

Context Based Diversification for Keyword Queries over XML Data Using Improved Information Retrieval Query Model

¹Saylee Suradkar, ²Mrs. S. A. Kinariwala, ³Mr. R. B. Patil

Abstract: Although keyword search terms can help people find a lot of information, the ambiguity of the keyword search query makes search query responses more efficient, especially for short keywords and obscure. By solving this challenge, in this article, we will propose ways to make it easier to find contextual XML keywords in XML data, short and vague keyword queries, and XML data. It comes from a search query using a simple property selection style. We then designed a powerful XML keyword analysis model to measure the quality of each candidate. Then propose two algorithms to determine the highest search order for the researchers. Two selection criteria are targeted: K. The selected research participant will be the most relevant for the search term while the maximum number of results is covered. At the end of the detailed evaluation, real data and summary data, we demonstrate the effectiveness of the proposed distribution model and our algorithm.

Index Terms: XML keyword search, context-based diversification

I. INTRODUCTION

All Keyword research in data structures and semi-structured structures has attracted a lot of research because users can retrieve data without having to learn the query language and the complex structure of the database looking for keywords in the search (IR) to search for related documents. Structures and semi-structures (such as DB and IR) will focus on specific content such as the smaller root part. Node v is considered a SLCA. If 1) the subtree that the root of this node contains all the keywords and 2, no successor node of the v_0 node of v is nested with v_0 with all the keywords. In other words, if the node is a SLCA, its ancestors will be completely separated from SLCA. SLCA can be used to identify specific results in the search for XML keywords. In this article, we also use the SLCA semantics. [2], [3], [4], [5] keyword search via XML data. In general, keywords with user search terms are easier to identify user search functions search is easy to find.

However, when you search for a specific keyword with a slightly uncertain keyword, it becomes very difficult to determine the user's search intent because of the ambiguity of keyword searches. Although it is sometimes useful to involve users in identifying keyword search intentions, user interaction can take time when the result set size is large. To solve this problem, we will develop various keyword suggestions for users in the context of specific keywords in the search terms. Thus, users can select specific search terms or modify existing search queries based on various search suggestions.

Consider the questionnaire $q = \{\text{query database}\}$ in the DBLP dataset. There are 21,260 publications or locations with the words "database" and publications or localizations 9,896 entries with the keyword "search" which contains 2,040 keywords, two words. Together, if we read directly the search results for the keywords, the process will be long and difficult to use because it produces a lot of effect. It takes about 54.22 seconds to calculate all SLCA results of q with XRank [2]. While system processing time is acceptable by accelerating keyword searches with efficient algorithms [3], [4], accidental search intentions and redundant results of large data retrieval will cause user disappointment. To solve the problem, we have a search term definition different from the original search term in relation to the XML data context, which can be used to search for different search objectives for the original search term.

II. RELATED WORK

In [17], Sarkas et al. provides a solution for creating interesting and meaningful k extensions for search terms by splitting additional words with interesting k values use extended queries to find more specific documents. The interesting thing is that the claws with the idea of surprise [19], [20], [21] in [18], Bansal and the faculty offer an effective algorithm to identify groups of keywords in the collection a large size of blog posts for temporary periods. Our work combines the ideas of both: we will measure the relationship of each pair of words using our common data model in Equation (1), a surprisingly easy metric. Next, we create a graph that has a long-term relationship that will retain all the terms and relationship values of these different words from [17], [18], very accurate for XML data.

Many researchers focus on the different dimensions of distribution problems. In this survey, we committed to provide a thorough review of a wide range of distribution techniques, including definitions for various consistent algorithms and distribution techniques for certain applications, stream application as well as the distribution system.

The first study of the investor relations community [6], [7], [8], [9], [10] is the first problem to be solved. Most of the operations are distributed as post-processing or the new document retrieval filing process is based on the analysis of the result set and / or the Agrawal query log and faculty [7]. The intent of the user model is at a particular level of taxonomy and Radlinski and Dumais [11]. Accept the search intent for possible searches by saving searches. However, it is not easy to obtain taxonomy and record this useful research.

Jianxin Li and faculty [1] offer a way to provide users with various search terms, in the context of keywords specified in the information to be searched. In doing so, the user can select the desired search term or modify the original search term according to the various search suggestions returned [2].

Jiaheng Lu "Comparative Analysis of TreePattern XML Query Processing" proposes here the problem of matching XML tree models and exploring some of the most recent works and algorithms. This comprehensive comparison compares five holistic algorithms and demonstrates performance and scalability. There is no clear winner in every situation of our experience, but TreeMatch has a good overall performance in terms of working time and ability to handle general tree models.

Jiaheng Lu, Tok Wang Ling "Extended XML tree model mapping: theory and algorithm" by introducing the concept of correspondence to solve the problems of optimizing algorithms match the XML tree model. In this case, specify the best large query class for three types of queries: $Q =; ==; _;$ $Q =; ==; _;$ $<and Q =; ==; _;$ $<;$, respectively, and proposes a new holistic algorithm called TreeMatch to get a class ask the most appropriate questions according to the theory. Finally, many experiments show the advantages of our algorithm and the validation of the theoretical results [3].

Mirella M. Moro, Zografoula Vagena, Vassilis J. Tsotras, "Tree Query on a Lightweight XML Processor" proposes a classification of algorithms for processing tree queries taking into account important characteristics such as data access and the process of Coupling queries also. Adjust the DFA method and improve efficiency by accessing nodes from the B + tree instead of purely sequential scans. Such improvements yield better results than ordinary DFAs and also suggest general methods for checking all deep plans on the left. [4]

Kamala Challa, E.Jhansi Rani "Algorithm for Matching XML Tree Templates and Query Processing" in this issue, provides problems mapping XML tree models and exploring the work and algorithms most recently, TreeMatch And Tjfasthave recommends two algorithms. TreeMatch has a good overall performance in terms of working time and ability to process general trees. [5]

M.Muthukumar "The Efficiency of the TreeMatch Algorithm in XML Tree Model Mapping" provides various analyzes to identify the effectiveness of the XML tree model matching algorithm. TreeMatch has good overall performance in terms of labeling, optimization, query handling, list of results, and ability to process extended XML tree models (branches). In this TreeMatch to get the best level of search, the TreeMatch template matching algorithm can respond to complex queries with good performance.

JT Yao M. Zhang "The Fast Tree Model Matching Algorithm for XML Query" provides a TreeMatch algorithm to directly find all query tree model mappings in the XML data source. Unlike previous searches on matching the query tree model, the TreeMatch algorithm does not need to break the tree into a linear model and does not generate any intermediate results that are not part of it the end result. The TreeMatch algorithm works when a non-leaf model node does not occur with autofill. Self-containment is rarely found in actual XML documents and can be easily identified.

III. PROPOSED SYSTEM

The With the q query and XML T data, our goal is to search for candidates with the above-mentioned query k in terms of high relevance and maximum diversity. For q in T here, each applicant presents a context or an opinion. Determined to find q in the T, examine the T-XML data and the relevant word-based dictionaries. The method of organizing W elements depends on the context of the application and will not affect our conversations later. For example, it may be a subset of all words made up of T-text or a set of appropriate vocabulary related to certain applications. In this work, different vocabulary pairs will be selected based on their mutual information. This combined information [15], [16] is used as a basis for selecting properties and changing the learning properties of the machine can be used to specify both the relevance and redundancy characteristics of variables, such as the selection of minimum repeatability properties. Suppose we have an XML T-tree and an R (T) output example.

We begin formal education on the problem of diversification in the XML keyword search, which allows to directly calculating a variety of results without having to extract all the relevant candidates. For this purpose, when receiving keyword queries, we will obtain a keyword associated with each search term based on the XML data collected in the probability theory used as a criterion for selecting attributes. Our specification selection is not limited to the XML element label. [23]

The combination of original words, properties, and search terms may indicate a variety of contexts. (Indicating that this is a specific research intention), we then evaluate each research intent received by measuring relevance to the original search term and the novelty of the generated results. In order to effectively compute different keyword searches, we propose a basic algorithm and two improved algorithms based on the observed properties of the different keyword search results. [23].

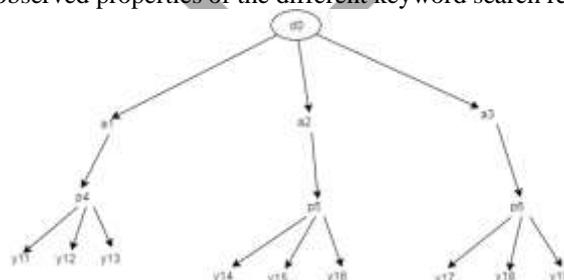


Fig.-1 Architecture of XML context diversified search

Feature Selection Model

Considering T-XML data and associated word dictionaries, the W method for organizing W elements depends on the context of the application. [1] This may be all or part of a vocabulary containing T-text or a set of words associated with certain applications. Different pairs of terms are selected on the basis of mutual information. Information sharing is used as a criterion for selecting features and modifying automatic learning properties. It can be used to specify both the relevance and redundancy characteristics of variables such as the selection of minimum redundancy properties [1].

Suppose we have an XML tree T and an example result R (T). Prob (x, T) is the probability that the word x appears In R (T) such that $Prob(x, T) = |R(x, T)| / |R(T)|$ where $|R(x, T)|$ is the number of results with x for Prob (x, y, T) is the probability that the terminologies x and y appear in R (T), such that $Prob(x, y, T) = |R(x, y, T)| / |R(T)|$ and there is free. Knowing that x will not give any information on y and vice versa. As a result, their shared data is zero. On the other hand, if the words x and y are identical, knowing that x determines the value of y and conversely. Therefore, simple measures can be used to measure the number of combined words to be observed, increase vocabulary dependence, functionality while reducing vocabulary duplication. In this case we use the common data model which is widely accepted as follows:

$$Michigan(x, y, T) = Prob(x, y, T) * Connect [Prob(x, y, T) / (Prob(x, T) * Prob(y, T))]$$

For each word of XML data, we need to find a set of vocabulary characteristics that can choose words in any way, such as the top word of a property or a word with a common value greater than a certain threshold Domain Application or Data Administrator. The specifications of the function can be calculated and stored in advance before the questionnaire evaluation process. So when you use a search term, you can get a property matrix for a search term using a pair of W words. The matrix represents the gap in research intent. (As the candidate for the query) XML data from the original query w.r.t. therefore, our first problem is to select a subgroup of queries with the highest probability of context interpretation of the original query. In this case, the selection of the candidate will depend on the approximate sample at the entity level of the XML data. [2]

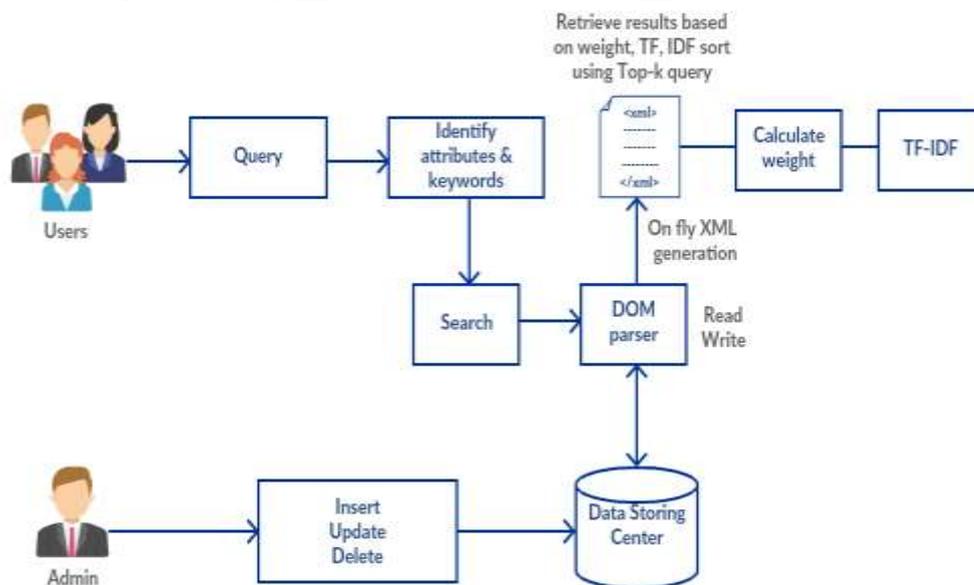


Figure 2 System Architecture

Top-k query

The retrieval system uses various methods to rank queries. Users are more concerned about the most important thing that is responding to popular queries in the vast answering area. Different emerging applications require powerful support for the top queries. For example, in the context of web performance and meta search engine performance, which are ranked by different search engines, are related to how effective rankings are. Similar applications exist in the view of data retrieval and data mining. Most of these applications compute queries related to joins and integration of multiple inputs to get maximum results.

The aim of popular search is to retrieve the best answer from a large collection of notes. Using the Top-k algorithm, we find the most accurate records from the record sets that match the keywords that are filtered and sorted by their scores. An XML tree is created with each result set, which is called a tree. The Steiner Tree Steiner created is a list of all the records sorted by their score, starting with the nearest result.

The most popular search algorithm uses score documents with keywords. Here, use this algorithm to score "Tuple Units." A unit is a set of highly relevant sets that contain search terms. Using these tuples, we create a set of tuples. Up when two tuples are directly related to each other.

The Top-K questionnaire score was based on two scoring methods: direct scoring and indirect scoring. TF-IDF stands for "Frequency Frequency, Inverse Document Frequency". The TF-IDE provides a way to rate the importance of words (or "terms") in a tuple, depending on the tuple. They repeatedly appear in the document several times.

Terms are called according to the following criteria:

- 1) If the word appears frequently in the tuples, this word is called significant and has a high score.
- 2) If a word appears in several sentences it is called a unique identifier and has a low rating.

So common terms like "the" and "for", which appear in many tuples, will be scaled down. Frequent appearances in one set are resized.

Top-K query processing algorithms work with age, weight, and high frequency word frequency and ignore low frequency words. This technique is called TF-IDF (Frequency Frequency - Inverse document frequency). TF -IDF consists of two words: First, the frequency (TF) is calculated, the number of times the word appears in the totem, divided by the total number of words in that word. The second word is Inverse Document Frequency (IDF). Logarithm of the number of tuples In the corpus divided by the number of tuples where specific words appear.

$$IDF(t) = \log(\text{Total number of tuple-rows} / \text{Number of tuple-rows with term } t \text{ in it}).$$

Consider examples in Figure 2 of XML describing vehicle data used in various forms: A car in the form of a car B, car in the place of sale, and C consists of cars arranged according to model and year. Take a look at Q1, a simple question. Ask for information (Usedcar) when considering certain conditions (do 'Ford' and prices do not reach "9000"). To create an XQuery to display this simple query, the user faces a challenge. That element is the main component of the used car, but the price may be either children or grandchildren. In addition, the custom definition of the XML data source can be very problematic for using XQuery.

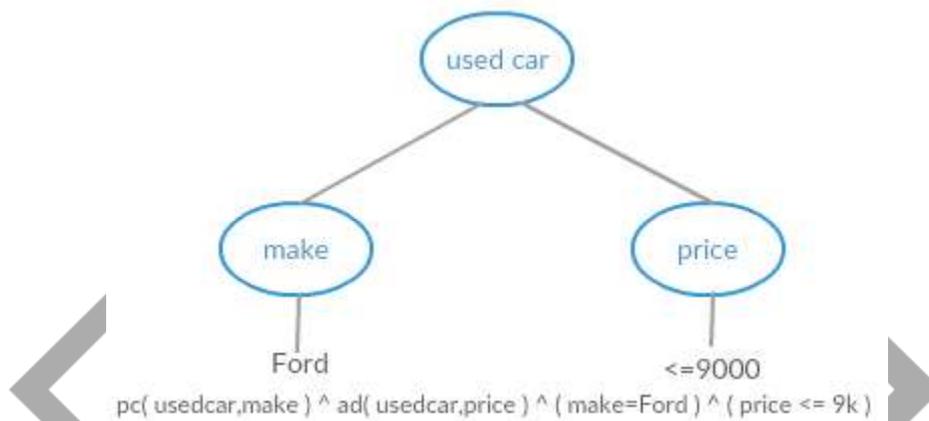


Figure 3 Query Generation

Query relaxation allows the system to narrow down query queries to meet user needs. Traditionally, user-submitted search queries are resolved in different ways and in different ways. Obviously, the estimated search can change the format from one query to another, and the change among them is based on two views: relaxation, structure and content relaxation. Let Q be a query that states if Q 'has a Q', which is due to a decrease in search constraints using query relaxation or similar content substitution, we will say that Q and Q 'are approximate queries and matches.

IV. SIMILARITY RELATIONS AND ASSESSMENT ON CONTENTS

Before In order to distinguish the different structural relaxation responses from the original questionnaire, we propose a similar relationship assessment approach to find the ranking factors (structures) used to find answers to the Find Top Answers.

For example, consider XML data, assuming that users submit simple queries to retrieve nested used car data within the car. The central node in the path from cars to used cars usually are not of major concern the middle of the user, like the mid-nodes version and the year.

The SimTree coefficient tree is a ranking factor used in ordering solutions to find the top-k solution on a structure to optimize online processing. Calculating the coefficient, the similarity between tree-like queries can be achieved in during offline processing.

Algorithm for Value Partitioning

Define Input: Attributes, Ai numbers, All values of Ai, and

Partition number n

Output: set of maximum limits of each partition s

a) min = getMinValue (val) // Get the smallest value

b) max = getMaxValue (val) // Get Maximum

c) total = getSize (val) // get 'number

d) avg = total / n

e) low = min up = max // set lower limit and upper limit

f) while (low < up)

g) c = query (low < Ai < up) // Run the query low < Authority and return

Number of results

- h) If ($c \leq \text{avg}$)
- i) add (up, s) // increase in set s
- j) low = up, up = max
- k) else up = low (up \leq low) / 2
- l) while
- m) return s

The obvious measure for distance on tree nodes could be a path metric [6], i.e. length of the shortest path between them, the similarity measure that is based on path metric then could be expressed as

$$s(v_i, v_j) = \frac{1}{1 + l(v_i, v_j)} = \frac{1}{1 + l(v_i, lca_{ij}) + l(v_j, lca_{ij})}$$

Procedure

Build Tree (T)

Input: XML Document T.

Output: Count-Stable synopsis S of T.

begin

1. $H := \emptyset; S := \emptyset$
 2. for each element $e \in T$ in post-order do
 3. $C := \{(ui, ci) : ui \text{ is a node in } S \text{ and } |\text{children}(e) \cap \text{extent}(ui)| = ci > 0\}$
 4. if ($H[\text{label}(e), C] = \emptyset$) then
 5. Add node u to S with $\text{label}(u) = \text{label}(e)$
 6. $H[\text{label}(e), C] := u$
 7. for $(ui, ci) \in C$ do add edge $u \xrightarrow{ci} ui$ to S
 8. endif
 9. $u := H[\text{label}(e), C]; \text{extent}(u) := \text{extent}(u) \cup \{e\}$
 10. end for
- end

In this model, not only But given the likelihood that newly created search terms, such as relevance, are also taken into account, we also take into account the new results and their differences, which is new.

1) The created query has the highest probability of interpreting the context of the original query queue, taking into account the data to search for.

2) The created query has the biggest difference from the previously created query set. Q . So we have a total notation function.

Points $(q_{\text{new}}) = \text{Prob}(q_{\text{new}} | q, T) * \text{DIF}(q_{\text{new}}, Q, T)$ where $\text{Prob}(q_{\text{new}} | q, T)$ represents the probability that q_{new} is the intention to seek when issued the original request q on the T data; $\text{DIF}(q_{\text{new}}, Q, T)$ represents the percentage of results generated by q_{new} , but not by queries previously created in Q

To evaluate the relevance of the probabilities of proposed research suggestions, the initial query that we can use

1. The Bay theorem
2. Properties of the SLCA semantics

V. EXPERIMENTAL SETUP

After Measuring a free distance on a tree node can be a trajectory indicator [6], such as the length of the shortest path between them, measuring similarities as a function of the trajectory indicators

$$s(v_i, v_j) = \frac{1}{1 + l(v_i, v_j)} = \frac{1}{1 + l(v_i, lca_{ij}) + l(v_j, lca_{ij})}$$

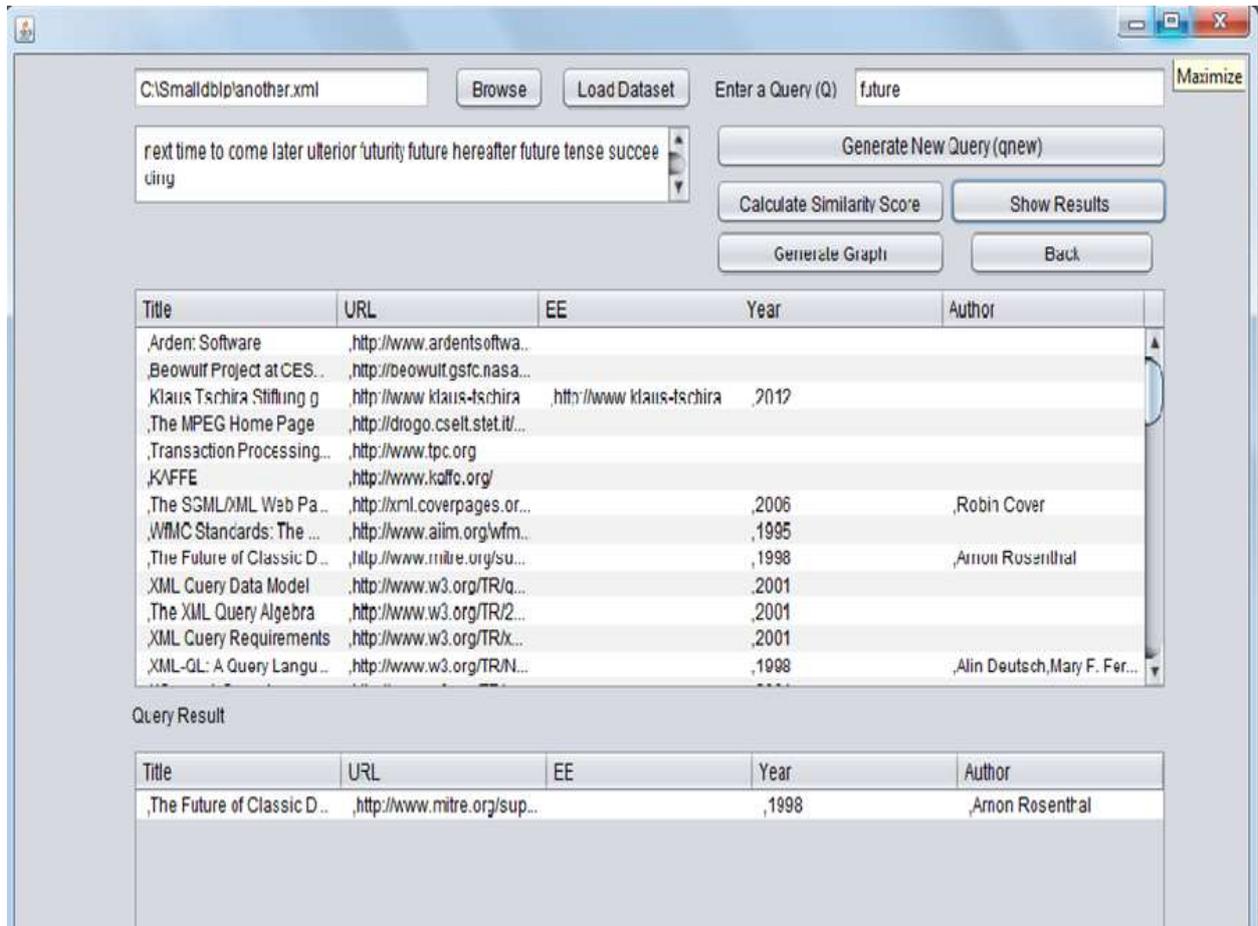


Figure 4 Data Loading and Parsing

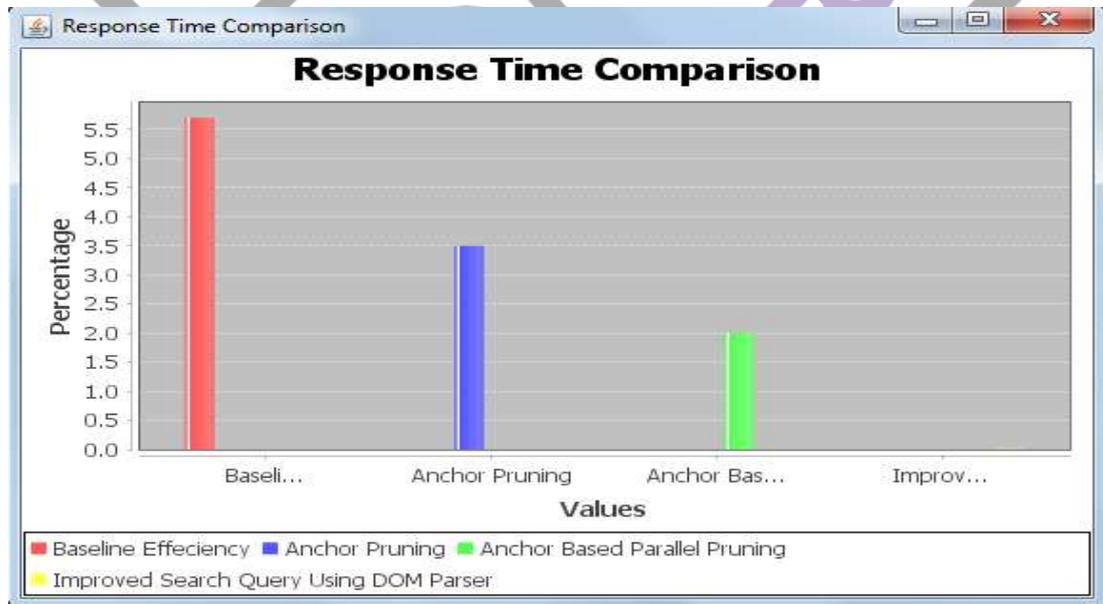


Figure 5 Result Analysis

Parameters	(BE)	(AE)	(ASPE)	Our Method
Response Time	5.7 sec	3.5 sec	2 sec	7.4 millisc

Table 1 Result Comparison with existing system

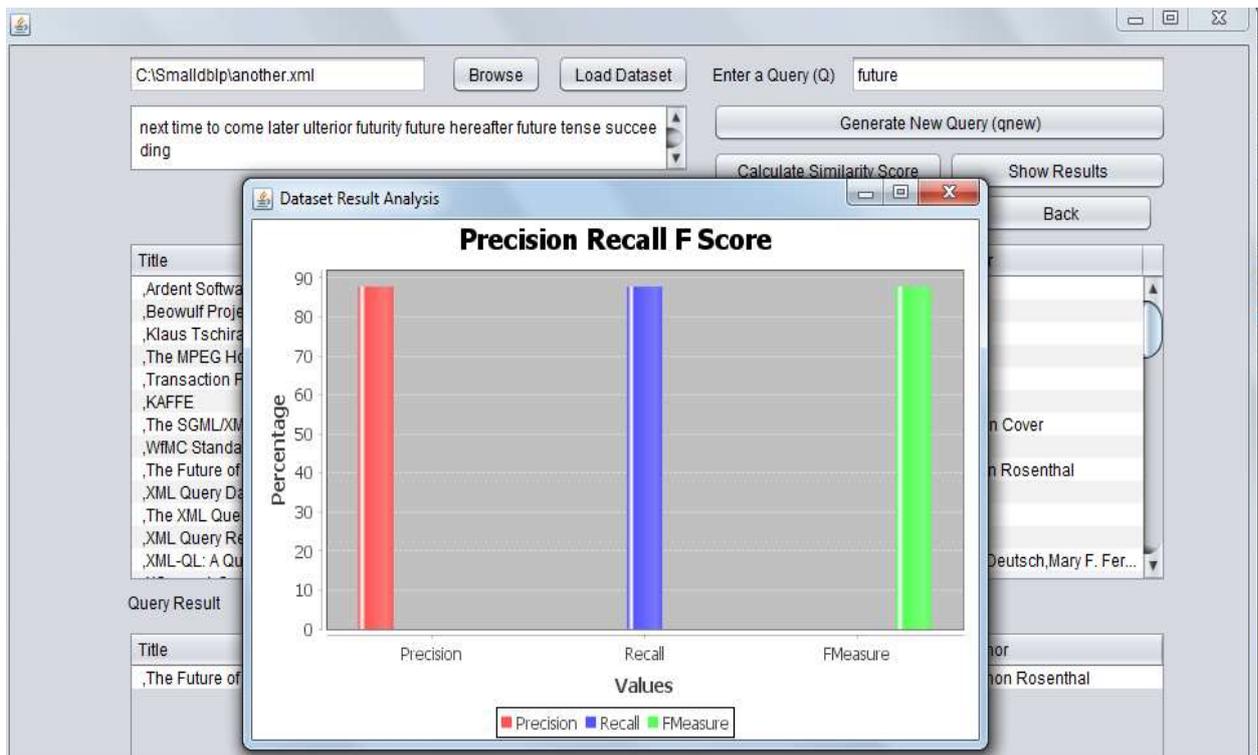


Figure 6 Key Index Parameters

VI. CONCLUSION

In this article, we focus on how to find more than one large XML dataset and provide a variety of output formats for contextual keyword queries. This work provides a way to search various search terms analyzed from data. The contextual XML of the query keywords in the data, the diversity of the context is measured by exploring the relevance to the original search terms and the novelty of the results. In addition, the structure is a powerful algorithm based on the observed properties of the XML query analysis. Our comparative study demonstrates the effectiveness of the proposed algorithm by executing queries. Much about the XMark dataset at the same time, we also monitor the effectiveness of our diversification model by analyzing the intent of the search returned for the specified query via the DBLP dataset, as well as the search intent and the results per user.

REFERENCES

List and number all bibliographical references in 10-point Times, single-spaced, at the end of your paper. When referenced in the text, enclose the citation number in square brackets, for example: [1]. Where appropriate, include the name(s) of editors of referenced books. The template will number citations consecutively within brackets [1]. The sentence punctuation follows the bracket [2]. Refer simply to the reference number, as in “[3]”—do not use “Ref. [3]” or “reference [3]”. Do not use reference citations as nouns of a sentence (e.g., not: “as the writer explains in [1]”).

Unless there are six authors or more give all authors’ names and do not use “et al.”. Papers that have not been published, even if they have been submitted for publication, should be cited as “unpublished” [4]. Papers that have been accepted for publication should be cited as “in press” [5]. Capitalize only the first word in a paper title, except for proper nouns and element symbols.

For papers published in translation journals, please give the English citation first, followed by the original foreign-language citation [6].

- [1] Jianxin Li, Chengfei Liu, Jeffrey Xu Yu, “ContextBased Diversification for Keyword Queries Over XML Data”, IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 27, NO. 3, MARCH 2015.
- [2] SimonFarande, Dr.D.S.Bhosale, “Survey Paper On Xml Context Diversified Search”, International Conference on recent innovation in engineering and management, 23 march 2016.
- [3] Jiaheng Lu, Tok Wang Ling, Senior Member, IEEE, Zhifeng Bao, and Chen Wang. “Extended XML TreePattern Matching: Theories and Algorithms”. IEEE transactions on knowledge and data engineering, vol. 23, no. 3, march 2011.
- [4] M. Moro, Z. Vagena, and V.J. Tsotras, “Tree-Pattern Queries on a Lightweight XML Processor,” Proc. Int’l Conf. Very Large DataBases (VLDB), pp. 205216, 2005.
- [5] Kamala Challa, E.Jhansi Rani “Algorithms for XML Tree Pattern Matching and Query Processing” Int.J. Computer Technology & Applications, Vol 3 (1), 447-451 JAN-FEB 2012.
- [6] J. Goldstein, “The use of MMR diversity-based reranking for reordering documents and producing summaries”, Proc. SIGIR, pp. 335-336, 1998.

- [7] R. Agrawal, S. Gollapudi, A. Halverson, S. Jeong, "Diversifying search results", Proc. 2nd ACM Int. Conf. Web Search Data Mining, pp. 5-14, 2009.
- [8] H. Chen, D. R. Karger, "Less is more: Probabilistic models for retrieving fewer relevant documents", Proc. SIGIR, pp. 429-436, 2006.
- [9] C. L. A. Clarke, M. Kolla, G. V. Cormack, O. Vechtomova, A. Ashkan, I. MacKinnon, "Novelty and diversity in information retrieval evaluation", Proc. SIGIR, pp. 659-666, 2008.
- [10] A. Angel, N. Koudas, "Efficient diversity-aware search", Proc. SIGMOD Conf., pp. 781-792, 2011.
- [11] S. T. Dumais, "Improving personalized web search using result diversification", Proc. SIGIR, pp. 691-692, 2006.
- [12] Z. Liu, P. Sun, Y. Chen, "Structured search result differentiation", J. Proc. VLDB Endowment, vol. 2, no. 1, pp. 313-324, 2009.
- [13] E. Demidova, P. Fankhauser, X. Zhou, W. Nejdl, "DivQ: Diversification for keyword search over structured databases", Proc. SIGIR.
- [14] J. Li, C. Liu, R. Zhou, B. Ning, "Processing xml keyword search by constructing effective structured queries", Advances in Data and Web Management, pp. 88-99, 2009, Springer.
- [15] H. Peng, F. Long, C. H. Q. Ding, "Feature selection based on mutual information: Criteria of maxdependency max-relevance and min redundancy", IEEE Trans. Pattern Anal. Mach. Intell., vol. 27, no. 8, pp. 1226-1238, Aug. 2005.
- [16] C. O. Sakar, O. Kursun, "A hybrid method for feature selection based on mutual information and canonical correlation analysis", Proc. 20th Int. Conf. Pattern Recognit.
- [17] N. Sarkas, N. Bansal, G. Das, N. Koudas, "Measuredriven keyword-query expansion", J. Proc. VLDB Endowment, vol. 2, no. 1, pp. 121-132, 2009.
- [18] N. Bansal, F. Chiang, N. Koudas, F. W. Tompa, "Seeking stable clusters in the blogosphere", Proc. 33rd Int. Conf. Very Large Data Bases.
- [19] S. Brin, R. Motwani, C. Silverstein, "Beyond market baskets: Generalizing association rules to correlations", Proc. SIGMOD Conf., pp. 265-276, 1997 [20] W. DuMouchel, D. Pregibon, "Empirical bayes screening for multi-item associations", Proc. 7th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, pp. 67-76, 2001.
- [21] A. Silberschatz, A. Tuzhilin, "On subjective measures of interestingness in knowledge discovery", Proc. 7th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, pp. 275-281, 1995.
- [22] K. Järvelin, J. Kekäläinen, "Cumulated gain-based evaluation of ir techniques", ACM Trans. Inf. Syst., vol. 20, no. 4, pp. 422-446, 2002.
- [23] Sivaji Yerraguntla, Borugadda.NagaRaju, "ContextBased Diversification For Keyword Queries Over Xml Data", International Journal of Applied Sciences, Engineering and Management, Vol. 04, No. 01, January 2015, pp. 57 – 61.
- [23] B. Noble, and I. N. Sneddon, "On certain integrals of Lipschitz-Hankel type involving products of Bessel functions," Phil. Trans. Roy. Soc. London, vol. A247, pp. 529–551, April 1955