ANL/DIS/CP--84374 Conb-9410141--4

Raising the IQ in Full-Text Searching via Intelligent Querying

Robert Kero, Lucian Russell, Craig Swietlik Argonne National Laboratory

Charles Morgan DOE/Environment, Safety & Health (EH-5)

Charlynn Clayton, Kelly Dunlap DOE/Office of Scientific and Technical Information

Abstract

Current Information Retrieval (IR) technologies allow for efficient access to relevant information, provided that user selected query terms coincide with the specific linguistical choices made by the authors whose works constitute the text-base. Therefore, the challenge is to enhance the limited searching capability of state-of-the-practice IR. This can be done either with augmented clients that overcome current server searching deficiencies, or with added capabilities that can augment searching algorithms on the servers. The technology being investigated is that of deductive databases, with a set of new techniques called cooperative answering. This technology utilizes semantic networks to allow for navigation between possible query search term alternatives. The augmented search terms are passed to an IR engine and the results can be compared. The project utilizes the OSTI Environment, Safety and Health Thesaurus to populate the domain specific semantic network and the text base of ES&H related documents from the Facility Profile Information Management System as the domain specific search space.

DISCLAIMER

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof.

MASTER

DISTRIBUTION OF THIS DOCUMENT IS UNLIMITED



The submitted manuscript has been authored by a contractor of the U.S. Government under contract. No. W-31-109-ENG-38 Accordingly, the U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of the contribution, or allow others to do so, for the submitted form of the contribution, or allow others to do so, for

.

DISCLAIMER

Portions of this document may be illegible in electronic image products. Images are produced from the best available original document.

Introduction

The Department of Energy, Office of Special Projects and Argonne National Laboratory (ANL), along with the Department of Energy, Office of Scientific and Technical Information (OSTI) have designed and implemented the Environment, Safety and Health Facility Profile Information Management System (FPIMS), which is a personal computer based information retrieval system. This system is the central repository for Department of Energy (DOE) oversight documents consisting of Baseline Compliance Assessments, Action Plans, follow-on Progress Assessments, site specific information (e.g., maps, fact sheets, profiles), etc.

The system itself is highly modular, consisting of two major components. The first component is the graphical user interface which displays a map of the United States with icons representing field offices and associated DOE sites that have been assessed. The documents available for each site are displayed in an intuitive bookcase, bookshelf and book level progression of screens. The second component is a robust, commercial information retrieval engine used to provide powerful searching capabilities over the many user-selectable sets of documents housed within the system.

The information retrieval engine, (Personal Library Software's PC LAN-based Personal Librarian), found within the original FPIMS facilitates searches on collections of documents by providing standard Boolean logic (i.e., AND, OR, NOT) as well as natural language style searching capabilities. These capabilities are further supplemented with the ability to dynamically formulate concept-based searches. These capabilities were provided initially in a stand-alone system, and then were migrated to a Local Area Network (LAN) [MORG93]. Advancing technologies on the Internet such as the World Wide Web and the emergence of an advanced commercial information retrieval product, Web-based Personal Librarian, have created new cost effective opportunities for disseminating important DOE information, as well as providing the tools necessary to support the migration of the original FPIMS to a Web-based FPIMS. This opportunity is currently being addressed.

With the opportunity, however, comes a challenge. As the number of documents in FPIMS grows and the number of potential users increases it becomes ever more a challenge to ensure effective use of the corpus of information. With this in mind Argonne National Labs has embarked on a project to raise the IQ capability of the FPIMS interface.

Motivation/Incentive

A major challenge is to provide all user classes a natural and easy means of using the FPIMS. The primary motivations are to improve specifically EH's, and broadly both DOE's and the public's accessibility to the data and to consolidate the information management needs of EH. This is consonant with EH's proactive management approach. It is also an initiative as called for in DOE's Mission Statement and the administration's National Information Infrastructure (NII). Furthermore, the proposed use of the SGML

standards, as a neutral format in which to save all EH documents is in full accord with the DOE mandate and requirement. Also, this approach allows for the exploitation of emerging advanced public domain technology, such as (Wide Area Information Servers (WAIS), World Wide Web (Web) servers and Mosaic clients), to fully leverage prior investments in FPIMS.

One technical challenge in meeting FPIMS dissemination and usage goals is communications. The migration of the FPIMS to the Internet through the use of the Web allows DOE environment, safety and health data to be also shared with outside organizations such as local, state and federal agencies, as well as the public. Hence, remote DOE operations offices and DOE sites will be provided a more consistent means of access to the information via the Internet instead of using slow phone connections to the original system's local area network. With this technical infrastructure in place, the attention is directed to the system and its potential for improving productivity and effectiveness.

From the perspective of EH, one of the purposes of FPIMS was to assist the DOE in planning future activities which might require access to ES&H documents. The idea was that any future activity of the DOE involving ES&H matters could be a less expensive and higher quality effort if plans could be made on the basis of already completed surveys and assessments. One example: a simplification of interactions with the field. Headquarters personnel could review all on-file records of interactions with the field offices regarding any specific area they are interested in. With this information in hand, activities involving headquarters and field offices might be conducted by telephone instead of needing in-person visits. In the same vein, new activities could eliminate duplicative effort by building upon already known and completed activities. Another clear example is the advantage of having a corpus of material that can be read by new employees, whether federal or contractor. FPIMS' documents, and the efficient means of locating specific areas of interest, can make these personnel much more effective in a shorter time frame.

The key to realizing this promise of increased efficiency is integrating the use of FPIMS into all relevant activities. In addition to the necessary conditions for this to take place (e.g. communications and software), there is one overriding condition: the government employee must be assured that the *time* spent interfacing with the system is *cost effective*. This means that the system, must provide a high degree of assurance that the necessary information will be found *without* the need to review large numbers of irrelevant documents. If this is so, the employee will continue to gain trust in the FPIMS knowing that its use is an efficient investment of their time. Any factors which mitigate this trust must be identified and eliminated.

One historical factor that mitigates the usefulness of any Information Retrieval system, is the difficulty in formulating search strategies which will obtain the *right resultant set* of documents. The Federal employee who must investigate ES&H matters is a professional in a related discipline, not necessarily a specialist in document searching. Therefore the exposure to FPIMS is likely to be intermittent and casual. Such exposure

puts a burden on the system that it contain aids for quickly navigating within the morass of detail that any such system must necessarily contain. To lift this burden, it appeared that the introduction of Cooperative Answering technology might be of assistance, and so experiments were conducted to test this hypothesis.

Advancing Efficient Full-Text Searching through Cooperative Query Construction

Locating specific full-text information throughout DOE increases in difficulty as the information continues to grow in size and complexity. Hypothetically speaking, a particular technology (which on average retrieves 7 irrelevant documents for every 3 relevant ones), is rendered unusable when document collections grow and precision to recall ratios remains static (30 %). This information explosion problem therefore makes it unreasonable to expect that users possess enough knowledge about existing information and its organization to accurately compose useful queries consistently - the burden must be shared by the information system they are interacting with!

Current Information Retrieval (IR) technologies do very little to facilitate this problem since they only allow efficient access to "relevant information", provided user selected query terms coincide with the specific linguistical choices made by the authors whose works constitute the corpus being searched. For casual or intermittent users of the system this cannot be guaranteed: they cannot be expected to know the context in which documents were indexed, categorized, or the terminology of the specialists who generated the documents. Cooperative Answering strategies assist these users by providing access to the conceptual framework of the specialists who created the documents.

Therefore by accessing the context used by the authors of the documents in FPIMS, the IQ software can guide the user during query construction. The benefit of this approach is that it permits earlier detection of poorly composed queries - e.g., queries that (1) are to general, or (2) contain misconceptions, or (3) are not linguistically equivalent to that of the corpus. The system can provide feedback and suggestions of alternative words that will make the queries more effective.

The proof-of-concept IQ system currently available utilizes the Department of Energy's Environment, Safety and Health thesaurus/dictionary (which includes taxonomies - broader/narrower, synonymous, and meronymous terms and phrases), to populate a domain specific semantic network. This semantic network is used to navigate between possible query search term alternatives. The domain specific HTML-based corpus is a subset of the FPIMS' documents.

The System Overview

The problem, as shown in Figure 1, is to navigate the maze of documents so as to locate the most likely required document subset for which to peruse - this subset could be found distributed across multiple IR engines. An integral part of the puzzle is the IQ-based

FPIMS user interface. By using a visual interface metaphor, of a map and bookcases located at various sites on the map, it is possible to provide user-selectable control over the total FPIMS document set. Therefore, the overall query space can be initially constrained by site, by document type or by combinations of both. The intelligence of the system, is found in the Semantic Net, a logical database of terms having one or more relationships to one another. This set of relationships contains at its core the terms developed by the OSTI team that created the DOE ES&H Thesaurus. The IQ software uses this semantic net to supply the user with possible new terms which might make their searches more effective. The following Figures are taken from the proof-of-concept IQ system; they are provided to allow for a better understanding of the process.

The example takes a compound search, "arm or joint" as the starting point as shown in Figure 2. In the original test there were over 100 records (portions of documents) which satisfied the original query - too much to peruse quickly. (This situation is typical in large document collections.) Note that the screen identifies the server as a Wide Area Information Server (WAIS) node. This is not necessary for IQ searching to be of value - any type of search engine could be substituted.

The next screen (Figure 3) shows how to make use of the Cooperative Answering capability - made narrower. The OSTI Thesaurus was searched for which words in the query could be changed, made either more general or more restricted or related in meaning, to refine the query's semantics. In the example shown both words could be "relaxed" - replaced by other related terms.

The system then displays (Figure 4) the choices which can be made. The word "joint" was relaxed through a narrower interpretation, and the 7 types of specific joints were provided for possible selection and subsequent query term replacement.

The selection of "shoulder" was made (Figure 5) and the resulting query (Figure 6) revealed that 21 possible documents were found, a savings of 500% in document perusal time!

This simple example shows that the incorporation of the IQ system can greatly improve the querying process.

Benefits and Value for DOE

The concept of using a discipline-specific thesaurus to assist a DOE user in formulating and enhancing queries, effectively leverages the experience of the subject matter experts in various disciplines for the refinement of the users query. The OSTI expertise invested in the thesaurus can now be used as an integral part of IQ-based information systems.

The IQ technology makes it possible for DOE and public users to access highly specific technical information within large document collections in an efficient and highly cost effective manner. Coupled with advances in Internet technologies, the dissemination of

the FPIMS' environment, safety and health oversight information is helping to achieve the DOE's openness initiatives. This technology enables DOE to work better and cost less.

Current and Future Research Directions

Although Cooperative Answering strategies have been successfully incorporated with deductive and relational databases in the past, our research will now investigate its incorporation with information retrieval. The overall goal of the project is to provide a more efficient user class independent (e.g., domain experts through the general public) query construction technology.

Acknowledgments

The authors wish to acknowledge the contributions of the organizations and staff who have been involved in the development of the systems. Department of Energy: Mark Gilbertson. Argonne National Laboratory: Michelle Bernard, Terry Gaasterland, Robert Haddad, Julia Lee, Wes Maciorowski, Michele McCusker-Whiting. Personal Librarian Software: Deborah Senay, Paul Campbell. The development of these FPIMS-based systems and experiments represent a close collaboration of Government, Department of Energy, National Laboratories and the commercial sector.

Work supported by the U.S. Department of Energy, Assistant Secretary for Environment, Safety and Health, Office of Special Projects under Contract W-31-109-Eng-38.

References

[AND93]

Marc Andreesan: NCSA Mosaic Technical Summary, NCSA Technical Document, NCSA Springfield, Champaign II, May 8, 1993.

[CLA93]

Charlynn Clayton, et. al: INFOTECH'93, "WAIS EVALUATION PROJECT REPORT", Office of Scientific and Technical Information U.S. Department of Energy, DOE Technical Information Meeting, Oak Ridge, Tennessee, Oct. 20-22, 1993.

[GAAS92]

Theresa Gaasterland, "Generating Cooperative Answers In Deductive Databases", University of Maryland, August 1992.

[KERO94]

Kero, R., C. Swietlik, Migrating the Facility Profile Information Management

System into the World Wide Web, DOE 11th Office Information Technology Conference, Chicago, Aug. 23-25, 1994.

[MORG93]

Morgan, C., C. Swietlik, B. Kissinger, Environmental, Safety and Health Facility Profile Information Management System (FPIMS), INFOTECH'93, Office of Scientific and Technical Information U.S. Department of Energy, DOE Technical Information Meeting, Oak Ridge, Tennessee, Oct. 20-22, 1993.

[PIT94]

James E. Pitkow & Mimi Recker: First International Conference on the World Wide Web, "RESULTS FROM THE FIRST WORLD WIDE WEB USER SURVEY", College of Computing, Georgia Institute of Technology, Atlanta, GA, May 25-27, 1994.

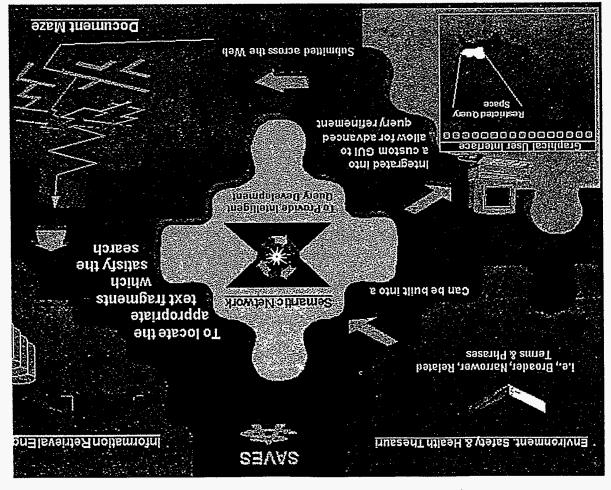


Figure 1 Navigating the Maze

ad Thesaurus	8800186	Clear Screen Help Information	gur
i etter Listerior Trent Query:	A amaginas servicios		
when seems to all a seems of a seems had	CONTRACTOR OF THE PARTY OF THE		-
	relax/2 are not together erms for queries are: o, operations_offices, o	coolants, monitors, offices, doe	· ·
Some sample sara, acts		coolants, monitors, offices, doe	· ·

انم

٠.

Φ.	WORKING.GU	in your control of the property of the party				na ma aménic anis anis	
		WELCOME to the	Cooperative \	VAIS Query	Build	er ;	
Loa	ad Thesaurus	oxuldet	Clear	Screen		<u>H</u> elp Informati	on CLUM
11		Select Direction of	Relaxation	Votienco	HE)		
			o temperature and make				
Cui	rent Query:	am foibles					
•	yes						
	?- split_qi	ery_set_array("arm,j	oint")	: ;			. ,
	7- relax qu Terms in	ery("arm, joint") current query are:		. :			;
	Enter Ch	1) arm 2) joint ice (1-2) in Input W	indow.	. ,.		,	, , , , , , , , , , , , , , , , , , ,
¥	yes			ne nedecominant in descrip	-inninnaire	**********	es property property and the second
	រប់រួចឥរ្រៈបវរស្រុន	nuncil y essessioni	alei (bligini) yy	avale ju	Y.D.		Walianiji ak

i

WORKING, GUITES	and the state of t			
	WELCOME to the Co	operative WAIS Query	Builder	
Load Thesaurus	balldat	Clear Screen	<u>H</u> elp Information	onic
	Select Direction of Rela	axation ↑ UELTON \$15,0000 \$10,0000 \$10,0000		
Current Query:	rm je interes en			
(1) (2) (3) (4) (5) (6)	n(2), relax current to rower words for >>>> sibilities are ankle elbow hip knee knuckle shoulder wrist e (1-7) in Input Wind			
yes ?-			one many the branch to the court of the	reament w
ECHCOSEWAISISON	Estat assess are (SA	in color devaled \$4	enpiradities sieve	lighining (

.

[o] WORKING.GUItext ○	※四 : . ·		ř.	•• .:
	VELCOME to the Coope	erative WAIS Query	Builder	
Load Thesaurus		Clear Screen		nation <u>Cluid</u>
Current Query:				
A : ?- select new word arm, shoulder				
•			***	
rciposeyas someth	005080000 SE000	conference Water	a pacononis	MEW Political

Controlled to the second of th

Xedit	it Save Load Foo.doclist Use Control-S and Control-R to Search.	le foo.doclist opened read - write.	foo.doclist Read - Write	D BY WAIS SEARCH	To select a document, click on "VIEW DUCUMENIS"	Amber Score Headline 1000 HS_CONCERN_FS_3_2_H1_C2_At_Westi 521 INEL_CONCERN_MA_3_2_H3_C1_The_EG 312 HS_CONCERN_MA_3_2_H3_C1_The_EG 312 HS_CONCERN_FS_2_2_H2_C1_At_Depart 195 HS_CONCERN_FS_3_3_H2_C2_Mestingh 145 HS_CONCERN_FS_3_3_H2_C1_At_Westi 129 INEL_CONCERN_FS_3_3_H2_C1_At_Westi 120 INEL_CONCERN_FS_3_3_H2_C1_At_Depart 120 HS_4_5_19_1_Overview 120 HS_4_5_19_1_Overview 120 HS_4_5_19_1_Overview 1210 HS_4_5_19_1_Overview 1229 HS_4_5_19_1_Overview 1220 HS_4_5_19_1_Overview 1220 HS_4_5_3_14_1_Overview 1220 HS_6_5_5_6_H1_SHMARY 1220 HS_6_5_5_6_H1_SHM
কি xedit		foo				