$, thus we got the matrix **A:**

$$
A = (a_{ik}) = \begin{array}{l} \text{word}_1 \\ \text{word}_2 \\ \dots \\ \text{word}_M \end{array} \begin{pmatrix} a_{11}, a_{12}, \dots a_{1N} \\ a_{21}, a_{22}, \dots a_{2N} \\ \dots \\ a_{M1}, a_{M2}, \dots a_{MN} \end{pmatrix}
$$

There may be a lot of words in the collection. Hence, the dimensionality of the feature space (matrix **A**) is very high. We can remove non-informative words from documents through the feature selection so as to improve effectiveness and reduce computational complexity.

The document frequency for a word is the number of documents in which the word occurs. The document frequency threshold can be chosen as the feature selection method to reduce the original feature set.

Applying the main idea above mentioned, we collect a set of tenders in 10 cities and do lexical analysis as following:

1) http://www.oshawa.ca/cit_hall/tenders.asp
2) http://www.city.brampton.on.ca/purchasing/rfp_2004-057.tml
3) http://www.city.brampton.on.ca/purchasing/rfp_2004-040.tml
4) http://www.city.cambridge.on.ca/cs_corporate/purchasing_tenders_list.php
5) http://www.city.kitchener.on.ca/tenders/tender.asp
6) http://www.city.london.on.ca/Purchasing/mainpg.htm
7) http://www.city.niagarafalls.on.ca/cityhall/qtenders.html
8) http://www.region.peel.on.ca/finance/purchasing/biddocs/index.htm
9) http://www.city.greatersudbury.on.ca/pubapps/tenders/
10) http://www.city.toronto.on.ca/tenders/tenders_to.htm
11)http://www.region.waterloo.on.ca/web/region.nsf/97dfc347666efede85256e590071a3d4/4a3823e3515da81e85256f11004b50c4!OpenDocument

Words whose occurring frequency larger than 10 are:

| Frequency | Word | Frequency | Word |
|---|---|---|---|
| 133 | click | 35 | mailing |
| 83 | tender | 34 | date |
| 85 | list | 34 | proposal |
| 76 | acrobat | 33 | opening |
| 75 | format | 31 | document |
| 74 | pending | 30 | purchasing |
| 63 | ad | 26 | city |
| 62 | services | 24 | supply |
| 40 | contract | 24 | tenders |
| 40 | addendum | 23 | bid |
| 40 | request | 23 | issued |
| 38 | september | 22 | sep |
| 35 | contact | 22 | peel |
| 35 | add | 20 | information |
| 34 | august | 20 | closed |

| Frequency | Word | Frequency | Word |
|---|---|---|---|
| 20 | works | 12 | proposals |
| 19 | wednesday | 12 | division |
| 19 | aug | 12 | road |
| 19 | extended | 12 | bids |
| 17 | rfp | 12 | site |
| 16 | quotation | 12 | Brampton |
| 16 | plan | 12 | hall |
| 16 | council | 12 | perform |
| 16 | emergency | 12 | takers |
| 16 | construction | 12 | links |
| 15 | street | 12 | view |
| 14 | deposit | 11 | public |
| 14 | form | 11 | pdf |
| 14 | friday | 11 | opportunities |
| 14 | replacement | 11 | regional |
| 14 | awarded | 11 | removal |
| 14 | search | 11 | community |
| 13 | download | 11 | register |
| 13 | open | 10 | west |
| 13 | transportation | 10 | centre |
| 13 | closing | 10 | delivery |
| 13 | buyer | 10 | full |
| 13 | home | 10 | floor |
| 12 | sept | 10 | fill |

*Table 4.1    Word Frequency*

From above results we can see that CFT content lexical analysis is important but importance of a word doesn't match exactly with its occurring frequency in a CFT document. For example, the first 8 high-frequency words are: click (133), list (85), tender (83), acrobat (76), format (75), pending (74), ad (63), services (62). If we evaluate a document with these nine words, we can hardly say the document is a CFT or not. On the contrary, some words play a very important role in a CFT document, though its occurring frequency is not too high. For example, occurring frequency of word "closing" has only thirteen. In fact, the word "closing" appears all in content of URL 1) to URL10) as an important sub-title "closing

date". Therefore CFT feature words should consider both CFT lexical analysis and CFT format (structure analysis). Different word has different tender format weight in a CFT document. From format point of view, a CFT normally has following features:

1) *token words*: tender(83), tenders(24), bid(23), bids(12), rfp(17), quotation (16), solicitation(2), proposal(12) etc.

2) *time tag*: closing(13), closed(20), opening(33), open(13), pending(74),date(34), issued(23) etc.

3) *action words*: request(40), contract(40),contact(35), purchasing(30), delivery(10) etc.

4) *requirements*: document(31), addendum(40), information (20)

5) *address words*: add. (35), city (26), street(15), road(12), mailing(35)

6) *value words*: deposit(14), awarded(13), buyer(13)

Although the basic CFT format is the same, the words above may vary with different area or different speciality. The most typical features of CFT are token words and time tag. As long as we combine these two features, the precision of query increase a lot. For example, suppose we want to search CFT in the city of Peel. A quick experiment is to query Google with terms "peel tender", "peel tender closing" and "peel tender request":

Query "peel tender":

1) www.taunton.com/finecooking/pages/c00161.asp   talking about cooking eggplant

2) www.baking911.com/howto_blanch.htm talking about cooking

3) home.comcast.net/~iasmin/ mkcc/MKCCfiles/candiedcitrus**peel**.html talking about cooking

4) www.extension.umn.edu/distribution/ nutrition/components/0555%5Bt01%5D.html cooking

5) www.nyapplecountry.com/edtastehowto.htm cooking

6) www.evergreen.edu/biophysics/ technotes/home/dry_food.htm    cooking Asparagus

7)   www.healthgoods.com/Education/Nutrition_Information/   Drying_Food/pretreating_vegetables.htm cooking mushrooms

8)   www.cooks.com/rec/search/0,1-00,baked_   sweet_potatoes_cinnamon,FF.html   cooking potatoes

9) www.wholehealthmd.com/refshelf/ foods_view/1,1523,29,00.html cooking Parsnips

10) www.geocities.com/Heartland/Acres/1012/marpre.html cooking

Query "peel tender closing":

1) www.region.**peel**.on.ca/housing/**peel**_living/ Peel_Living_Minutes/1980s/1980/plmin19801030.htm
   a CFT in region Peel

2) www.spc.lk/seventhTender2.html a CFT in Sri Lanka

3) purchasing.**peel**schools.org/documents/2601t.pdf   a CFT in region Peel

4) purchasing.**peel**schools.org/documents/2469t.pdf   a CFT in region Peel

5) www.nationalparks.nsw.gov.au/npws.nsf/Content/
Tender+advert+roads+fire+trial+maintenance+and+fire+suppression a CFT in region Peel

6) www.pakwatan.com/main/businessnew/ **tender**international.php3   a CFT

7) www.pakwatan.com/main/businessnew/**tender**int.php3 a CFT

8) https://www.raqsb.mto.gov.on.ca/raqs_contractor/
raqscont.nsf/0/25EF082A2079207D85256A1B00626DAF?Opendocument a CFT in region Peel

9) www.city.pg.bc.ca/city_services/supply/2003_documents/T03-
17_Construction_Services_Roof_Repairs.pdf a CFT in region Peel

10) www.dpcdsb.org/purchasing/program.cfm      a CFT


Query "peel tender request":

1) www.region.**peel**.on.ca/housing/**peel**_living/
Peel_Living_Minutes/1980s/1980/plmin19801030.htm          Bussiness meeting in region Peel

2) purchasing.**peel**schools.org/documents/2589t.pdf      a CFT in region Peel

3) www.ipc.on.ca/scripts/index_.asp?action=31& P_ID=15033&N_ID=1&PT_ID=759&U_ID=0
   An appeal in Peel

4) www.thriftyfun.com/readers**request**/tf876797_rea.html Cooking

5) www.extension.umn.edu/distribution/ nutrition/components/0555%5Bt01%5D.html
   Vegetables

6) starbulletin.com/97/07/30/features/**request**.html   Fruit

7) www.region.**peel**.on.ca/housing/**peel**_living/
Peel_Living_Minutes/1990s/1998/plmin19980813.htm          A Meeting in region Peel

8) www.region.**peel**.on.ca/housing/**peel**_living/
Peel_Living_Minutes/1980s/1984/plmin19841204.htm          A Meeting in region Peel

9) www.ext.nodak.edu/extpubs/yf/foods/he187w.htm Vegetables

10) 01wholesale.com/ import-export-directory--201-Import-export---wood.html

International    trade

From this quick experiment, we can see that there is huge weight difference between CFT feature words. We adopt three CFT feature lists in our domain specific system to evaluate a web page. The first is a CFT token words list, including "tender", "bid", "RFP", "RFQ", "solicitation". The second is a CFT time tag list, including "closing", "open", "closed", "issued", "date". And the last is a hybrid CFT list, including "document", "information", "contract", "contact", "specification". If a web page doesn't contain one word in the first list, it is not a CFT page. Otherwise, it must also contain one word in the second list. After that the web page should also contain word in the third list. Actually these evaluations check not only lexical content but also CFT format information.

In this chapter, we discuss some methods for domain specific search. Based on these, we developed a domain specific search system. In chapter 5, we will present the system application and implementation.

Chapter 5

# Implementation

In this chapter, we'll first introduce some software packages as the system construction tools. Then we discuss and compare different commercial search engines (CSE). Finally we talk about the system functions and some detail of implementation.

## 5.1    Survey of construction tools

The designing and construction of modern software is eased by the availability of high quality public domain system. For our applications, the following software packages were used: the HttpUnit package to be used in simulating a web browser to complete an URL connection; the Lucene will be introduced as a very good indexing tool for lexical analyzing; the HTMLParser package to be used in parsing and analyzing those downloaded Web pages; and the Package Regexp to be used in regular expression and the package HSQL Data Base to be used in data store.

### 5.1.1    HttpUnit

Automated testing is a great way to ensure that code being maintained works. HttpUnit is a suite of Java classes to test Web applications over HTTP. Coupled with Junit, HttpUnit is a powerful tool for creating test suites to ensure the end-to-end functionality of Web applications. One great aspect of HttpUnit is that it can test entire Web applications, not just single pages. HttpUnit makes it easy for a program to bypass the browser and access a web site from a program. Written in Java, HttpUnit emulates the relevant portions of browser behavior, including form submission, JavaScript, basic http authentication, cookies and automatic page redirection, and allows Java test code to examine returned pages either as text, an XML DOM, or containers of forms, tables, and links. Because it works closely with information and cookies, we can write tests to cover a whole session. For example, if our Web application includes a shopping cart, we

could write a test to try logging in, selecting an item, placing it the shopping cart, and checking out. Since the tests are written in Java, there's no limit to how in-depth the tests can be. We introduce the main parts of HttpUnit in the following.

### 1) Junit

JUnit is a simple framework to write repeatable tests. JUnit is a program used to perform unit testing of virtually any Java software. JUnit testing is accomplished by writing test cases using Java, compiling these test cases and running the resultant classes with a JUnit Test Runner.

### 2) Making a request

Since HttpUnit can emulate an entire session and not just a single request, the system uses a class, called WebConversation, to manage the requests, handle the cookies, and resolve relative URLs. As we write more complicated tests or simulate a browser to connect with an URL covering a whole session, the WebConversation class will become more important.

### 3) Parsing a response

Once we've made a successful request to a Web server, it's time to parse the result of the request. HttpUnit makes use of the JTidy package, which is included in the HttpUnit distribution, to parse the resulting HTML into a Document Object Model (DOM) tree. For those who aren't familiar with DOM trees, they offer a uniform way to manipulate a document in a hierarchical data structure. JTidy provides a standardized way to manipulate the HTML result.

It should be noted that, however, the DOM tree that JTidy builds represents the structure that a document should have, not necessarily what it does have. That means that JTidy may add structural elements to the document tree that aren't in the HTML source but should be. That includes head and body markers, paragraph marks, font tags, and more. Iterating through a DOM tree and visualizing the output, can be helpful.

### 4) Shortcuts through the DOM

Navigating through the DOM can be difficult and time-consuming. Fortunately, HttpUnit includes tools to make quick work of dealing with some HTML elements. HttpUnit contains similar shortcuts for navigating the DOM tree for forms and tables.

### 5) Navigating links

HttpUnit offers much more than just parsing the result of a single connection. The real power comes from being able to make multiple requests through the WebConnection object. The easiest way to make multiple requests is by following HTML links.

### 6) Posting forms

By far, the most common use of HttpUnit is to test the results of filling out a form. We can fill out a form of interactive Web page to achieve the desired effect. The WebConversation object takes care of making sure that all the links go to the correct endpoint and that all cookies are honored through the session. Each request uses the same WebConversation object, which ensures that the cookie values stay the same, simulating an entire session. The output of the test will include all the headers from the server so we can see that the cookies are being sent only once by the server.

In a word, all these features make sure that HttpUnit meet not only the need of testing, but also the need of a Web page analyzing. The chosen version is HttpUnit 1.5.4, published at 21 Aug, 2003, sizing around 2.59M.

### 5.1.2 Lucene

Lucene is a Java library that adds text indexing and searching capabilities to an application. Lucene offers two main services: text indexing and text searching. These two activities are relatively independent of each other, although indexing naturally affects searching. The fundamental Lucene classes for indexing text are: *IndexWriter*, *Analyzer*, *Document*, and *Field*. IndexWriter is used to create a new index and to add Documents to an existing index. Before text is indexed, it is passed through an Analyzer. Analyzers are in charge of extracting indexable tokens out of text to be indexed, and eliminating the rest. Lucene comes with a few different Analyzer implementations. Some of them deal with skipping stop words (frequently-used words that don't help distinguish one document from the other, such as "a," "an," "the," "in," "on," etc.), some deal with converting all tokens to lowercase letters, so that searches are not case-sensitive, and some like PorterStemAnalyzer will perform Porter stemming (or "Porter

stemmer") on its input, which is a process for removing the more common morphological and inflectional endings from words in English. For example, a group of words CONNECT, CONNECTED, CONNECTING, CONNECTION, CONNECTIONS will usually have similar meanings. After Porter stemming, the group of words are left CONNECT only. Its main function is to be part of a term normalization process that is usually done when setting up Information Retrieval systems. An index consists of a set of Documents, and each Document consists of one or more Fields. Each Field has a name and a value. Think of a Document as a row in a RDBMS, and Fields as columns in that row.

### 5.1.3   HTMLParser

HTML Parser is a library, written in Java, which allows to parse HTML (HTML 4.0 supported). It has been used by people on live projects. The architecture is flexible, allowing to extend it easily.   It is a super-fast real-time parser for real-world HTML. What has attracted most users to HTMLParser has been its simplicity in design, speed and ability to handle streaming real-world html. The chosen version is version 1.3, published on 24 May, 2003, sizing around 1.8M bytes.

### 5.1.4   Regexp

Regexp is a 100% Pure Java Regular Expression package that was graciously donated to the Apache Software Foundation by Jonathan Locke. He originally wrote this software back in 1996 and it has stood up quite well to the test of time. It includes complete Javadoc documentation as well as a simple Applet for visual debugging and testing suite for compatibility. The chosen version is Regexp 1.3, sizing around 195K bytes. It is a sub-project of Jakarta.

### 5.1.5   Hsql Data Base

HSQL Database Engine (hsqldb) Project hosted at SourceForge. The hsqldb is a relational database engine written in Java, with a JDBC driver, supporting a rich subset of ANSI-92 SQL (BNF tree format). It offers a small (less than 160k), fast database engine which offers both in memory and disk based tables. Embedded and server modes are available. Additionally, it includes tools such as a minimal web server, in-memory query and management tools (can be run

as applets) and a number of demonstration examples. The chosen version is HSQLDB version 1.7.1.

## 5.2    Choice of a commercial search engine (CSE.)

Normally, a web search application needs connect to URLs, downloads web pages, and analyzes them (Figure 5.2.1). But as a search research project, there is no need to implement the search application from scratch.    The part of accessing HTML pages can be replace by selected CSE to download HTML pages in the Web.



Figure 5.2.1 *Access HTML pages through CSE*

Currently, there are a lot of CSEs can be chosen. I tested 5 well known CSEs for both coverage and precision. These 5 CSEs are ALTAVISTA, GOOGLE, LYCOS, MSN, and YAHOO. In order to test the coverage of the five search engines, we first choose some typical query terms usually used in the "call for tender"(CFT) domain so as to check the searching results of different search engines. For out tests,  we used the following term: "Ottawa tender" and "Ottawa RFQ" (request for quotation). The results are listed in table 5.2.1. The tables is ordered by the number of hits:

| Order | CSE | Query      Term | Coverage |
|---|---|---|---|
| 1 | GOOGLE | Ottawa tender | 57,000 |
|   |   | Ottawa RFQ | 2,490 |
| 2 | YAHOO | Ottawa tender | 36,400 |
|   |   | Ottawa RFQ | 2,250 |
| 3 | LYCOS | Ottawa tender | 26,466 |
|   |   | Ottawa RFQ | 552 |
| 3 | ALTAVISTA | Ottawa tender | 20,196 |
|   |   | Ottawa RFQ | 496 |
| 3 | MSN | Ottawa tender | 19,647 |
|   |   | Ottawa RFQ | 637 |

*Table 5.2.1 Comparison of CSE (Coverage)*

From *Table 5.2.1* we can see, for the different CSE, although the quantity of coverage varies with the query term, it roughly reflects the coverage of each CSE. Google and Yahoo cover respectively much more pages than the other three CSEs (ALTAVISTA LYCOS & MSN).

It is not enough to get many hits (coverage); a user wants useful and pertinent hits.

| Order | CSE | Query    Term | Coverage | Precision (at 20) |
|---|---|---|---|---|
| 1 | GOOGLE | Manitoba tender | 33,700 | 35% |
| 2 | YAHOO | Manitoba tender | 22,700 | 30% |
| 3 | ALTAVISTA | Manitoba tender | 10,350 | 25% |
| 4 | MSN | Manitoba tender | 10026 | 25% |
| 5 | LYCOS | Manitoba tender | 13,396 | 10% |

*Table 5.2.2    Comparison of CSE (Coverage & precision)*

To test this, we examined precision manually, examining the first 20 URLs returned by each search engine. We respectively input a query term "Manitoba tender" to each CSE (see the *Table 5.2.2.*). Google or Yahoo is also much higher than the other three. Based on these experiments data, the Google and Yahoo were chosen as the selecting CSE for the domain specific system.

## 5.3 Introduction of the system

The system was developed with Java 1.4.It amounts to 8179 lines of code and includes 29 classes in addition to the public library package. The appearance looks like the Figure 5.3.1. The overall view of menus are composed by four parts: "URL"(choose CSE) and "City" are set for the system start; "Depth" and "Option" are set for crawler; "Analysis" is experiment tools; Searching results can be checked in menu "CFT storage", "CFT sites" , "CFT URL" & "Interactive CFT. We will first run the system then introduce the system functions. Before running the system, some parameters of the system need to be configured properly.
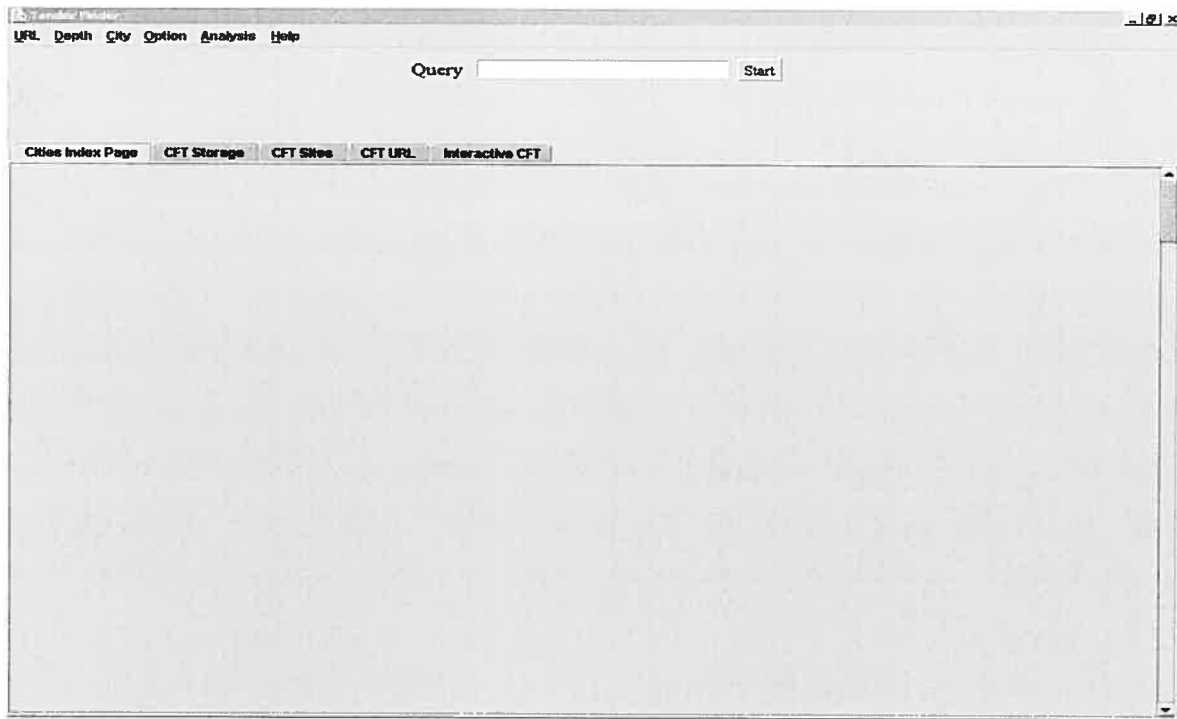


Figure 5.3.1  *Appearance   of   the   system*

### 5.3.1 Parameters Configuration

Figure 5.3.1 shows various menus: URL, Depth, City, Option, Analysis and Help choice to be set. Let's examine them one by one.

#### 1). URL

Both Google and Yahoo can be chosen as the CSE for initial request (Figure 5.3.2).



Figure 5.3.2  *URL configuration*

The setting process is, double clicking the menu URL, or pressing ALT+U, then appearing a pop-menu, double clicking the CSE in the window again.

#### 2). Depth

In our system, the crawler will examine all interesting pages within specified distance (depth) from the initial links returned by the CSE. See the figure 5.3.3, under the Depth configuration there are two sub-items: "crawling depth" and "page number". Like the URL

configuration, we can choose the "crawling depth" or "page number" respectively. The configuration of "crawling depth" decides layers a crawler further digging in. This parameter is the most time consuming factor. The more the "crawling depth" being chosen, the more time you need waiting, of course the more search results are provided.



Figure 5.3.3   *Depth   configuration*

The page number decides how much the CSE response pages will be chosen as seed pages.

### 3). City

In order to simplify the use and experimentation of our system in the Canadian context we have created a database of Canadian cities. Our system concentrates on one city at a time. To select a city, we can double click a "City" in city choosing pop-menu. Or we can select "No choice" to input in the query field later. When choosing a city, the chosen city will be combined with other query terms provides to the search system. See the figure 5.3.4 in the following page.

58



Figure 5.3.4   *City   configuration*

## 4). Option

The configuration of "option" includes two sub-items: "Tender extension" and "Tender Filter".    Tender extension aim to extend the query term so as to avoid tunneling and get more search results. You can input a word and add it to a list of "candidate words" and then press a button to forward the word to a list of   "Extension Terms", and vice versa. See the following figure 5.3.5

Figure 5.3.5 *Option    configuration --- tender extension*

There are another option is "tender filter". You can add some word to a list of "tender filter" to filter the searching result. For example, you can add a word "project" to see how many frequencies appearing in some group of documents. This can be combined used with lexical analysis to check a word occurring frequency.    See figure 5.3.6 in next page.

### 4). Analysis

This menu leads to the sub-system responsible for analyzing the HTML structure of a page to extract the portion that leads with the domain of interest. As shown in figure 5.3.6, analysis includes structure analysis and lexical analysis.

Figure 5.3.6 Option   *configuration --- tender filter*

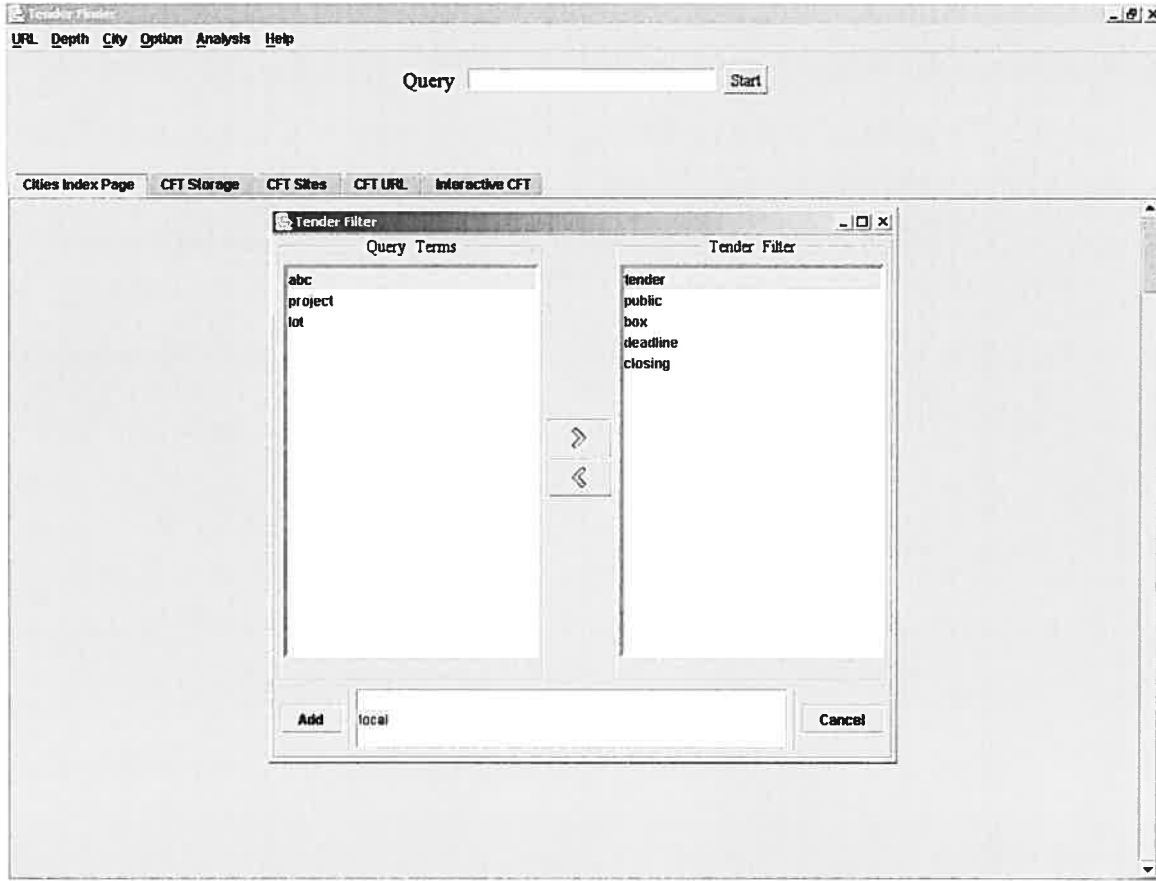"Structure analysis" like figure 5.3.7 shown, you input an URL like: http://www.cityofkingston.ca/cityhall/tenders/index.asp and press the button "Begin". The system will explore that web page and analyze it. The right side of the pop window is content of that web page; the left side of the pop window is structural analyzing results, which marks sub-table with an order. One of these sub-tables is a selected one, which marks with capital letter like SUB_TABLE with a number 23. See figure 5.3.7 in next page. The SUB_TABLE_23 looks like figure 5.3.8 shown.

Necessity of analyzing logical structure for a CFT page based on the following considerations. The more exploring paths a crawler has the better searching coverage it has. On the other hand, while the crawler meets a qualified CFT page, it should bring the current exploring process to a temporary end, downloads this page, and starts another exploring process.

Figure 5.3.7 *Analysis configuration --- structural analysis*



Figure 5.3.8     *Selected sub-table content*

It is easy to process a CFT page which contains mainly text or mainly links. But for a mixed page like Figure 5.3.9, extracting a middle part of text content can improve accuracy of a lexical evaluation. Besides, such cutting process can also be used for analyzing a CFT page part by part.



Figure 5.3.9    *Mixed page*

The detail of decomposing mixed page to several sub-tables has been explained in chapter 4.

Figure 5.3.10 in next page shows the lexical analysis process.

Figure 5.3.10 *Analysis configuration --- lexical analysis*

The lexical analysis is a statistic tool for analyzing various web pages. In the above window, the top window is for analyzing a web pages or a set of pages, the bottom window is for analyzing a saved example web file or a local directory which contains many web files. After inputting files in a URL or a Directory field, pressing the button "Start" or "Begin" will produces results as the above figure showing. These results count all the words in documents and their frequency of occurrence.

The directory d:\htmllucene\special_tenders includes three downloaded web pages of Vancouver area. Their contents are shown in the next page, respective figure 5.3.11.a, figure 5.3.11.b and figure 5.3.11.c.

Specification Number I-02.02.31.04

Two identical 15 degree dipole magnets are required for the delivery of 500~MeV protons to the new ISAC Target Facility. TRIUMF will supply the copper for the excitation coil manufacture and the steel for the magnet yoke and pole manufacture. The invitation to tender is for the manufacture and assembly of completed magnets.

Base Specification:

Bend Angle: 15 degrees for 500 MeV protons

Maximum Field: 10.241 kG

Effective Length: 37.477 inches along arc

Pole Gap: 4.000 inches

Pole Shape: Rectangular pole; entrance and exit angle = 7.50 degrees. No curvatures

Weight: Each magnet, including copper coils, weighs approximately 11,000 lbs.

TRIUMF will provide the detailed specifications as part of the tender package.

Figure 5.3.11.a    *Saved tender sample of Vancouver*

TENDER SPECIFICATION # I-0202.31100.1

A tender package will be issued by December 20. 1996 for the manufacture of 10 small dipole electromagnets for use as steering correction magnets in the beam transport of 500 MeV protons in beam line 2A from the TRIUMF cyclotron to the new ISAC facility.
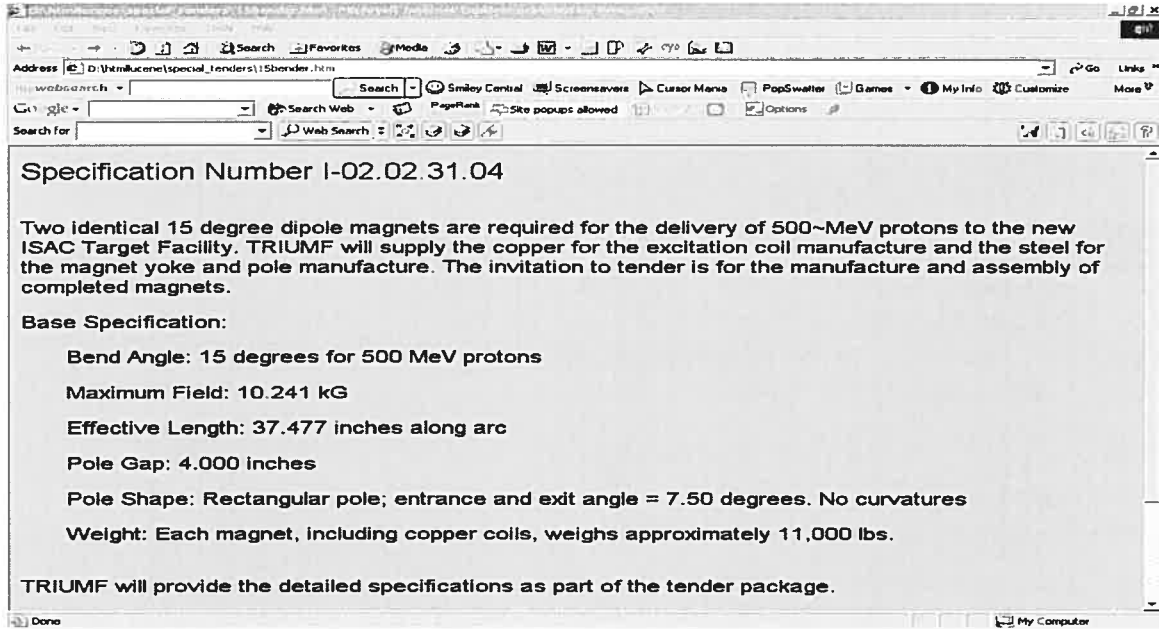
The initial specifications are as follows:

The steering magnets consist of steel sub-assemblies and electrical excitation coils of 1000 turn windings of 0.0808 inches square wire giving a coil resistance of about 3 ohms. The coils will be energized with a direct current of not more than 5 amperes.

TRIUMF will provide the detailed specifications as part of the tender package.

As a federally funded Canadian National Research Laboratory, TRIUMF encourages contractors, suppliers and manufactures to optimize Canadian content, provided it is achieved in a technically competitive and cost effective manner.

Companies interested in bidding on the above mentioned items should fax a letter of interest to Roy Moore (Contract Administrator) (604) 222-7307 no later than December 13. 1996.
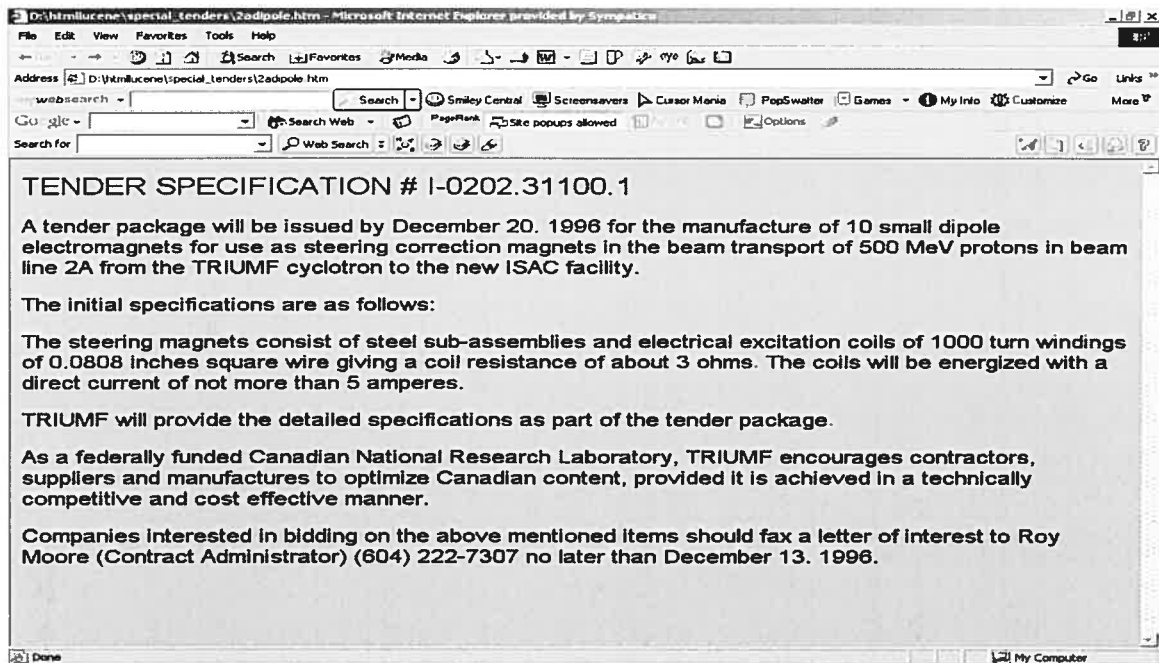
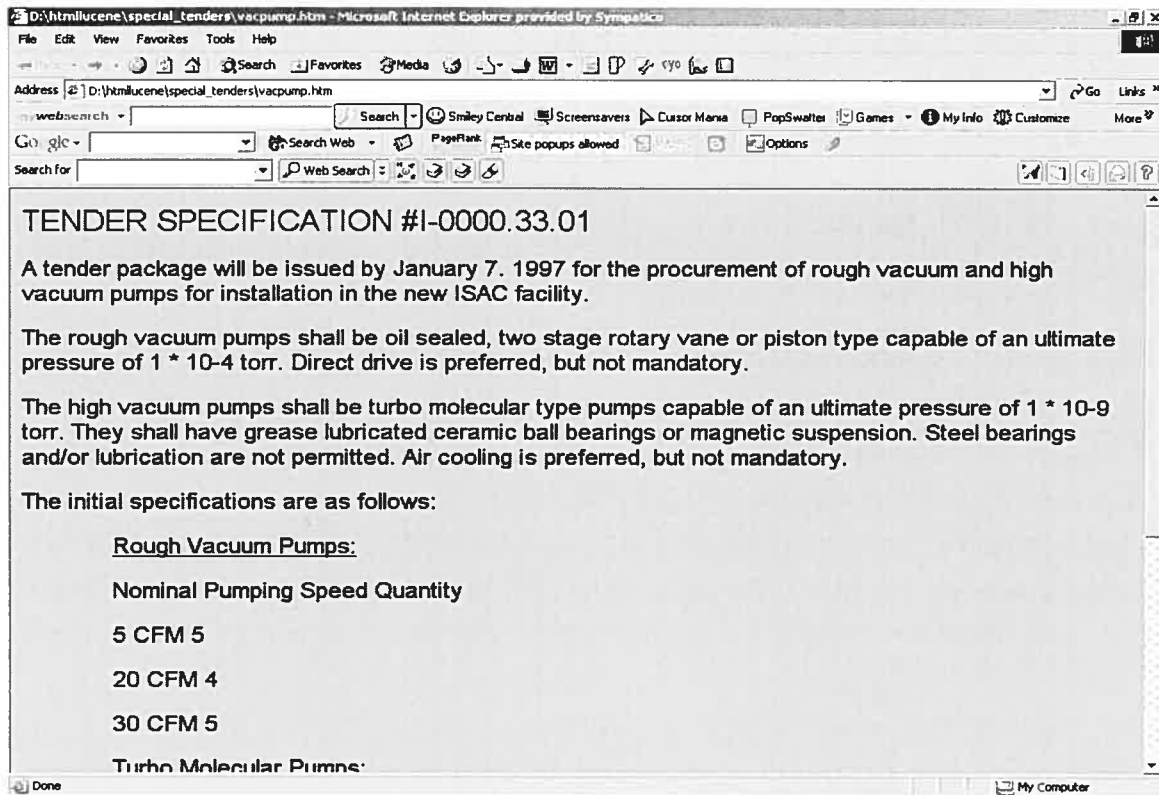Figure 5.3.11.b    *Saved tender sample of Vancouver*

Figure 5.3.11.c    Saved tender sample of Vancouver

The lexical analysis also provides "T-Start" and "T-Begin" buttons which trigger analysis according to pre-selected key words pertinent, in our case, to tenders specifically, the system will output the frequency in the document of keywords entered previously under the "tender filter" option. See the following figure 5.3.12a and figure 5.3.12b.   In the figure 5.3.12a, we choose "lot", "tender", "project" and "deadline" as our filter words. In the figure 5.3.12b, the result shows that only "lot", "tender" and "project" occurring frequencies. Same result when input a directory and pressing button "T-Begin".
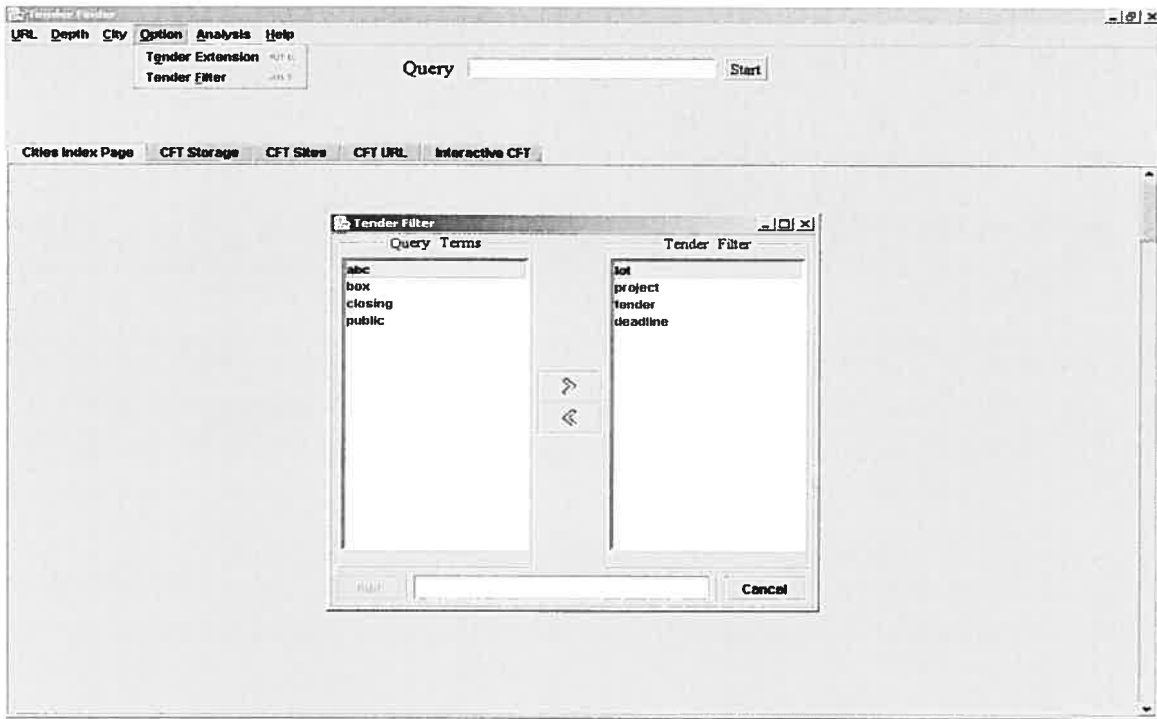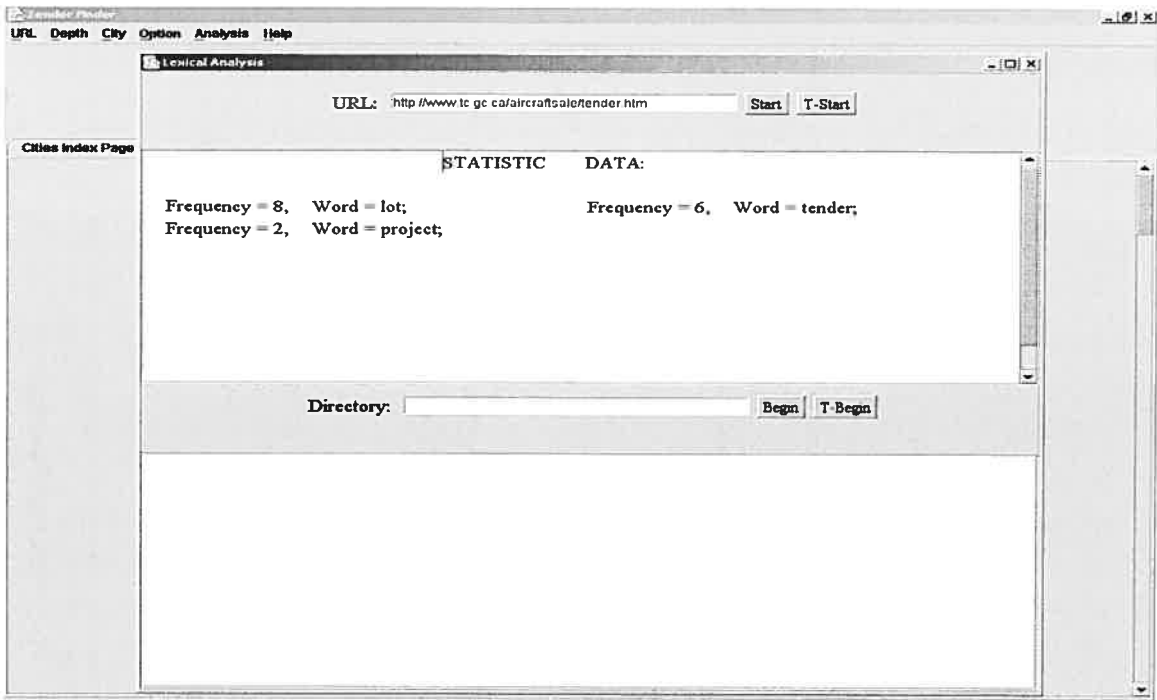
Figure 5.3.12a    *Choice of tender filter*
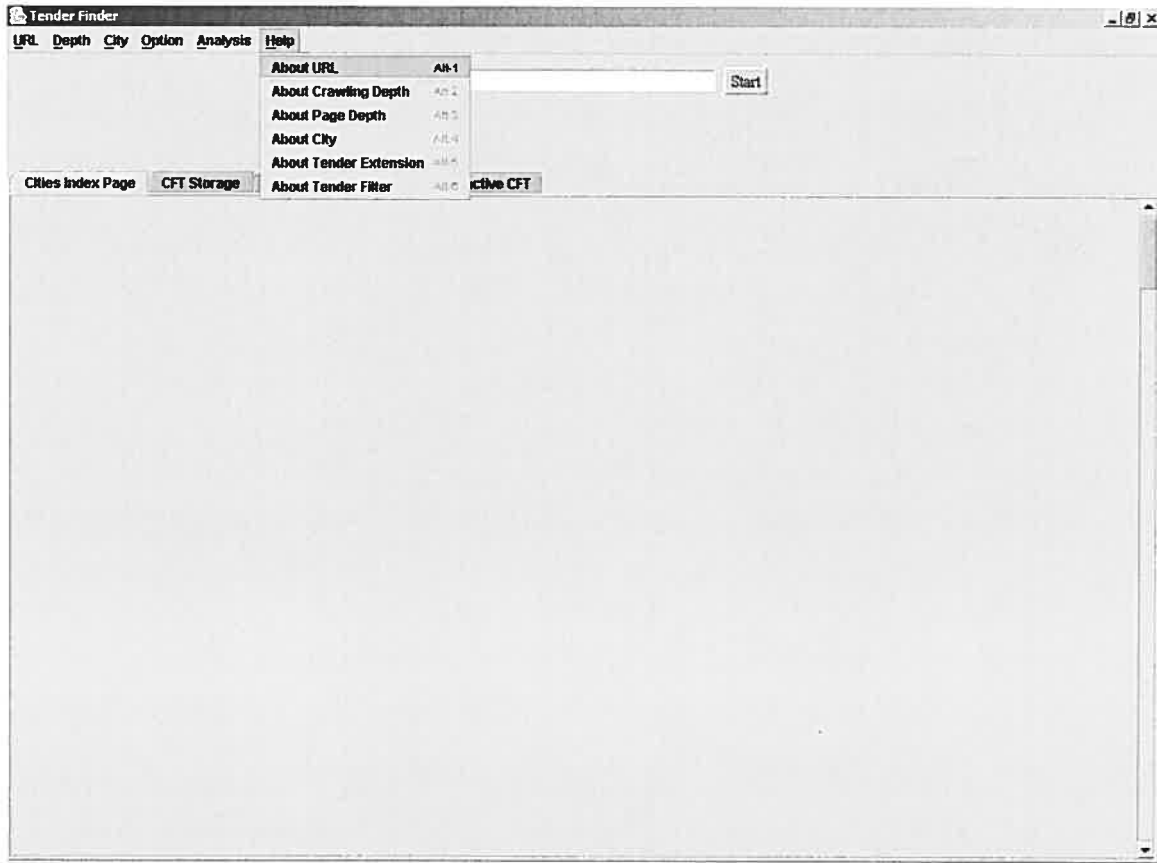


Figure 5.3.12b    *T-Start result*

## 5). Help



Figure 5.3.13    *Help menu*

The "help" provides explanations for all the parameter configurations. It helps user to configure properly.

## 5.3.1 Running the system

We start a session by typing some keywords into the query field and pressing "start". Typically query would be a city with a term "tender" (you can input others if you choose "No choice" in the city pop-menu). This causes initial seed URLs to be obtained from Google or Yahoo then our crawler explores all interesting links to the specified depth limit and stores its results in the local database. When it has finished its exploration, the results can be examined via the TABBED sub-menus of the bottom window. We now illustrate the functionality of

our system by going through a typical session The system includes 5 tabbed sub-menus which represent 5 different functions: Cities index page, CFT storage, CFT sites, CFT URL and interactive CFT. We'll introduce them one by one.

### 1). Cities index page

The index page is a site which can help us to navigate further to each province or city's home page (see figure 5.3.1a, figure 5.3.1b).
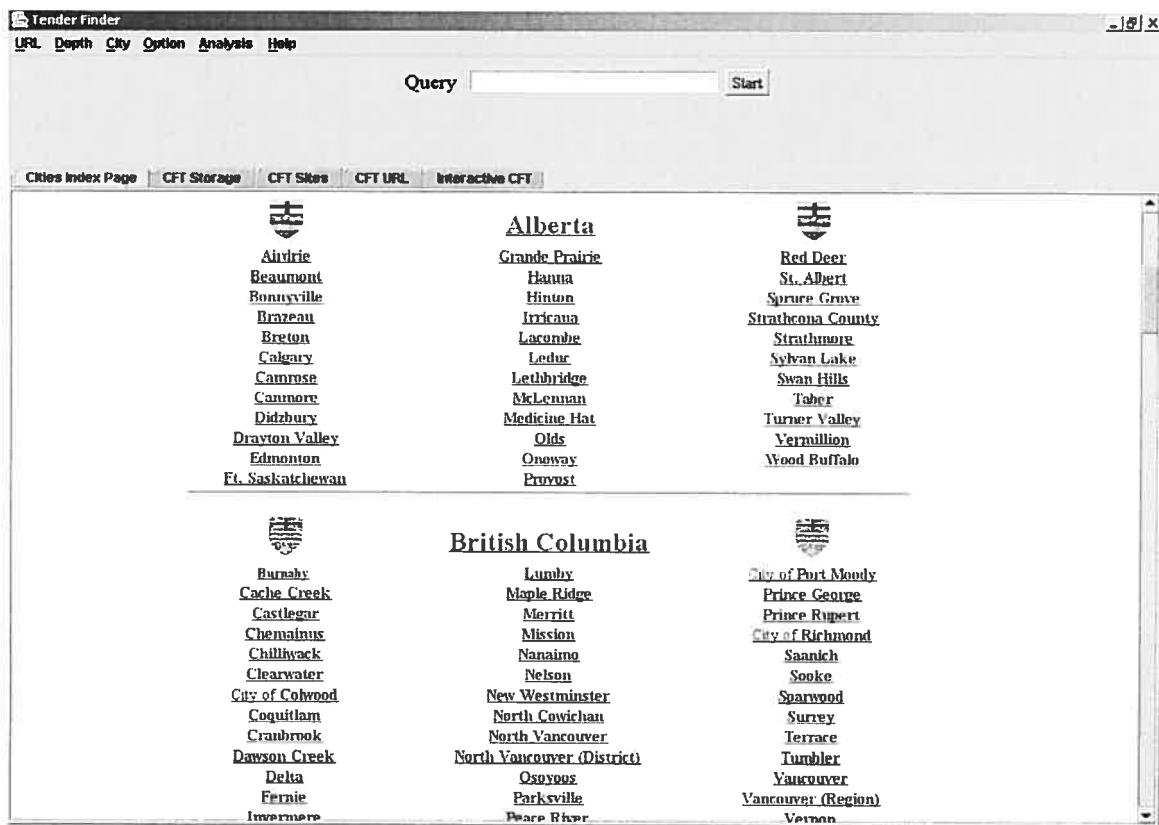


Figure 5.3.1a   *City Index Page*

Figure 5.3.1b     *City Index Page*

This is an auxiliary function for user to navigate to Canadian cities so as to access home page of each city. Figure 5.3.1b shows the result of selecting Hanna, Alberta.

## 2).    CFT Storage

Clicking on the CFT storage tab shows the contents of the database (see Figure 5.3.2a) and clicking on an entry opens a further window showing the contents of the selected file. The crawler stores all the tender related web pages to a database according to content of pages or a URL string of a hyperlink. If a web page is a CFT page or "tender" is contained in the URL string, then the URL will be store in a database. See the Figure 5.3.2a in next page.

Figure 5.3.2a    *CFT Storage    --- a CFT pertinent page*

The clicked URL string in the figure 5.3.2a is:

http://www.gracekennedy.com/grace/newstendernotice.htm.    It is a CFT page.

Figure 5.3.2b, Figure 5.3.2c and Figure 5.3.2d are other three random examples within the CFT storages. The first two are highly pertinent to CFT pages. The third is not. The check of most URLs within these results shows that those pointed pages are highly pertinent to CFT content. Thus, a URL string can be used in a fast content check before a pointed web page being downloaded.

Figure 5.3.2b     *CFT Storage --- another CFT page*

The first example (Figure 5.3.2b) comes from Niagarafalls City which is a CFT.



Figure 5.3.2c     *CFT Storage    --- CFT pertinent page*

The second example (Figure 5.3.2c) from Vaxxine is also pertinent a CFT pertinent page (hub page).

Most of CFT storages are tender pertinent pages, this is because a part of them are collected by content lexical analysis, these pages are CFT pertinent pages with sure; others are collected by their URL strings which contain the term "tender", these pages seem to be CFT or CFT pertinent pages, but not for sure. Figure 5.3.2.d shows such an exception.



Figure 5.3.2.d     CFT Storage --- not a CFT page

The last example (Figure 5.3.2d) from Vancouver is of no interest.

### 3).    CFT sites

To do this, we reorder URLs in the CFT storage with common prefix as nodes in a tree. This function aims at finding some CFT hubs, where a lot of CFTs are gathered.

Through analyzing each candidate URL string, try to find out those CFT Hubs which are marked as middle nodes of a Jtree in    figure 5.3.3.a and figure 5.3.3.b



Figure 5.3.3a          *CFT sites*

The site in figure 5.3.3a is http://www.city.toronto.on.ca including following leaf nodes which are all CFT pertinent. Specially, 4) is a CFT hub where we are looking for.

1) /tenders/toronto_zoo

2) /tenders/proposal.htm

7)  /tenders/quotation.htm

8)  /tenders

Again our system allows examination of the content of entries. The site in figure 5.3.3b and 5.3.3c is http://contractscanada.gc.ca, its leaf nodes are:    1)/en/tender-e.htm and 2)/en/cmhc-e.htm. 1) is an indirect CFT hub where lead to a CFT hub http://www.merx.com .

Figure 5.3.3.b          *CFT sites*



Figure 5.3.3.c          *CFT sites*

Through summarizing huge amount of separate web pages, we can find some URL path where usually gather some up to date CFTs. These URL paths normally are more meaningful and stable than some separate CFTs.　Just like Figure 5.3.6.c shows.

**4).　CFT URL**

CFT URL is a result that only based on the current searching phase. In the figure 5.3.4.a, we choose the city as Calgary first, and　then query "tenders". We check the list in the figure 5.3.4.a . The figure 5.3.4a, figure 5.3.4b and figure 5.3.4c　all satisfy our requirements. We'll compare the CFT URL results with Google and Yahoo in next chapter.



Figure 5.3.4a　　　　*CFT URL ---- (1)*

(1)　http://www.calgarywinterclub.com/constructionmar01.html

Figure 5.3.4b          *CFT URL --- (2)*

(2)   http://www.gasandoil.com/contracts/samples/2004issue22.htm



Figure 5.3.4c          *CFT URL --- (3)*

(3)   http://surplus.gov.ab.ca

**5). Interactive CFT**

Within the whole set of CFT pages, there are pages with various forms. The most interesting act as gateways to database of CFTs. But usually these pages have "forms" for user to communicate interactively. These are so called interactive CFT pages. To access these pages requires fill and submit forms, register as a member or create an account etc. All these steps increase the complexity for a crawler to access them automatically. Our system collects these pages, characterized as "data" pages with HTML forms. They are shown by clicking "Interactive CFT" tab. The Figure 5.3.5a shows such an interactive CFTs. collection during one search process. The Figure 5.3.5b is another example of interactive CFT formats.



Figure 5.3.5a   *Interactive CFT*

Figure 5.3.5b   *interactive CFT*

In this chapter, we introduced the system configuration, running and functions. In the following chapter we will test the system and also evaluate its general performances by comparing with the CSE the system based on.

# Chapter 6

# Testing and Evaluation

In this chapter we'll test the domain specific system with some sets of web pages. Testing will cover system tools test such as the structural test, the lexical testing as well as system function test including CFT storage, CFT sites and CFT URL etc.

## 6.1    System tools test

To evaluate web pages, our system also provides two system tools to support the analysis of web pages. These two tools are structural analysis and lexical analysis.

## 6.1.1   Structural analysis

The structural analysis focuses on extracting a text area from a web page with a complex logical structure. Considering most complex web pages are composed with tables within tables, this extracting process sets out from decomposing a web page to several tables and then analyzes these tables. We choose two examples web page for structural analysis testing. The examples are:

1)  http://www.cityofkingston.ca/cityhall/tenders/index.asp
2)  http://www.tc.gc.ca/aircraftsale/tender.htm

**Test 1:**

The page concerns Kingston * looks like figure 6.1.1a and figure 6.1.1b in next page. Previously, in section 5.4, we used this page to show our structural analysis tool.

* http://www.cityofkingston.ca/cityhall/tenders/index.asp

Figure 5.3.6 showed that it was composed of about 30 HTML tables and table 23 was judged to contain the most pertinent data. Figure 6.1.1c and Figure 6.1.1d show this part of the page in more detail.



Figure 6.1.1a   *Mixed page top side*



Figure 6.1.1b   *Mixed page bottom side*

After extracting:



Figure 6.1.1c   *Extracted sub-table top side*



Figure 6.1.1d   *Extracted sub-table bottom side*

## Test 2:

The original web page looks like figure6.1.2a and extracted one like figure 6.1.2b shown. There are not much differences after extracting.



Figure 6.1.2a    *Original web page top side*



Figure 6.1.2b    *Extracted sub-table top side*

These results show that the structural analysis is especially proper for those pages which contain a complex logical structure.

## 6.1.2 Lexical analysis

The lexical analysis provides for both web files and local directory/files. The analysis can also accumulate the words and frequencies during an analyzing session. Therefore, we can use it to analyze a group of web pages.

For example, we have 2 web pages:

1)Edmonton City:

http://www.city.quintewest.on.ca/services/finance/tenders/PW04-17.htm

2)    St John City: http://www.stjohns.ca/csj/TendersDetails?id=156&loc=W

add them one by one to lexical analysis, we get the result as following:



Figure 6.1.3a    *Page 1) lexical analysis*

Figure 6.1.3b    *Page 2) lexical analysis*



Figure 6.1.3c.    *Page 1) & 2) lexical analysis*

From the results we can see that the 5 most high frequency words in first page are: tremur (4), lake (4), tender (3), decommissioning (3) and station (2). For the second page these high frequency words are tender (5), department (4), hall (3), home (3), and lot of words with frequency equal 2, such as outdoor, employment, news, date, city etc. The figure 6.1.3.c shows the cumulative frequencies for the two pages. The first 5 high frequency words are: tender (8), trmur(4), lake(4), department(4), hall(3), services(3), city(3), decommissioning (3), home(3), station(2) etc. So, if we want to test a set of web page, just add them to the URL field and press start button.

## 6.2    System test

To test the system functions, we first configure the system parameters. For example, we configure page_number=1, crawl_depth=2; city=Calgary; no tender extension and no tender filter.

### 6.2.1    CFT storage and CFT sites



Figure 6.1.4a *CFT storage -test*

Figure 6.1.4b *CFT storage -test*

Both CFT storage and CFT sites accumulate CFT pertinent URLs each searching process. In figure 6.1.4b, the current searching city is Calgary, and the marked URL: http://www.ucalgary.ca. is added in the end of the CFT sites. At present, there are 829 saved URLs in CFT storage and 456 URLs in CFT sites. About 75% of the CFT storage or CFT sites are CFT pertinent URLs.

## 6.2.2    CFT URL

This part will evaluate the searching results of our system. At the same time input the same query term to Google. Comparing the top 10 responses with each other .

### 1)  Top 10 of Domain specific searching system

CFT URLs are current searching results. There are 30 URL listed in the window. The first 10 are listed in the following page (see figure 6.1.4cb, there are more CFT URL tests, see appendix 2).

1) http://www.calgarywinterclub.com/constructionmar01.html

2) http://www.lexum.umontreal.ca/csc-cc/en/pub/1987/vol2/texte/1987scr2_0757.txt

3) http://www.gasandoil.com/contracts/samples/2004issue22.htm

4) http://surplus.gov.ab.ca/attachments/TNCALOUT_2003_044__High_River_.pdf

5) http://www.lexum.umontreal.ca/csc-scc/en/pub/1987/vol2/html/1987scr2_0757.html

6) http://www.cbe.ab.ca/boes/buspurch/procedures.asp

7) http://content.calgary.ca/CCA/City+Hall/Business+Units/Finance+and+Supply/Tenders

8) http://www.calgarysun.com/htdocs/static/calsun/ratecards/general.pdf

9) http://www.businessedge.ca/viewnews.asp?id=6680

10) http://www.millerthomson.ca/issue.asp?Print=Yes&NL=2&Year=1997&Season=1

Figure 6.1.4c *CFT URLs –test*

Let's evaluate this system's response. Figure 6.1.5a shows the content of

http://www.calgarywinterclub.com/constructionmar01.html



Figure 6.1.5a *CFT URLs —response evaluation*

2)  http://www.lexum.umontreal.ca/csc-cc/en/pub/1987/vol2/texte/1987scr2_0757.txt file is not found at testing moment;

3)  http://www.gasandoil.com/contracts/samples/2004issue22.htm is a good URL point to a CFT in Calgary.

Figure 6.1.5b *CFT URLs –response evaluation*

4) http://surplus.gov.ab.ca/attachments/TNCALOUT_2003_044__High_River_.pdf is a PDF file. The content is a CFT document in Calgary. See the PDF content in figure 6.1.5c in the following.

A **SURPLUS SALES CALGARY**
**CORPORATE SERVICE CENTRE TN.CALOUT.2004.044**
**3850 MANCHESTER RD, SE**
**CALGARY, AB. T2G 3Z9**
**PHONE: (403) 297-6430**
**FAX: (403) 297-4576**
**INVITATION TO TENDER**
**TENDER No:**
Sealed Tenders, subject to the General Terms and Conditions – Public Tenders and any Special Terms and Conditions,
will be received at the office of Alberta Corporate Service Centre, Surplus Sales Section at the above address, **NOT**
**LATER THAN: 14:00:59 Alberta Time, Tuesday, September 21, 2004**
Tenders will be opened in public at the above noted address, time and date, FOR THE SALE OF:
**Lot No. Description Quantity**

1 HOUSE, 3 BEDROOM, TWO BATHS, AND
FIREPLACE, APPROX 1,700 SQ FT
1 OFFERS
2 THREE STALL GARAGE, APPROX 720 SQ
FT
1 OFFERS
3 STORAGE BUILDING, APPROX 10' X 12' 1 OFFERS
4 STORAGE BUILDING, APPROX 8' X 16' 1 OFFERS
5 HORSE SHELTER, APPROX 8' X 12' 1 OFFERS
6 CORRAL, APPROX 400 SQ FT 1
7 LOT OF 12, PLANTED SPRUCE TREES 1 LOT
8 LOT OF 40, MATURE POPLAR TREESS 1 LOT OFFERS
## GST WILL BE ADDED TO THE BID PRICE
used, it is taken as confirmation that the bidder wants to purchase only one lot.
**For pictures and additional Bid Forms, please visit us on the Internet at:**
**http://surplus.gov.ab.ca**
LOCATED AT:
NW9-19-28-4, HWY 2 AND HWY 543
NORTH OF HIGH RIVER
VIEWING IS BY APPOINTMENT ONLY
FOR FURTHER INFORMATION AND APPOINTMENT
CONTACT: DICK DAVISON, (403) 382-4077
TN.CALOUT.2004.044

Figure 6.1.5c *CFT URLs – response evaluation*

5) http://www.lexum.umontreal.ca/csc-scc/en/pub/1987/vol2/html/1987scr2_0757.html is tender related page, but strictly speaking, not a CFT (Figure 6.1.5d).

6) http://www.cbe.ab.ca/boes/buspurch/procedures.asp is a CFT in Calgary (Figure 6.1.5e).

7) is a tender hub where indicate to *Commodities Tenders, Construction Tenders and Services Tenders* (see four Figure 6.1.5f - Figure 6.1.5f (3)).

http://content.calgary.ca/CCA/City+Hall/Business+Units/Finance+and+Supply/Tenders

Figure 6.1.5d *CFT URLs – response evaluation*



Figure 6.1.5e *CFT URLs – response evaluation*

Figure 6.1.5f  *CFT URLs – response evaluation*



Figure 6.1.5f (1)  *Commodities tenders*

Figure 6.1.5f (2)    *Construction tenders*



Figure 6.1.5f (3)    *Services tenders*

8)   http://www.calgarysun.com/htdocs/static/calsun/ratecards/general.pdf is just a advertisement of a newspaper. Similar to a tender of advertisement, but not a real CFT (Figure 6.1.5g ).



Figure 6.1.5g   *Services tenders*

9) http://www.businessedge.ca/viewnews.asp?id=6680    is shown in Figure 6.1.5h in the next page. It's a tender pertinent report which published at 8/19/2004 –Vol. 1, No.17. The title is **"Bid shopping' accusation stings province"**, and sub-title is "**Contractors applaud ruling that forbids hybrid process By Monte Stewart - Business Edge**". Actually this is a typical example of CFT confusion type. It is really a difficulty one for crawler to distinguish this page as a CFT or not. Because from lexical view, it meet nearly all the need of CFT, such as tender, value, date, time etc. we'll discuss this type of page in next chapter.

Figure 6.1.5h *CFT URLs – response evaluation*

10) http://www.millerthomson.ca/issue.asp?Print=Yes&NL=2&Year=1997&Season=1

is an article talking about construction law which pertinent with the construction tender.



Figure 6.1.5i *CFT URLs – response evaluation*

The title of the 10) pointed article is "Construction Contracts - Off the Rack or Custom Made?"

Summarizing above top 10 searching results gets following table 6.1.1.

| Category | Quantity | Percentage | Comment |
|---|---|---|---|
| CFT in Calgary | 6 | 60% | 1),3),4),5),6),7) |
| Similar to CFT | 1 | 10% | 8); in Calgary |
| File not Found | 1 | 10% | 2) |
| Others | 2 | 20% | 9) CFT news;10)Construction law |

Table 6.1.1 *Summarizing of evaluation*

## 2) Top 10 of Google

At the same time, inputting the "Calgary tenders", top 10 search results provided by Google are :

a) http://content.calgary.ca/CCA/City+Hall/Business+Units/Finance+and+Supply/Tenders/Tenders.htm

b) www.westernwheel.com/WW-PAGE%2041.pdf

c) surplus.gov.ab.ca/attachments/ TNCALOUT_2003_044__High_River_.pdf

d) www.lexum.umontreal.ca/csc-scc/ en/pub/1987/vol2/html/1987scr2_0757.html

e) www.lexum.umontreal.ca/csc-scc/ en/pub/1987/vol2/texte/1987scr2_0757.txt

f) www.cca.cc/about_us/practice.htm

g) www.cbe.ab.ca/boes/buspurch/procedures.asp

h) directory.teradex.com/Regional/Calgary/Government

i) ciqs.org/jobs/res-2003-10-17.pdf

j) www.pc.gov.bc.ca/data/docs/02-05-03_511.6.pdf

Let's evaluating   Google's top 10 response results one by one:

a)  same as 7) a CFT hub point to CFTs in Calgary;

b)  see the figure 6.1.6a, a tender advertisement nested in a newspaper.



Figure 6.1.6a    *CFT URLs — comparison*

c)  same as 4), is a CFT in Calgary;

d)  File not found;

e)  File not found;

f)  See the figure 6.1.6b in next page, it is a "CODE OF PRACTICE" for members of the Calgary Construction Association.

g)  same as 6), is a CFT page in Calgary;

Figure 6.1.6b  *CFT URLs — comparison*

h)   http://directory.teradex.com/regional/calgary/government   is not a CFT page.



Figure 6.1.6c *CFT URLs — comparison*

i) ciqs.org/jobs/res-2003-10-17.pdf   is not a CFT PDF file;

```
EXPERIENCE

Free -Lance Estimator                        April 2002 to Present
    • Consultant on Structural Steel and Miscellaneous Metal Projects
    • Producing bid proposals for various clients

COLONY Management Ltd., Vancouver BC         July 2001 to April 2002
Business Development Manager
    • Business development of new markets primarily in the USA and Alberta.
    • Marketing and sales of Pre-Engineered Metal Building and Stress Skin Building.
    • Developing new product line, the custom conventional steel.
    • Company yearly sales is $10 million
    • Duties include prospecting for new clients, researching for new markets.
    • Producing tender proposals, negotiation with prospective clients.

AMEC Dynamic Structures, Port Coquitlam BC   January 2000 to July 2001
Sales
    • Sales representative for the territories of British Columbia and Northwest USA (Seattle and
      Portland Regions)
    • Marketing and sales of Conventional Structural Steel
    • Company yearly sales is $30 million
```

Figure 6.1.6d *CFT URLs — comparison*

j)   www.pc.gov.bc.ca/data/docs/02-05-03_511.6.pdf   is a CFT useful page, but not special for Calgary CFT;

```
The Agreement on
Internal Trade

The Agreement on Internal Trade (AIT) is an intergovernmental agreement between the Federal Government, the Provinces, the
Northwest Territories and the Yukon to reduce and eliminate barriers to the free movement of people, goods, services and
investments within Canada. Under the Agreement, these governments have agreed to apply the principles of non-discrimination,
transparency, openness and accessibility with respect to their procurement opportunities and those of their municipalities and
municipal organizations, school boards and publicly funded academic, health and social services entities. The Agreement covers
only those tenders where the procurement value exceeds a specified amount. To learn about these threshold amounts and other
provisions of the Agreement, suppliers should consult the web site of the Internal Trade Secretariat, found at
www.intrasec.mb.ca.

Suppliers can also visit MARCAN, a new initiative under the Agreement, found at www.marcan.net. MARCAN provides links
to Internet sites, including the sites indicated below, that are known to publish tender notices by the Canadian public sector and
para-public sector for the purchase of goods and services and for construction projects.

GOVERNMENT PROCUREMENT INFORMATION
```

| | Who to contact<br>1   for enquiries and complaints<br>2   for information on how to register on a<br>     source list. | Where tenders may be advertised |
|---|---|---|
| Government of Canada | ·   General enquiries | www.merx.bmo.com |

Figure 6.1.6e *CFT URLs — comparison*

Summarizing above top 10 searching results get following table 6.1.2.

| Category | Quantity | Percentage | Comment |
|---|---|---|---|
| CFT in Calgary | 3 | 30% | a),c),g) |
| Similar to CFT | 2 | 20% | b) in Alberta, j) in Canada |
| File not Found | 2 | 20% | d),e) |
| Others | 3 | 30% | f),h),i) less relation to CFT |

Table 6.1.2 *Summarizing of evaluation*

### 6.2.3   Interactive CFT

There are a lot of interactive web pages in CFT domain. Usually they appear as various formats with different forms. It is not easy to auto-fill those forms with an unique way. Anyway, collecting and recording them may help use to grasp some hints for digging further CFT resources.

Interactive CFT same as the CFT URL is a one session searching. At the beginning, a user inputs a query term and presses the "Start" for starting a searching session, after that double click the tabbed menu "Interactive CFT" and wait for a moment. The appearance looks like the following figure.



Figure 6.1.7   *Interactive CFT example*

Chapter 6 introduced and tested the domain specific system's configurations functions and performances. In chapter 7, I'll give a brief conclusion and some future works to do.

# Chapter 7

## Conclusions

Performing a search in a specific domain with a traditional, general-purpose search engine is tedious. For this reason, domain specific search engines are becoming increasingly popular, but techniques for successful operation are not fully established. In this thesis, a domain specific searching system (DSS) was developed to complement CSEs (Commercial Search Engine), such as Google and Yahoo. The system employs some specific technologies in web page logical structure analysis and lexical analysis. The prototype has been successful in discovering call for tender pages on the web sites of the cities we tested it with. To achieve this goal, the crawler has to combine many techniques and heuristics. Here are some of the ideas that were found necessary to achieve success:

1) Use of existing search engines to find good starting points for further exploration.

2) Analysis of page structure before examining content. It is critical to distinguish HUB, CONTENT and MIXED pages and to treat each differently.

3) Pages with a complex HTML structure (i.e. Frames and Tables) must be analyzed specially to isolate the information core from repetitive decoration.

4) A keyword profile is important to quickly evaluate the relevance of pages.

5) Development of advanced web software can benefit from the availability of high quality public domain packages.

Preliminary comparisons with Google using additional terms for maximum precision, show that our system's precision is twice as good as Google (about 60% to 30%, see the test data in Chapter 6: Table 6.1.1 & Table 6.1.2). Full evaluation of the performance will have to wait for the next phase of the project; but already it is apparent that the main problem will not be in smart navigation of the web; but rather in establishing classifiers that can detect satisfactorily what is a relevant page and what is not.

Typically all pages found and accepted by our system deal with CFTs in the given city. But most are not CFTs: some are news about contracts, others refer to the terms and conditions of doing business, others are HUBs etc... An automatic classifier is at the heart of any intelligent discovery agent, but more research is required to design a classifier that could learn the subtle differences between these documents.

In retrospect, one search system might function best as a semi-automatic work-bench to allow domain experts to discover sites and write collector scripts for domain specific intelligences gathering systems.

Now that we have a working prototype, we will concentrate on the design of a good classifier as well as implementing the learning techniques proposed by Aggarwal [Aggarwal et al. 2001] in his description of intelligent crawling. We also want to study the best ways to tunnel through to hidden web pages by automatically filling out forms. Finally, we will design experiments to properly evaluate the performance of the system.

# Bibliography

[Aas et al. 1999] "Text Categorisation:A Survey" Kjersti Aas and Line Eikvil, http://citeseer.ist.psu.edu/aas99text.html

[Aggarwal et al. 2001] "Intelligent Crawling on the World Wide Web with Arbitrary Predicates", Charu C. Aggarwal, Fatima Al-Garawi & Philip S. Yu, 2001. WWW10 May 1-5,2001. ACM 1-58113-348-0/01/0005

[Apache Jakarta Project 2003] Apache Jakarta Project
http://jakarta.apache.org/lucene/docs/index.html

[ Ask Jeeves, Inc. 2004] Ask Jeeves, Inc. http://www.ask.com/

[Baroni et al. 2002] "Using textual association measures and minimum edit distance to discover morphological relations", Marco Baroni, Johannes Matiasek & Harald Trost, Collocations Workshop, July,23,2002.

[Bergmark et al. 2002] "Focused Crawls, Tunneling, and Digital Libraries"
Donna Bergmark, Carl Lagoze, Alex Sbityakov, in CiteSeer.IST 2002.
http://citeseer.ist.psu.edu/bergmark02focused.html

[Bra 1994] P. De Bra, R. Post. Searching for Arbitrary information in the WWW :the Fish-Search for Mosaic. WWW Conference, 1994

[Brin et al. 1998] "The Anatomy of a Large-Scale Hypertextual Web Search Engine" by Sergey Brin and Lawrence Page, 1998

[Carchiolo V. et al. 2003] "Extracting Logical Schema from the Web", Longheu A. & Malgeri M., in "Applied Intelligence" Vol.18, pp341-355, 200, May, 2003

[Carchiolo et al. 2003] "Extacting Logical Schema form the Web", Vinceza Carchiolo, Alessandro Longheu and Michele Malgeri, in "Applied Intelligence" Vol. 18, Number 3, May 2003. ISSN:0924-669X

[Chakrabarti et al. 1999] "Focused Crawling: A New Approach to Topic-Specific Web Resource Discovery", S. Chakrabarti, M. van den Berg and B. Dom. In Proceedings of the 8th International WWW Conference, Toronto, Canada, May 1999.
http://www8.org/w8-papers/5a-search-query/crawling/index.html

[Cho et al. 1998] "Efficient Crawling Through URL Ordering", J. Cho, H. Garcia-Molina, L. Page. In Proceedings of the 7th International WWW Conference, Brisbane, Australia, April 1998. http://www7.scu.edu.au/programme/fullpapers/1919/com1919.htm

[Cho et al. 2002] "Parallel Crawler", Junghoo Cho & Hector Garcia-Molin, WWW2002, May7-11,2002. ACM 1-58113-449-5/02/0005

[Cleverdon et al. 1966] "Aslib-Cranfield research project",C. W. Cleverdon and E. M. Keen, Cranfield institute of technology, Cranfield; England, , 1966

[Crimmins 2002] "Search System Overview", Francis Crimmins, 29 May 2002.
http://dev.panopticsearch.com/system-overview.html

[F. Paradis et al. 2004] " MBOI :Scientifique et Technologique (VSST) " , Toulouse, France, 2004.

[Karger et al. 1999] "Web Caching with Consistent Hashing", by David Karger, Alex Sherman, Andy Berkheimer, Bill Bogstad, Rizwan Dhanidina, Ken Iwamoto, Brian Kim,

Luke Matkins, Yoav Yerushalmi. Published by Elsevier North-Holland, Inc. New York, NY, USA 1999. ISSN:1389-1286

[Kleinberg 1998] J. Kleinber. Authoritative Sources in a Hyperlinked Environment. SODA, 1998.

[Koster 1999] "The Web Robots Pages", M. Koster. 1999.
*http://info.webcrawler.com/mak/projects/robots/robots.html*

[Liddy 2001] "How a Search Engine Works" Elizabeth Liddy, May, 2001. FindArticles.com http://www.findarticles.com

[Lancaster,F.W 1968] " Information Retrieval Systems: Characteristics, Testing and Evaluation" published by Wiley, New York, 1968.

[McCallum et al. 1999] "Building Domain-Specific Search Engines with Machine Learning Techniques", A. McCallum, K. Nigam, J. Rennie, and K. Seymore. In 1999 AAAI Spring Symposium on Intelligent Agents in Cyberspace, Stanford University, USA, March 1999. http://www.ri.cmu.edu/pubs/pub_2716.html

[Mladenic 1999] "Text-learning and related intelligent agents" by Dunja Mladenic, in "Applications of Intelligent Information Retrieval", July-August 1999

[Oyama et al. 2003] "Domain-Specific Web Search with Keyword Spices", Satoshi Oyama, Takashi Kokubo, and Toru Ishida, Fellow, IEEE, in IEEE Transaction on knowledge and data engineering, Jan.2004, page(s):17-27,Vol. 16,Issue:1 ISSN: 1041-4347

[Page et al. 1998] "The PageRank citation ranking:Bringing order to the Web", Stanford Digital Library Technologies Project, 1998.

[Sherman et al. 2003] "The Invisible Web", Chris Sherman &Gary Price, CyberAge Books 0-910965-51-X/softbound, 2003. http://www.searchwise.net/p/iw-fla2003.pdf

[Sullivan April,2003] "Major Search Engines and Directories" Danny Sullivan in "Search Engine Watch" April, 2003 http://www.searchenginewatch.com/links/article.php/2156221

[Sullivan September, 2003] "Search Engine Sizes", Danny Sullivan, in Search Engine Watch, Sept.2003.    http://searchenginewatch.com/reports/article.php/2156481

[Thorsten 1996] "A Probabilistic Analysis of the Rocchio Algorithm with TFIDF for Text Categorization (1996)", Thorsten Joachims in proceedings of ICML-97, 14th International Conference on Machine Learning.

# Appendices

## 1. CFT Feature Words Experiment Data

### 1) CFT URLs

http://www.oshawa.ca/cit_hall/tenders.asp

http://www.city.brampton.on.ca/purchasing/rfp_2004-057.tml

http://www.city.brampton.on.ca/purchasing/rfp_2004-040.tml

http://www.city.cambridge.on.ca/cs_corporate/purchasing_tenders_list.php

http://www.city.kitchener.on.ca/tenders/tender.asp

http://www.city.london.on.ca/Purchasing/mainpg.htm

http://www.city.niagarafalls.on.ca/cityhall/qtenders.html

http://www.region.peel.on.ca/finance/purchasing/biddocs/index.htm

http://www.city.greatersudbury.on.ca/pubapps/tenders/

http://www.city.toronto.on.ca/tenders/tenders_to.htm

http://www.region.waterloo.on.ca/web/region.nsf/97dfc347666efede85256e590071a3d4/4a38
23e3515da81e85256f11004b50c4!OpenDocument

### 2) Words and Frequencies

Statistic data of all the words and its frequency in the CFT set:

| | | | |
|---|---|---|---|
| Frequency = 133, | Word = click; | | |
| Frequency = 85, | Word = list; | Frequency = 83, | Word = tender; |
| Frequency = 76, | Word = acrobat; | Frequency = 75, | Word = format; |
| Frequency = 74, | Word = pending; | Frequency = 63, | Word = ad; |
| Frequency = 62, | Word = services; | Frequency = 40, | Word = contract; |
| Frequency = 40, | Word = addendum; | Frequency = 40, | Word = request; |
| Frequency = 38, | Word = september; | Frequency = 35, | Word = contact; |
| Frequency = 35, | Word = add; | Frequency = 35, | Word = mailing; |
| Frequency = 34, | Word = august; | Frequency = 34, | Word = date; |
| Frequency = 34, | Word = proposal; | Frequency = 33, | Word = opening; |
| Frequency = 31, | Word = document; | Frequency = 30, | Word = purchasing; |
| Frequency = 26, | Word = city; | Frequency = 24, | Word = supply; |
| Frequency = 23, | Word = issued; | Frequency = 23, | Word = bid; |
| Frequency = 22, | Word = sep; | Frequency = 22, | Word = peel; |
| Frequency = 20, | Word = information | Frequency = 24, | Word = tenders; |
| Frequency = 20, | Word = closed; | Frequency = 20, | Word = works; |
| Frequency = 19, | Word = wednesday; | Frequency = 19, | Word = aug; |
| Frequency = 19, | Word = extended; | Frequency = 17, | Word = rfp; |
| Frequency = 16, | Word = quotation; | Frequency = 16, | Word = plan; |
| Frequency = 16, | Word = council; | Frequency = 16, | Word = emergency; |
| Frequency = 16, | Word = construction; | Frequency = 15, | Word = street; |

Frequency = 14, Word = deposit; Frequency = 14, Word = form;
Frequency = 14, Word = friday; Frequency = 14, Word = replacement;
Frequency = 13, Word = awarded; Frequency = 13, Word = search;
Frequency = 13, Word = download; Frequency = 13, Word = open;
Frequency = 13, Word = transportation; Frequency = 13, Word = closing;
Frequency = 13, Word = buyer; Frequency = 13, Word = home;
Frequency = 12, Word = sept; Frequency = 12, Word = proposals;
Frequency = 12, Word = division; Frequency = 12, Word = road;
Frequency = 12, Word = bids; Frequency = 12, Word = site;
Frequency = 12, Word = brampton; Frequency = 12, Word = hall;
Frequency = 12, Word = perform; Frequency = 12, Word = takers;
Frequency = 12, Word = links; Frequency = 11, Word = view;
Frequency = 11, Word = public; Frequency = 11, Word = pdf;
Frequency = 11, Word = opportunities; Frequency = 14, Word = regional;
Frequency = 11, Word = removal; Frequency = 11, Word = community;
Frequency = 11, Word = register; Frequency = 10, Word = west;
Frequency = 10, Word = centre; Frequency = 10, Word = delivery;
Frequency = 10, Word = full; Frequency = 10, Word = floor;
Frequency = 10, Word = fill; Frequency = 9, Word = officially;
Frequency = 9, Word = district; Frequency = 9, Word = charge;
Frequency = 9, Word = mail; Frequency = 9, Word = fax;
Frequency = 9, Word = department; Frequency = 9, Word = pwd;
Frequency = 9, Word = office; Frequency = 9, Word = ca;
Frequency = 9, Word = opportunity; Frequency = 9, Word = mississauga;
Frequency = 9, Word = ontario; Frequency = 9, Word = kitchener;
Frequency = 9, Word = jul; Frequency = 9, Word = business;
Frequency = 9, Word = number; Frequency = 9, Word = snow;
Frequency = 8, Word = quotations; Frequency = 8, Word = service;
Frequency = 8, Word = vendor; Frequency = 8, Word = south;
Frequency = 8, Word = police; Frequency = 8, Word = installation;
Frequency = 8, Word = documents; Frequency = 7, Word = sudbury;
Frequency = 7, Word = october; Frequency = 7, Word = summary;
Frequency = 7, Word = maps; Frequency = 7, Word = required;
Frequency = 7, Word = corporate; Frequency = 7, Word = page;
Frequency = 7, Word = process; Frequency = 7, Word = news;
Frequency = 7, Word = upgrades; Frequency = 7, Word = status;
Frequency = 7, Word = related; Frequency = 7, Word = shown;
Frequency = 7, Word = awards; Frequency = 7, Word = health;
Frequency = 6, Word = bidders; Frequency = 6, Word = wellington;
Frequency = 6, Word = development; Frequency = 6, Word = management;
Frequency = 6, Word = july; Frequency = 6, Word = approval;
Frequency = 6, Word = logo; Frequency = 6, Word = web;
Frequency = 6, Word = provide; Frequency = 6, Word = sewer;
Frequency = 6, Word = housing; Frequency = 6, Word = received;
Frequency = 6, Word = trucks; Frequency = 6, Word = sd;
Frequency = 6, Word = details; Frequency = 6, Word = water;
Frequency = 6, Word = privacy; Frequency = 6, Word = oct;
Frequency = 6, Word = axle; Frequency = 6, Word = online;
Frequency = 6, Word = map; Frequency = 6, Word = avenue;
Frequency = 6, Word = fire; Frequency = 6, Word = program;
Frequency = 6, Word = website; Frequency = 6, Word = st;
Frequency = 5, Word = posted; Frequency = 5, Word = accordance;
Frequency = 5, Word = living; Frequency = 5, Word = watermain;
Frequency = 5, Word = employment; Frequency = 5, Word = releases;
Frequency = 5, Word = cheque; Frequency = 5, Word = properties;
Frequency = 5, Word = calendar; Frequency = 5, Word = asphalt;

| | | | |
|---|---|---|---|
| Frequency = 5, | Word = canada; | Frequency = 5, | Word = questions; |
| Frequency = 8, | Word = project; | Frequency = 5, | Word = treatment; |
| Frequency = 5, | Word = email; | Frequency = 5, | Word = economic; |
| Frequency = 5, | Word = review; | Frequency = 5, | Word = drive; |
| Frequency = 4, | Word = watermains; | Frequency = 4, | Word = frequently; |
| Frequency = 4, | Word = time; | Frequency = 4, | Word = region; |
| Frequency = 4, | Word = jan; | Frequency = 4, | Word = library; |
| Frequency = 4, | Word = caledon; | Frequency = 4, | Word = results; |
| Frequency = 4, | Word = sign; | Frequency = 4, | Word = ted; |
| Frequency = 4, | Word = statement; | Frequency = 4, | Word = rental; |
| Frequency = 4, | Word = sale; | Frequency = 4, | Word = facilities; |
| Frequency = 4, | Word = refundable; | Frequency = 4, | Word = nov; |
| Frequency = 4, | Word = june; | Frequency = 4, | Word = pre; |
| Frequency = 4, | Word = property; | Frequency = 4, | Word = mar; |
| Frequency = 4, | Word = jun; | Frequency = 4, | Word = al; |
| Frequency = 4, | Word = tri; | Frequency = 4, | Word = apr; |
| Frequency = 4, | Word = improvements; | Frequency = 4, | Word = materials; |
| Frequency = 4, | Word = qualification; | Frequency = 4, | Word = winter; |
| Frequency = 4, | Word = plant; | Frequency = 4, | Word = directory; |
| Frequency = 4, | Word = events; | Frequency = 4, | Word = requests; |
| Frequency = 4, | Word = dec; | Frequency = 4, | Word = thursday; |
| Frequency = 4, | Word = municipal; | Frequency = 4, | Word = address; |
| Frequency = 4, | Word = justin; | Frequency = 4, | Word = departments; |
| Frequency = 4, | Word = rating; | Frequency = 4, | Word = complete; |
| Frequency = 4, | Word = asked; | Frequency = 4, | Word = hours; |
| Frequency = 4, | Word = building; | Frequency = 4, | Word = fuel; |
| Frequency = 4, | Word = maintenance; | Frequency = 4, | Word = wastewater; |
| Frequency = 4, | Word = ministry; | Frequency = 4, | Word = tickets; |
| Frequency = 4, | Word = system; | Frequency = 4, | Word = feb; |
| Frequency = 4, | Word = greater; | Frequency = 4, | Word = current; |
| Frequency = 4, | Word = rejected; | Frequency = 4, | Word = traffic; |
| Frequency = 4, | Word = official; | Frequency = 4, | Word = general; |
| Frequency = 4, | Word = livingston; | Frequency = 4, | Word = programs; |
| Frequency = 4, | Word = procedure; | Frequency = 4, | Word = storm; |
| Frequency = 3, | Word = phone; | Frequency = 3, | Word = fee; |
| Frequency = 3, | Word = cash; | Frequency = 3, | Word = listed; |
| Frequency = 3, | Word = call; | Frequency = 3, | Word = install; |
| Frequency = 3, | Word = mayor; | Frequency = 3, | Word = projects; |
| Frequency = 3, | Word = file; | Frequency = 3, | Word = hot; |
| Frequency = 3, | Word = copy; | Frequency = 3, | Word = include; |
| Frequency = 3, | Word = refer; | Frequency = 3, | Word = interest; |
| Frequency = 3, | Word = tr; | Frequency = 3, | Word = cold; |
| Frequency = 3, | Word = law; | Frequency = 3, | Word = square; |
| Frequency = 3, | Word = main; | Frequency = 3, | Word = computer; |
| Frequency = 3, | Word = taxation; | Frequency = 3, | Word = bridge; |
| Frequency = 3, | Word = accepted; | Frequency = 3, | Word = working; |
| Frequency = 3, | Word = grants; | Frequency = 3, | Word = souza; |
| Frequency = 3, | Word = corporation; | Frequency = 3, | Word = facility; |
| Frequency = 3, | Word = registration; | Frequency = 3, | Word = control; |
| Frequency = 3, | Word = order; | Frequency = 3, | Word = live; |
| Frequency = 3, | Word = operative; | Frequency = 3, | Word = long; |
| Frequency = 3, | Word = repairs; | Frequency = 3, | Word = social; |
| Frequency = 3, | Word = march; | Frequency = 3, | Word = submitted; |
| Frequency = 3, | Word = exterior; | Frequency = 3, | Word = heritage; |
| Frequency = 3, | Word = close; | Frequency = 3, | Word = power; |
| Frequency = 3, | Word = financial; | Frequency = 3, | Word = column; |

| Frequency = 3, | Word = telephone; | Frequency = 3, | Word = affordable; |
|---|---|---|---|
| Frequency = 3, | Word = yard; | Frequency = 3, | Word = submit; |
| Frequency = 3, | Word = local; | Frequency = 3, | Word = items; |
| Frequency = 3, | Word = vendors; | Frequency = 3, | Word = esd; |
| Frequency = 3, | Word = renovations; | Frequency = 3, | Word = intersection; |
| Frequency = 3, | Word = purchase; | Frequency = 3, | Word = arena; |
| Frequency = 3, | Word = parking; | Frequency = 3, | Word = supplies; |
| Frequency = 3, | Word = jobs; | Frequency = 3, | Word = note; |
| Frequency = 3, | Word = publications; | Frequency = 3, | Word = planning; |
| Frequency = 3, | Word = work; | Frequency = 3, | Word = fleet; |
| Frequency = 3, | Word = rehabilitation; | Frequency = 3, | Word = signed; |
| Frequency = 3, | Word = care; | Frequency = 3, | Word = citizen; |
| Frequency = 3, | Word = recreation; | Frequency = 3, | Word = advanced; |
| Frequency = 3, | Word = recycling; | Frequency = 3, | Word = quotes; |
| Frequency = 3, | Word = archives; | Frequency = 3, | Word = disclaimer; |
| Frequency = 3, | Word = systems; | Frequency = 3, | Word = certificate; |
| Frequency = 3, | Word = provided; | Frequency = 3, | Word = fair; |
| Frequency = 2, | Word = digital; | Frequency = 2, | Word = top; |
| Frequency = 2, | Word = visual; | Frequency = 2, | Word = expression; |
| Frequency = 2, | Word = operation; | Frequency = 2, | Word = payable; |
| Frequency = 2, | Word = section; | Frequency = 2, | Word = depot; |
| Frequency = 2, | Word = trunk; | Frequency = 2, | Word = provision; |
| Frequency = 2, | Word = times; | Frequency = 2, | Word = child; |
| Frequency = 2, | Word = diesel; | Frequency = 2, | Word = supplier; |
| Frequency = 2, | Word = owned; | Frequency = 2, | Word = manor; |
| Frequency = 2, | Word = decisions; | Frequency = 2, | Word = thedomain; |
| Frequency = 2, | Word = biodiesel; | Frequency = 2, | Word = prequalification; |
| Frequency = 2, | Word = staff; | Frequency = 2, | Word = finance; |
| Frequency = 2, | Word = sports; | Frequency = 2, | Word = visitors; |
| Frequency = 2, | Word = important; | Frequency = 2, | Word = schedule; |
| Frequency = 2, | Word = receiving; | Frequency = 2, | Word = hamilton; |
| Frequency = 2, | Word = onmousedown; | Frequency = 2, | Word = tom; |
| Frequency = 2, | Word = nile; | Frequency = 2, | Word = animal; |
| Frequency = 2, | Word = groups; | Frequency = 2, | Word = data; |
| Frequency = 2, | Word = wards; | Frequency = 2, | Word = sidewalks; |
| Frequency = 2, | Word = guide; | Frequency = 2, | Word = planting; |
| Frequency = 2, | Word = reserved; | Frequency = 2, | Word = roadwatch; |
| Frequency = 2, | Word = streets; | Frequency = 2, | Word = plowing; |
| Frequency = 2, | Word = qualifications; | Frequency = 2, | Word = contents; |
| Frequency = 2, | Word = responsibility; | Frequency = 2, | Word = interior; |
| Frequency = 2, | Word = calling; | Frequency = 2, | Word = east; |
| Frequency = 2, | Word = accommodations; | Frequency = 2, | Word = rights; |
| Frequency = 2, | Word = restoration; | Frequency = 2, | Word = shepherd; |
| Frequency = 2, | Word = gm; | Frequency = 2, | Word = solicitation; |
| Frequency = 2, | Word = technology; | Frequency = 2, | Word = agent; |
| Frequency = 2, | Word = français; | Frequency = 2, | Word = mapping; |
| Frequency = 2, | Word = roads; | Frequency = 2, | Word = successes; |
| Frequency = 2, | Word = event; | Frequency = 2, | Word = bidding; |
| Frequency = 2, | Word = wage; | Frequency = 2, | Word = ed; |
| Frequency = 2, | Word = originally; | Frequency = 2, | Word = prices; |
| Frequency = 2, | Word = title; | Frequency = 2, | Word = transhelp; |
| Frequency = 2, | Word = specific; | Frequency = 2, | Word = check; |
| Frequency = 2, | Word = mix; | Frequency = 2, | Word = certified; |
| Frequency = 2, | Word = dump; | Frequency = 2, | Word = marked; |
| Frequency = 2, | Word = measures; | Frequency = 2, | Word = human; |
| Frequency = 2, | Word = tt; | Frequency = 2, | Word = insurance; |

| | | | |
|---|---|---|---|
| Frequency = 2, | Word = oxford; | Frequency = 2, | Word = newspaper; |
| Frequency = 2, | Word = capreol; | Frequency = 2, | Word = formal; |
| Frequency = 2, | Word = rec; | Frequency = 2, | Word = arts; |
| Frequency = 2, | Word = moonah; | Frequency = 2, | Word = minutes; |
| Frequency = 2, | Word = lakeview; | Frequency = 2, | Word = transit; |
| Frequency = 2, | Word = tax; | Frequency = 2, | Word = require; |
| Frequency = 2, | Word = reference; | Frequency = 2, | Word = king; |
| Frequency = 2, | Word = celebrate; | Frequency = 2, | Word = storage; |
| Frequency = 2, | Word = architectural; | Frequency = 2, | Word = entrance; |
| Frequency = 2, | Word = submission; | Frequency = 2, | Word = international; |
| Frequency = 2, | Word = dundas; | Frequency = 2, | Word = hardware; |
| Frequency = 2, | Word = davies; | Frequency = 2, | Word = term; |
| Frequency = 2, | Word = consulting; | Frequency = 2, | Word = faq; |
| Frequency = 2, | Word = clarkson; | Frequency = 2, | Word = line; |
| Frequency = 2, | Word = lowest; | Frequency = 2, | Word = copyright; |
| Frequency = 2, | Word = application; | Frequency = 2, | Word = irene; |
| Frequency = 2, | Word = downtown; | Frequency = 2, | Word = envelope; |
| Frequency = 2, | Word = tandem; | Frequency = 2, | Word = attractions; |
| Frequency = 2, | Word = pioneer; | Frequency = 2, | Word = notices; |
| Frequency = 2, | Word = instructions; | Frequency = 2, | Word = make; |
| Frequency = 2, | Word = pick; | Frequency = 2, | Word = initiatives; |
| Frequency = 2, | Word = toronto; | Frequency = 2, | Word = landscape; |
| Frequency = 2, | Word = murphy; | Frequency = 2, | Word = joint; |
| Frequency = 2, | Word = performing; | Frequency = 2, | Word = hss; |
| Frequency = 2, | Word = laws; | Frequency = 2, | Word = hill; |
| Frequency = 2, | Word = sealed; | Frequency = 2, | Word = statements; |
| Frequency = 2, | Word = software; | Frequency = 2, | Word = advertised; |
| Frequency = 2, | Word = virus; | Frequency = 2, | Word = monday; |
| Frequency = 2, | Word = alterations; | Frequency = 2, | Word = windows; |
| Frequency = 2, | Word = high; | Frequency = 2, | Word = bidder; |
| Frequency = 2, | Word = policy; | Frequency = 2, | Word = place; |
| Frequency = 2, | Word = media; | Frequency = 2, | Word = wp; |
| Frequency = 2, | Word = cls; | Frequency = 2, | Word = theatre; |
| Frequency = 2, | Word = shane; | Frequency = 2, | Word = electrical; |
| Frequency = 2, | Word = training; | Frequency = 2, | Word = asphaltic; |
| Frequency = 2, | Word = competitive; | Frequency = 2, | Word = agendas; |
| Frequency = 2, | Word = area; | Frequency = 2, | Word = theemail; |
| Frequency = 2, | Word = mechanical; | Frequency = 2, | Word = bloom; |
| Frequency = 2, | Word = box; | Frequency = 2, | Word = communities; |
| Frequency = 2, | Word = funds; | Frequency = 2, | Word = resources; |
| Frequency = 2, | Word = cao; | Frequency = 2, | Word = checkbutton; |
| Frequency = 2, | Word = landfill; | Frequency = 2, | Word = phase; |
| Frequency = 2, | Word = limited; | Frequency = 2, | Word = tourism; |
| Frequency = 2, | Word = technical; | Frequency = 2, | Word = mentioned; |
| Frequency = 2, | Word = images; | Frequency = 2, | Word = reports; |
| Frequency = 2, | Word = surplus; | Frequency = 2, | Word = hospital; |
| Frequency = 2, | Word = opened; | Frequency = 2, | Word = peelregion; |
| Frequency = 2, | Word = land; | Frequency = 2, | Word = design; |
| Frequency = 2, | Word = pollution; | Frequency = 2, | Word = mailed; |
| Frequency = 2, | Word = dearness; | Frequency = 2, | Word = courier; |
| Frequency = 2, | Word = company; | Frequency = 2, | Word = december; |
| Frequency = 2, | Word = porcarelli; | Frequency = 2, | Word = entry; |
| Frequency = 2, | Word = reconstruction; | Frequency = 2, | Word = leisure; |
| Frequency = 2, | Word = aluminum; | Frequency = 2, | Word = signage; |
| Frequency = 2, | Word = seasonal; | Frequency = 2, | Word = procedures; |
| Frequency = 2, | Word = canadian; | Frequency = 2, | Word = parks; |

Frequency = 2, Word = gas; Frequency = 2, Word = calls;
Frequency = 2, Word = gr; Frequency = 2, Word = master;
Frequency = 2, Word = reader; Frequency = 2, Word = necessarily;
Frequency = 2, Word = sand; Frequency = 2, Word = rfps;
Frequency = 2, Word = tourists; Frequency = 1, Word = tracing;
Frequency = 1, Word = wl; Frequency = 1, Word = tw;
Frequency = 1, Word = budgets; Frequency = 1, Word = notified;
Frequency = 1, Word = paid; Frequency = 1, Word = taxes;
Frequency = 1, Word = purchasingpolicy; Frequency = 1, Word = village;
Frequency = 1, Word = greatersudbury; Frequency = 1, Word = shelter;
Frequency = 1, Word = energy; Frequency = 1, Word = publicly;
Frequency = 1, Word = tool; Frequency = 1, Word = explorer;
Frequency = 1, Word = location; Frequency = 1, Word = body;
Frequency = 1, Word = fail; Frequency = 1, Word = chief;
Frequency = 1, Word = extension; Frequency = 1, Word = administrative;
Frequency = 1, Word = located; Frequency = 1, Word = sunforest;
Frequency = 1, Word = wheelchair; Frequency = 1, Word = easan;
Frequency = 1, Word = terrace; Frequency = 1, Word = lookout;
Frequency = 1, Word = bank; Frequency = 1, Word = pump;
Frequency = 1, Word = oshawa; Frequency = 1, Word = de;
Frequency = 1, Word = residents; Frequency = 1, Word = singles;
Frequency = 1, Word = railway; Frequency = 1, Word = fiber;
Frequency = 1, Word = accessibility; Frequency = 1, Word = urb;
Frequency = 1, Word = speed; Frequency = 1, Word = core;
Frequency = 1, Word = courts; Frequency = 1, Word = norris;
Frequency = 1, Word = assessment; Frequency = 1, Word = fees;
Frequency = 1, Word = comment; Frequency = 1, Word = surina;
Frequency = 1, Word = pay; Frequency = 1, Word = viewed;
Frequency = 1, Word = incomplete; Frequency = 1, Word = server;
Frequency = 1, Word = introduction; Frequency = 1, Word = bovaird;
Frequency = 1, Word = installed; Frequency = 1, Word = 0dz7udyj4;
Frequency = 1, Word = untitled; Frequency = 1, Word = audit;
Frequency = 1, Word = civic; Frequency = 1, Word = sales;
Frequency = 1, Word = trades; Frequency = 1, Word = consultant;
Frequency = 1, Word = enterprises; Frequency = 1, Word = playground;
Frequency = 1, Word = symbol; Frequency = 1, Word = acceptable;
Frequency = 1, Word = cable; Frequency = 1, Word = cellular;
Frequency = 1, Word = assorted; Frequency = 1, Word = finlayson;
Frequency = 1, Word = middleton; Frequency = 1, Word = abc;
Frequency = 1, Word = steeles; Frequency = 1, Word = parts;
Frequency = 1, Word = curb; Frequency = 1, Word = valves;
Frequency = 1, Word = heat; Frequency = 1, Word = caradoc;
Frequency = 1, Word = environment; Frequency = 1, Word = utility;
Frequency = 1, Word = apparatus; Frequency = 1, Word = successful;
Frequency = 1, Word = north; Frequency = 1, Word = internet;
Frequency = 1, Word = cancelled; Frequency = 1, Word = jalna;
Frequency = 1, Word = vehicles; Frequency = 1, Word = industry;
Frequency = 1, Word = ew; Frequency = 1, Word = price;
Frequency = 1, Word = visa; Frequency = 1, Word = logging;
Frequency = 1, Word = send; Frequency = 1, Word = don;
Frequency = 1, Word = archive; Frequency = 1, Word = valley;
Frequency = 1, Word = airport; Frequency = 1, Word = crime;
Frequency = 1, Word = leo; Frequency = 1, Word = rob;
Frequency = 1, Word = language; Frequency = 1, Word = electronic;
Frequency = 1, Word = ski; Frequency = 1, Word = roof;
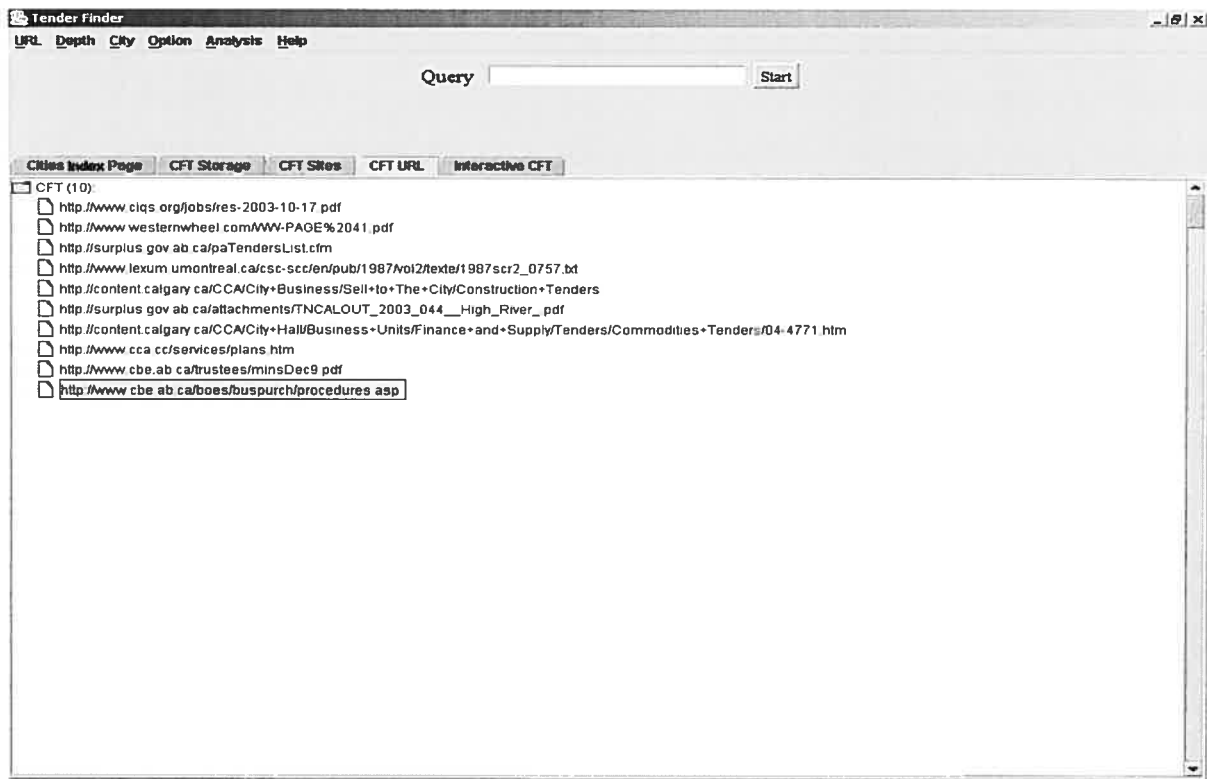Frequency = 1, Word = marriage; Frequency = 1, Word = outfall;

Frequency = 1,　　Word = therapy;　　　　Frequency = 1,　　Word = walden;
Frequency = 1,　　Word = workplace;

## 2.　CFT Searching in Five Cities

(Yes: CFT pertinent; No: not a CFT page; Crawler depth:2; Page:2)

**Calgary**

http://www.ciqs.org/jobs/res-2003-10-17.pdf　No

http://www.westernwheel.com/WW-PAGE%2041.pdf　Yes

http://surplus.gov.ab.ca/paTendersList.cfm Yes

http://www.lexum.umontreal.ca/csc-scc/en/pub/1987/vol2/texte/1987scr2_0757.txt No

http://content.calgary.ca/CCA/City+Business/Sell+to+The+City/Construction+Tenders Yes

http://surplus.gov.ab.ca/attachments/TNCALOUT_2003_044__High_River_.pdf　Yes

http://content.calgary.ca/CCA/City+Hall/Business+Units/Finance+and+Supply/Tenders/Commodities+Tenders/04-4771.htm Yes

http://www.cca.cc/services/plans.htm Yes

http://www.cbe.ab.ca/trustees/minsDec9.pdf

http://www.cbe.ab.ca/boes/buspurch/procedures.asp Yes



Appendix Figure　*Test in Calgary*

## 2) Vancouver

http://www.translink.bc.ca/About_TransLink/Business_Opportunities/Bid_Details/Q4-0014.asp Yes

http://www.crd.bc.ca/tenders/041510_trenchless_pipe_rehab.pdf Yes

http://www.bchydro.com/rx_files/info/info13924.pdf Yes

http://www.gvrd.bc.ca/gvrdtenders/TenderDetails.aspx?tenderId=274&tenderType=open Yes

http://www.gvrd.bc.ca/gvrdtenders/TenderDetails.aspx?tenderId=279&tenderType=open Yes

http://www.grainscanada.gc.ca/newsroom/misc/tenders/surplusvan-e.htm Yes

http://www.dnv.org/upload/documents/cpolicy/c209501.pdf Yes

http://www.bestwesterncoqinn.com/pantrylunch.html No

http://www.bchydro.com/rx_files/info/info9912.pdf No

http://www.city.vancouver.bc.ca/ctyclerk/ADS/ADcontract2004-01.htm Yes

http://www.885westgeorgia.com/extranet/news_emp.html No

http://www.city.richmond.bc.ca/webnews/tenders/docs/0818_park1570.pdf Yes

http://www.translink.bc.ca/About_TransLink/Business_Opportunities/default.asp Yes

## 3) Winnipeg

http://www.mpi.mb.ca/salvage/auctionprint.asp?salenm=156 Yes

http://www.winnipeg.ca/matmgt Yes

http://www.mpi.mb.ca/salvage/auctionprint.asp?salenm=167 Yes

http://www.winnipeg.ca/finance/findata/matmgt/documents/2004/196-2004/196-2004_Part_A-Submission.pdf Yes

http://www.prha.mb.ca/news.htm CFT News No.

http://www.mhca.mb.ca/HeavyNews/heavy_news_june3_04.pdf No

http://www.winnipeg-chamber.com/pdf/policy/fldsol.pdf Yes

http://www.smartwinnipeg.mb.ca/Smart_Recommendations.htm CFT News, No

http://www.cwb.ca/en/publications/farmers/pdf/nov-dec-2003.pdf No

http://www.gov.mb.ca/tgs/gs/contracts/tenders/ads/1198_04_ad.pdf Yes

http://www.wsd1.org/Board/policy_minutes/Policy/Policy_DJC.pdf Yes

http://www.cwb.ca/en/publications/students_researchers/pdf/2002-03_full_english_statistics.pdf    No

http://www.cwb.ca/en/movement/elevator_managers/pdf/2004-05_elevatorguide.pdf    No

http://www.mhca.mb.ca/HeavyNews/heavy_news_april29.pdf Yes

http://www.canoe.ca/NewsStand/Columnists/Winnipeg/Frank_Landry/2004/07/30/562424.html No

http://www.canadinns.com/confbanq/ambmenuwin.pdf    No

http://www.gov.mb.ca/tgs/gs/contracts/tenders/ads/0212sr04_ad.pdf Yes


**3) Yellowknife**

http://www.canlii.org/nt/laws/regu/f-3/20040512/whole.html    Yes

http://www.recycle.ab.ca/2003Proceedings/JaredBuchko.pdf    No

http://www.amec.com/earthandenvironmental/where/Officeprofiles.asp?PageID=1081&OfficeID=88    No

http://www.ozemail.com.au/~kshapley/flash.html No

http://city.yellowknife.nt.ca/userfiles/page_attachments/Library/1/358200_LAND_ADMINISTRATION_BY-LAW_NO_3853.pdf Yes

http://www.diavik.ca/dialogue/dialogue15.pdf    No

http://collection.nlc-bnc.ca/100/201/300/first_perspective/2001/05-10/employment/emp1050102.html    Yes

http://www.yukonweb.com/community/yukon-news/apr5.htmld Yes

http://www.wcb.nt.ca/AboutWCB/pdf/Goverance%20Council/Dec2002.pdf No

http://city.yellowknife.nt.ca/userfiles/page_attachments/Library/1/3143357_CAPITAL_UPDATE_JULY_30__2004_-_version_2.pdf Yes

http://www.pwgsc.gc.ca/reports/text/pdfs/DPR_PWGSC_2003_e.pdf    Yes

http://www.nnsl.com/yir/yir01/yirdehcho01.html No

http://www.cbsc.org/english/search/display.cfm?code=4066&coll=FE_FEDSBIS_E    Yes

**4) Halifax**

http://www.hrwc.ns.ca/whats_new/tenders.html Yes

http://www.bids.ca/protect/tndr290.htm Yes

http://www.herald.ns.ca/stories/2004/09/07/fEntertainment127.raw.html No

http://www.region.halifax.ns.ca/procurement/tenders3.asp Yes

http://www.novaforestalliance.com/media/documents/tender_offer_Mar26-02.pdf Yes

http://www.acoa.ca/e/media/press/press.shtml?1401 Yes

http://www.nseia.ns.ca/EnviroNews%20%20March%2026,%202004.pdf No

http://www.cbsc.org/english/search/display.cfm?code=4066&coll=FE_FEDSBIS_E
      Yes

http://www.bids.ca/protect/tndr670.htm Yes


**5) Ajax**

http://www.georgeortiz.com/3D/220 NO

http://www.ipc.on.ca/scripts/index_.asp?action=31&P_ID=7589&N_ID=1&PT_ID=9
23&U_ID=0 Yes

http://www.townofajax.com/government/minutes/2004/gg-july8.pdf No

http://www.trainandemploy.qld.gov.au/client/resources/about/research_publications/st
rategy_policy/pdf/example_plan.pdf Yes

http://www.bids.ca/tenderlist_agency.php?AGENCY=AJAX Yes

http://news.moneycentral.msn.com/ticker/sigdev.asp?Symbol=IBN CFT report No

http://www.rodney.govt.nz/council/minutes/fullcouncilsuppagenda29JAN04.pdf Yes

http://news.bbc.co.uk/sport/hi/english/football/champions_league/newsid_1449000/14
49507.stm No