

# Topic-based Modeling using Query Based Approach with Independent XML Structure Data

<sup>[1]</sup> Heena Malani, <sup>[2]</sup> S. R. Ghungrad

<sup>[1][2]</sup> Department of Computer Engineering, MSS College of Engineering Jalna, Aurangabad, <sup>[2]</sup> Professor

**Abstract:** The rapid adoption of XML as a standard for representation and exchange of information is an enormous amount of XML data retention and archiving on the Internet or incorporate data repositories. This will lead to the development of online decision support systems where users and analysts can interact large XML datasets through open query interfaces (such as XQuery or XSLT). Estimated responses are effective mechanisms to reduce response time and provide feedback to users. This approach has been successfully used in relational and more confident systems in the XML world. Complex evaluations of structured data are more expensive. Ranking and the most relevant search results became the most popular paradigm for processing XML messages. However, existing proposals did not adequately consider the structure, and therefore were not rational. Link

structure to content to answer a relaxing question. To solve this problem, we offer a sophisticated query framework to support an approximate query of XML data. Responses under this framework do not need to be strictly adhered to. It may use a view that was subtracted from the original query. So we've developed a new top-of-the-line search method that can generate the most likely response in a ranking order relative to rank. We work with a comprehensive set of experiments to demonstrate the effectiveness of our approach in terms of precision and recall.

**Keywords:** - XML, approximate queries, query relaxations, top-k.

## I. INTRODUCTION

In the last few years, XML, which originally offered to represent, exchange, and publish information on the web, has become widely used in many applications: as a solution for disseminating inherited data in order to retain information that can be Show with no other original information. To ensure interoperability between applications to integrate web services to retrieve web and other data, these applications create strong demand for both repositories, XML document repositories, and XML query languages. It is designed for data retrieval and restructuring of XML data. But all languages that actually apply for XML data [BC00] only provide answers to search queries. When used with large XML or warehouse repositories, accurate queries may take time to respond. In order to overcome this problem in traditional relational warehouses, we will support an approximate query based on concise statistical data generated using a histogram or sampling technique. We believe that the current trend of XML claims for extending the approach is also a very big question of the XML data set. The basic idea for the estimated answer is to store the pre-calculated summary of the XML warehouse, also known as Comprehensive data collection and retrieval of those data represent the original database, saving time and IT costs. The basic idea for an estimated answer is to store a set of predefined XML bundles, also known as bundles. (To gather concise information), and to find out these data instead of the original database, saving time and money on the computer. Users must submit a request to the system, which should use

the information instead of the original information to respond to the query. The result should come in a very short time at the expense of precision loss. The ideal way to answer XML queries should be to benefit from database query style and IR style queries, since the IR query improves query value by allowing for higher query volumes of the query. Text content in IR style search by specifying the context to be searched[6].

### The contributions of this document are as follows:

1) We offer a way to relax queries that incorporate structure and content, as well as factors that users are more concerned about, to support their estimated XML queries, especially how our users assume that they are concerned about Analyze the original query of the user to facilitate the relaxation of the query. In addition, our approach is different from the relaxation sequence, rather than equally important for each node to relax. In particular, the first relaxed structure to be considered is what has the highest similarity coefficients with the original query and the first node to be released as the least important node.

2) We define a similar assessment by analyzing the inherent meaning presented in the XML source. To solve our problem, we divide the node into two groups: the attribute node and the numerical attribute node, and design the corresponding method. Similarly, in evaluating relationships, the similarity of the attribute nodes and the nodal attributes. We also designed a non-circular indexed (CDAG) chart to create and structure relaxation and develop effective evaluation coefficients for similar structural evaluations.

Based on the similarity assessment offered and the priority, we will relax the query with an automated recovery method that can generate the most effective response possible.

## II. RELATED WORK

Extensive searches were carried out on structured queries as well as text search on XML data and graph data [4]. Structural Join [4] decomposes the tree model queries into a set of binary components, and then the matches of each individual component are assembled to obtain the final results. Holistic Join [12] adopts a technique of a linked stack string to compactly represent partial results on the query paths, and these paths can be compounded to get the final matches for the tree query. By introducing a fuzzy marking scheme, work [15] introduced a holistic joining algorithm to match twigs involving OR / NOT predicates. Based on the fuzzy label flows, the problem of the ordered tree pattern Concordance to fuzzy XML data has been set in the following work [12]. A common characteristic of these above approaches is that they deal exclusively with tree model queries with precise query conditions. They do not process an approximate query when users can not specify their query conditions accurately. A system called XRANK, which is designed for keyword searches in XML documents, is used in [10]. XRANK has a ranking mechanism and returns documents as answers. Other search engines based on keywords for XML, called XSearch, are shown in [16]. In XSearch, query responses are categorized using extended retrieval techniques and built in the same order in the ranking. However, as mentioned above, structured query queries are rarely found in real-world applications, as these methods rely on schema information to avoid unintended consequences. Trust or may yield unsatisfactory results. Recent efforts have included the ability to manage database structured data with the ability to find effective keywords for querying. To provide structured XML query responses effectively, unstructured query mechanisms [11], keyword query mechanisms [15], and query relaxation mechanisms [14]. In [10], Brodianskiy et al. Propose a framework for deriving self-correcting queries on XML. Within their framework, satisfactory similar requests are generated when the given request is unsatisfactory. The user will then choose a satisfiable questioning of interest, and receiving exactly satisfactory answers to this request. In [3], Mandreoli et al. Present a model that combines structures and vocabularies to support approximate queries. In [13], Liu et al. Provide a content-oriented approach to XML content and structure requests. They first decompose a query of content and structure into multiple query fragments, and get the

intermediate results, and then combine and rank the results according to the relevance of the original query.

In [11], Buche et al. Developed techniques for fuzzy marking and XML queries in [13]. Based on the WordNet dictionary and fuzzy set theory [1, 2, 9], Campi et al. The Keif query constraint in XML [18], Damiani et al., Introduces a flexible query model through contextual crawl options to support an approximate XML query [17] based on semantics. Of the Cohen and Shil horizon. each presents a query language for XML that includes both value and structure requirements.

In Termehchy and Winslett, they provide a way of ranking for XML keyword queries, which categorize candidate responses according to the statistical criteria of their correspondence. Their approach can enable users to take advantage of XQuery to process tree queries directly, without the need for knowledge of the checkbook.

In [18], Wang and his team propose ways to add value to a given keyword by introducing vague or ambiguous keywords or keywords that have a semantic relationship in the vocabulary database to improve their performance. Inquiries only. However, using the terminology database alone can not find similar substitutes that may be useful to users, so their approach does not solve the problem of finding the most relevant answers.

In [15], Yan et al. Indicate a preference ranking pattern to handle the approximate search queries in XML. Their approach is based on user-defined attention scores to generate a top response. Add a load to the user's search query. Unable to discover possible responses that users may be interested in because there is no similarity assessment designed to derive the inherent meaning presented in the XML source.

In [21], based on flexible semi-flexible matching and matching, kanzet al offers two definitions adapted to the ontological critique of semi-structured data. Structuring and querying text in response to an estimated query has generated considerable interest in [6], Amer-Yahia et al. Launch a framework called FlexPath, which combines structure and text.

Compared with the previous proposal, our approach is different from the above. To solve our problem, we have developed a method of relaxing search terms, including structure and content. It also includes factors that the user is concerned about. (Inferred from the original query analysis by identifying the relaxation sequence and structure), answering the XML query. In addition, the previous work

often emphasized the importance of each node to be relaxed. However, our method of determining the first relaxation structure is the structure with the highest tree similarity coefficient to the original query, and the first node to be the least significant node. The first innovation of our solution is to introduce factors that the user is more concerned with and to identify the sequence of rest. Another innovation is the similarity assessment designed to find potential responses that users may be interested in that have been refined by analytics. The inherent meaning presented in the XML source. In our solution, we group the nodes in two groups: the attribute node and the attribute node. We then designed a way to evaluate the similarity of cluster-based attributes by estimating the percentage of ANV pairs associated with the tree's meaning contacts, and developing a similar approach for evaluating the similarity. The node of the numerical attribute is Evaluation function

### III. DATA AND QUERY MODEL

**A. Data Models** We consider the data model for the XML data that is displayed as data sets. Basically, a data structure is part of the real world via an entity. (Usually a set of attributes) values and relationships. A simple XML snippet with different car collections is shown in Figure 1.

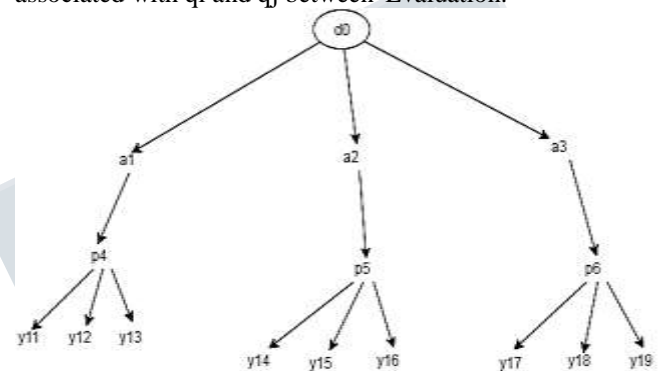
#### XML Data Model

XML data format We follow the XML document in the form of a large  $T(V, E)$  tree, labeled by each node. The node  $u \in V$  corresponds to the XML element and looks like a unique object identifier (oid). ) And the label (or tag) that comes from the literal strings of the literal, which will compile the definition of the edge element (ei, ej). E is used to capture the confinement of (sub) element ej under ei in the base. Information (we use the label (ei), child (ei) to Defines labels and subset nodes for element nodes  $ei \in V$ ). For example, FIG. 1 shows an example of an XML data structure containing bibliographic data. The document contains elements of the author, each with several titles and sub-elements of paper and books. Each article has one year's published title and at least one search term, while a book recently gave the title. Note that the element node in the structure is named with the first character of the element label plus the unique identifier. Elements of  $Le$  in  $T$  are generally valued. But our main goal in this task is to capture and look up the tag structure of the XML data structure, rather than spreading the relevant values.

The XML query templates we focus on are XML branch queries, which represent the basic blocks of declarative language declarations for XML (including XQuery [4] and XSLT [7] standards). In short, branch search describes the

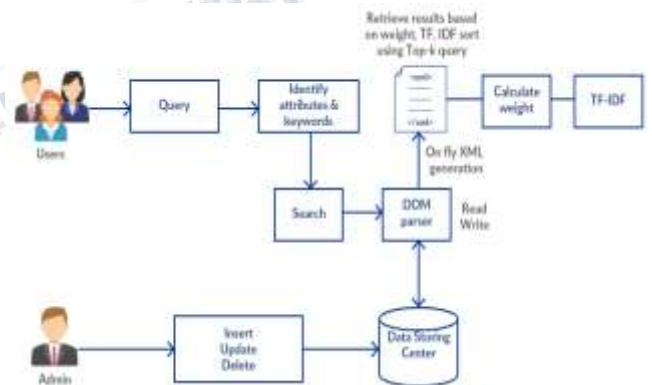
path. The complexity in the XML data structure and the return of structured XML results are structures that are generated through the interpreted values. (That is, the integration of the structure) of the multiple path representation (shown in XPath [8]).

We create a query Q as a Tree. The query is labeled by the TQ node, where (1) each node of TQ has a qi label in Q (where q0 is the different root node to the root of the XML document) and (qi, qj) of TQ has a description of the XPath expression path (qi, qj) describing the constraints of the specific structure specified in Q between the data elements associated with qi and qj between 'Evaluation.'



**Figure 1 XML Document**

#### Proposed System



**Figure2 System Architecture**

There are two users using the system

- 1) Data user
- 2) Master Admin

The user's XML query uses knitting and wording techniques to identify keywords and attributes. The DOM analyzer will search for attributes and keywords from the XML file. To answer the XML query, calculate how often the frequency (TF) and inverse document frequency (xx) of the xml file are



extracted to the user. Sort the weighted frequencies (TFs) and inverse frequencies (IDFs) according to the most popular search algorithms. The retrieval system uses a variety of ways to rank queries. Users are more concerned about the most important answer, such as the answer to the top-k question in the answering area is very large. Different Apps Need Effective Support for Top Queries. Administrators write to insert, update, and delete data from the repository. This information is used by the DOM parser to retrieve the output.

### 3.1 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 recordsets 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.

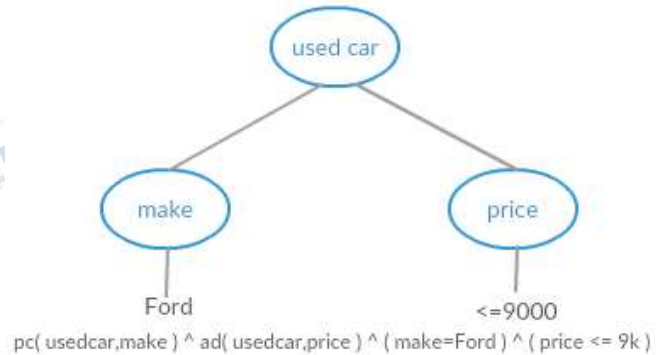
So common terms like "the" and "for", which appear in many tuples, will be scaled down. Frequent appearances in one set are resized.

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.

### 3.2 Query Generation



**Figure 3 XML Data**

### 3.3 Perform Query Relaxation

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.

### 3.4 Similarity Relations and Assessment on Contents

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

Input: Attributes,  $A_i$  numbers, All values of  $A_i$ , and Partition number  $n$

Output: set of maximum limits of each partition  $s$

- a)  $\min = \text{getMinValue}(\text{val})$  // Get the smallest value
- b)  $\max = \text{getMaxValue}(\text{val})$  // Get Maximum
- c)  $\text{total} = \text{getSize}(\text{val})$  // get 'number
- d)  $\text{avg} = \text{total} / n$
- e)  $\text{low} = \min$   $\text{up} = \max$  // set lower limit and upper limit
- f) while ( $\text{low} \leq \text{up}$ )
- g)  $c = \text{query}(\text{low} \leq A_i \leq \text{up})$  // Run the query low Authority and return Number of results
- h) If ( $c \leq \text{avg}$ )
- i) add ( $\text{up}, s$ ) // increase in set  $s$
- j)  $\text{low} = \text{up}$ ,  $\text{up} = \max$
- k) else  $\text{up} = \text{low} + (\text{up} - \text{low}) / 2$
- 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

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

#### Procedure

Build Tree (T)

**Input:** XML Document T.

**Output:** Count-Stable synopsis S of T.  
 begin

1.  $H := \phi$ ;  $S := \phi$
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 ( $H[\text{label}(e), C] = \phi$ ) 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  $u \text{ --} \rightarrow ui$  to  $S$
8. endif
9.  $u := H[\text{label}(e), C]$ ;  $\text{extent}(u) := \text{extent}(u) \cup \{e\}$
10. end for
- end

### Experimental Setup

To provide a thorough assessment of our approach in terms of impact on top search query processing, we conducted our experiments on synthetic synthesized XML datasets and collected data on second hand cars. Crawl from autos.yahoo.com In order to provide relaxation, the query is created using a different DTD. For our evaluation, the generated document contains information that represents the node, the entity, the warranty, and the description as entity nodes, and Make a model year, price, position, length, and color as the attribute nodes, which in year, place, and color format. The attribute and price attribute and the mileage attribute are considered as nodes. We then tested differently used car data sets, expressed as AD1, AD2, and AD3 (the default data set is AD2) by compiling documents in terms of the number of nodes between 150,000 and 350,000. Used ED1 vehicle, which contains data about 2,000,000 nodes, will collect data from www.edmunds.com An entity node is an entity node. Options and descriptions are the connection nodes. The generation, year, price, distance, location, and color attribute are the nodes. The place, year, and color format is the attribute node, and the cost and distance are the numerical attributes nodes.

To showcase our improvised work, we have selected queries to display the query relaxation structure, which is possible in real situations. We used JAVA ranking and query methods, and tested on a 280GHz I5 processor and 8GB RAM running on Windows 7.

### Key index parameters for results

To retrieve data with a binary precision classification. (Called positive predictive value) is the portion of the instance that retrieves the relevant data while it is being restored. (Also known as sensitivity) is the portion of the associated instance that is restored. Accuracy and recall are based on understanding and measuring relevancy.

In simple terms, high accuracy means that the algorithm yields results that are associated with more irrelevant results, while higher recall means that the algorithm yields the most relevant results.

**Table 1.0 Precision, Recall, F measure Calculations**

Precision	Recall	F Measure
$P = TP/(TP + FP)$	$R = TP/(TP + FN)$	$2*P*R/(P+R)$

**Table 2.0 Confusion Matrix**

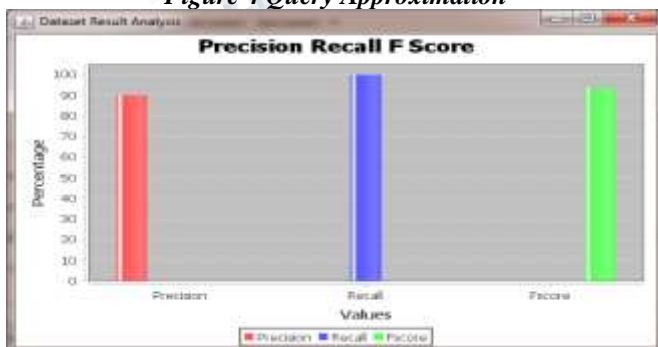
Confusion Matrix	Predicted True	Predicted False
Actual True	7	17
Actual False	19	2

**Table 3.0 Tabular Top-K Results**

Top K Count	Accuracy	Precision	Recall	F measure
5	83.00 %	82.00 %	83.00 %	89.00%
10	85.00%	85.12%	85.00%	88.00%
15	85.80%	91.90%	91.39%	90.36%



**Figure 4 Query Approximation**



**Figure 5 Graphical Results**

Based on the results above, we will summarize the following two points. First, the value of the similarity estimation presented in the XML data source should be set. Use d to evaluate the XML query, as it can provide effective support to find potential answers. Second, guessing about the factors that users are more concerned about is a key source for finding the most relevant answers when answering the XML queries.

#### IV. CONCLUSION

In this article, we present the concept and architecture of a system that collects data online on XML data. This system has the ability to provide quick feedback on the collection request by approaching and refining the final response throughout the process. Query processing It also provides a precise guarantee by attaching the weight data to the estimate. We've introduced a new query processing method that breaks out query patterns into search query patterns. Powerful search query processing is performed by the new service provider to select and join a pattern route along with the appropriate index structure. In order to obtain accurate estimates, we have modified the principle of DOM segmentation in XML query processing. During in-depth evaluation, we have shown that our systems refer to explicit assumptions in response to Finally, before the original system can be produced. In addition, good estimation can be obtained quickly for queries that do not have branch nodes. We have tried to improve our relaxed approach and ranking to become a friendly improvement guide in a dynamic environment. In future we also plan to improve our approach by incorporating emerging semantic technologies to address estimated queries with unstructured data structures and associated data.

#### REFERENCES

- [1] Ashraf Aboulnaga, Alaa R. Alameldeen, and Jeffrey F. Naughton. "Estimating the Selectivity of XML Path Expressions for Internet Scale Applications". In Proceedings of the 27th Intl. Conf. on Very Large Data Bases, 2001.
- [2] P. Buneman, M. Grohe, and C. Koch. "Path Queries on Compressed XML". In Proceedings of the 29th Intl. Conf. on Very Large Data Bases, 2003.
- [3] Kaushik Chakrabarti, Minos Garofalakis, Rajeev Rastogi, and Kyuseok Shim. "Approximate Query Processing Using Wavelets". In Proceedings of the 26th Intl. Conf. on Very Large Data Bases, 2000.

- [4] Don Chamberlin, James Clark, Daniela Florescu, Jonathan Robie, Jerome Simon, and Mugur Stefanescu. "XQuery 1.0: An XML Query Language". W3C Working Draft, 2001.
- [5] Zhimin Chen, H.V. Jagadish, Laks V.S. Laksmanan, and Stelios Pappas. "From Tree Patterns to Generalized Tree Patterns: On Efficient Evaluation of XQuery". In Proceedings of the 29th Intl. Conf. on Very Large Data Bases, 2003.
- [6] Zhiyuan Chen, H. V. Jagadish, Flip Korn, Nick Koudas, S. Muthukrishnan, Raymond Ng, and Divesh Srivastava. "Counting Twig Matches in a Tree". In Proceedings of the Seventeenth Intl. Conf. on Data Engineering, 2001.
- [7] James Clark. "XSL Transformations (XSLT), Version 1.0". W3C Recommendation, November 1999.
- [8] James Clark and Steve DeRose. "XML Path Language (XPath), Version 1.0". W3C Recommendation, November 1999.
- [9] Juliana Freire, Jayant R. Haritsa, Maya Ramanath, Prasan Roy, and Jerome Simon. "StatiX: Making XML Count". In Proceedings of the 2002 ACM SIGMOD Intl. Conf. on Management of Data, 2002.
- [10] Yannis E. Ioannidis and Viswanath Poosala. "Histogram-Based Approximation of Set-Valued Query Answers". In Proceedings of the 25th Intl. Conf. on Very Large Data Bases, Edinburgh, Scotland, September 1999.
- [11] Raghav Kaushik, Pradeep Shenoy, Phillip Bohannon, and Ehud Gudes. "Exploiting Local Similarity for Efficient Indexing of Paths in Graph Structured Data". In Proceedings of the Eighteenth Intl. Conf. on Data Engineering, 2002.
- [12] L. Lim, M. Wang, S. Padmanabhan, J.S. Vitter, and R. Parr. XPathLearner: An On-Line Self-Tuning Markov Histogram for XML Path Selectivity Estimation. In Proceedings of the 28th Intl. Conf. on Very Large Data Bases, 2002.
- [13] Jason McHugh and Jennifer Widom. "Query Optimization for XML". In Proceedings of the 25th Intl. Conf. on Very Large Data Bases, 1999.
- [14] Tova Milo and Dan Suciu. "Index structures for Path Expressions". In Proceedings of the Seventh Intl. Conf. on Database Theory (ICDT'99), Jerusalem, Israel, January 1999.
- [15] N. Polyzotis and M. Garofalakis. "Statistical Synopses for Graph Structured XML Databases". In Proceedings of the 2002 ACM SIGMOD Intl. Conf. on Management of Data, 2002.
- [16] N. Polyzotis and M. Garofalakis. "Structure and Value Synopses for XML Data Graphs". In Proceedings of the 28th Intl. Conf. on Very Large Data Bases, 2002.
- [17] Neoklis Polyzotis, Minos Garofalakis, and Yannis Ioannidis. Approximate XML Query Answers. 2004.
- [18] Neoklis Polyzotis, Minos Garofalakis, and Yannis Ioannidis. "Selectivity Estimation for XML Twigs". In Proceedings of the Twentieth Intl. Conf. on Data Engineering, 2004.
- [19] C. M. Procopiuc. Geometric Techniques for Clustering: Theory and Practice. PhD thesis, Duke Univ., 2001. [20] D. Sasha and K. Zhang. Fast algorithms for the unit cost editing distance between trees. Jnl. of Algorithms, 11, 1990.
- [20] Wei Wang, Haifeng Jiang, Hongjun Lu, and Jeffrey Xu Yu. Containment join size estimation: Models and methods. In Proceedings of the 2003 ACM SIGMOD Intl. Conf. on Management of Data, 2003.
- [21] Yuqing Wu, Jignesh M. Patel, and H.V. Jagadish. "Estimating Answer Sizes for XML Queries". In Proceedings of the 8th Intl. Conf. on Extending Database Technology, 2002.