

FLUX: Fuzzy Content and Structure Matching of XML Range Queries

Hua-Gang Li
Department of Computer Science
University of California
Santa Barbara, CA 93106.
huagang@cs.ucsb.edu

Divyakant Agrawal
Department of Computer Science
University of California
Santa Barbara, CA 93106.
agrawal@cs.ucsb.edu

S. Alireza Aghili
Department of Computer Science
University of California
Santa Barbara, CA 93106.
aghili@cs.ucsb.edu

Amr El Abbadi
Department of Computer Science
University of California
Santa Barbara, CA 93106.
amr@cs.ucsb.edu

ABSTRACT

Range queries seek the objects residing in a constrained region of the data space. An XML range query may impose predicates on the numerical or textual contents of the elements and/or their respective path structures. In order to handle content and structure range queries efficiently, an XML query processing engine needs to incorporate effective indexing and summarization techniques to efficiently partition the XML document and locate the results. In this paper, we describe a dynamic summarization and indexing method, FLUX, based on Bloom filters and B^+ -trees to tackle these problems. We present the results of extensive experimental evaluations which indicate the efficiency of the proposed system.

1. INTRODUCTION

In recent years, XML has gained wide acceptance as an emerging standard and is being employed as a key technology for data exchange, integration and storage of semi-structured data. The XML data model, due to its rich presentation (content and semi-structuredness) poses unique challenges to effectively support complex queries. Numerous research efforts have been conducted [1, 2, 6, 8, 13, 16, 17, 18, 19, 20, 22] to provide powerful and flexible query capabilities to extract patterns from XML documents. Queries on such ordered trees generally impose predicates on the content of **ELEMENT labels** (keyword search) and/or their corresponding *structural relationships* (structure pattern search). These queries require the presence of some keywords in the document tree along with the conformation of the result with some structural predicates, which might be a specific linear path structure (*path expression query*, e.g. Q_1) or a subtree pattern structure (*twig query*, e.g. Q_2) in the underlying data. For instance, $Q_1 = /dblp//article[/year = '2005']$ represents a simple *path expression query* which matches all the journal articles published in the year 2005 from the **dblp** [15] bibliography database.

XML query languages [7, 11] provide support for *content-and-structure (CAS)* class of queries. Additionally, full-text keyword search techniques [5] have been added to XML query languages to support more sophisticated full-text content retrieval. Furthermore, the XQuery and XPath query languages provide support for *exact* range queries which are one of the basic functionalities supported by general database query processing engines.

In this paper, the class of *content-and-structure (CAS)* single path queries are extended to include (i) *range* predicates, and (ii) *fuzzy content-and-structure* predicates and, furthermore, efficient techniques are proposed for processing them. We call this class of holistic queries *Fuzzy-Range (FuR)*, since they provide efficient support for *exact* and *fuzzy (approximate)* matching of queries with *path*, *content* and *range* predicates. The fuzzy/approximate matching feature is supported without requiring any specific instructions from the user. For example, the query

$$Q_2 = /dblp//article[2004 \leq /year \leq 2005]$$

matches the journal articles published in the year range [2004–2005] from **dblp** database. *Fuzzy structure matching* feature of the system additionally reports those instances that match the range predicate of the query, but whose structure *resembles* the query's structure, for instance

$$p_1 = /dblp/article/year/2005, \text{ and}$$
$$p_2 = /dblp/articles//year//2004.$$

Due to the heterogeneity of the XML data, it is essential to provide the support for *fuzzy matching* in current XML query engines. In the real world, users may pose a query such as Q_2 on a collection of bibliography XML databases. However, the datasets of the collection may neither have nor conform to a common DTD (Document Type Definition). There may exist journal article records in the database published in the year 2004 to 2005, however, not necessarily conforming to the order and the structure imposed by the Q_2 . An XML query system should be able to hide the details and complexity of the underlying data from the user while

returning both the exact answers to Q_2 , as well as, approximate answers assuming that the user might not have had complete/exact knowledge of the low-level representation of the underlying data in the database.

The efficient evaluation of such FuR queries is determined by the choice of an efficient execution and data access plan which is one of the critical responsibilities of the database optimizer. For instance, consider a possible query plan for Q_2 where the query engine has to perform the FuR query

$$Q'_2 = \text{/dblp//article[/year = 2004 OR /year = 2005]}$$

to find all the journal articles published in the year 2004 or 2005. The `dblp` dataset contains 111,609 instances of the `[/article/year]` path structure and only 259 instances of `[/year/2005]` (see Figure 4). That is the `[/year = '2005']` predicate will return 259 instances while the structure predicate `[/dblp//article/year]` results in 111,609 instances. Hence, it is essential to utilize the *selectivity*¹ of the structural elements and their contents for efficient evaluation of the FuR queries. As a result, an efficient query execution plan should apply the evaluation starting at the more selective segments of the query. However, one of the main challenges involved in such execution plans for FuR queries is that numerical predicates, which happen to be more selective in this case, typically involve the leaf level of the XML document tree. Therefore, pushing the evaluation down to the leaves of the tree should be accompanied with the appropriate leaf indexing techniques to avoid inspecting a large number of leaf nodes. It is clear that plans such as Q'_2 do not utilize the common optimization technique of pushing down the `selection` operation down to the leaves of the query plan tree. Ignoring the selectivity of the path elements results in the exponential growth of the intermediate result set which must be retrieved from the database. We argue that it is essential to utilize effective summarization and indexing techniques to eliminate the space based on the numerical and most *selective* attributes of the XML document collections.

In this paper, we develop an XML query processing system for FuR queries named FLUX. FLUX employs an efficient B^+ -tree based index structure to locate the leaf matches n_i to the range predicate of a query in its initial stage. Each leaf match n_i of the document tree stores a *compact path signature* of the root-to-leaf path structure of n_i , using the notion of Bloom Filters [4]. In the next step, the path signatures of each matched leaf instance n_i is compared with the query's path signature to eliminate those instances whose path signature is very different from that of the query. To the best of our knowledge, this is the first work to specifically address exact and approximate matching of FuR class of queries in XML document collections. The main features of FLUX are summarized as follows:

- FLUX is an efficient and effective system for exact and approximate matching of *Fuzzy-Range (FuR)* class of queries with *range* and *fuzzy structure* predicates on XML document collections.
- An efficient B^+ -tree based indexing scheme is constructed on the *indexable* (e.g., textual, numerical, date,

¹The fraction of the structural elements that satisfy the predicate.

etc.) attributes of the XML document for effective retrieval and matching of the query's range predicate.

- FLUX incorporates a novel bit-wise hashing scheme based on the notion of Bloom Filters [4] on `ELEMENT` and `ATTRIBUTE` contents of XML document trees. A family of hash functions are applied on the path content components where each path is summarized to a compact bit vector signature. As a result, the *fuzzy path matching* is performed very efficiently through the comparison of path signature bit vectors.
- Extensive experimental evaluations depict the effectiveness of FLUX for complex FuR queries on real and synthetic XML document datasets.

The rest of the paper is organized as follows: Section 2 presents the problem definition. Sections 3 and 4 provide the descriptions of range matching and path matching procedures, respectively. Section 5 finalizes the FLUX algorithm followed by Section 6 which provides the experimental analysis and results. Section 7 concludes the work.

2. PROBLEM FORMULATION

XML documents are rooted and ordered tree structures where each node in the document tree corresponds to the document's `ELEMENT`, `ATTRIBUTE`, or `TEXT` nodes. The `TEXT` nodes represent the values of their parent `ELEMENT` nodes, and `ATTRIBUTE` nodes introduce branches off their parent `ELEMENT` nodes. For now, we only consider simple Numerical-Path expressions (NaP) which are defined as follows:

DEFINITION 1. (Numerical-Path² Expression, NaP). A simple path expression $p = e_1t_1e_2t_2\dots e_kN$ is called a Numerical-Path expression (NaP), where each e_i denotes an Ancestor-Descendant (AD, //) or Parent-Child (PC, /) edge, t_i denotes the tag of an `ELEMENT` or `ATTRIBUTE`, and N represents a range predicate or a numerical value, respectively.

Example 1. $q_1 = \text{/dblp//article/[2004}\leq\text{year}\leq\text{2005]}$, and $q_2 = \text{/management//employee/salary/72,000}$ represent NaP expressions on `dblp` [15] and an employee database, respectively. For instance in q_1 : $e_1 = /, t_1 = \text{dblp}, e_2 = //$, $t_2 = \text{article}, e_3 = /$, and $N = [2004\leq\text{year}\leq\text{2005}]$.

Given an XML dataset and a numerical-path expression, we need to locate and retrieve all the qualifying matching instances. Matching the query against an NaP instance of the dataset involves comparing their corresponding path structure and the numerical sentinels of the their path expression (*range* or *single value*, N). The numerical predicate match of the sentinel N of the query expression seeks all the corresponding instances in the database having numerical sentinel N' such that: $N' = N$ or $N' \in N$. For instance considering the query Q_1 , this phase corresponds to locating all the instances of the `year` attribute being equal to 2004. Furthermore, the path structure signatures of all the matching instances are compared against the query's path

²The techniques proposed in this paper are generally applicable to all **indexable** (e.g., numerical, textual, date, ...) attributes. The term "*Numerical*" is solely used for better explanation of the examples and the comparison schemes.

structure signature. The exact and approximate (fuzzy) instances are finally reported based on their level of similarity to the query. In an offline phase, a hash function f maps each tag of each individual path structure of the database onto a hash value. The hash values of the element tags of each path are collectively combined to construct a single bit-vector signature for the path structure. The similarity among the bit-vector signatures p_i of path structures of the database and the path structure signature of the query, is used to measure the similarity of the path structures of the database p_i to the query.

DEFINITION 2. (Matching Instance). *Given the query expression $q = e_1t_1e_2t_2 \dots e_tN$, and any NaP matching instance $p = e'_1t'_1e'_2t'_2 \dots e'_tN'$ of the database, let $\text{hash}()$ denote a hash function which maps a path structure onto a bit-vector. Moreover, let $h:u \rightarrow 2^{\aleph}$ denote a function on bit-vectors which returns the set of all the “set” bit indices of any bit-vector u . Then, the NaP matching instance p is called a matching instance to q , if*

$$h(\text{hash}(e_1t_1 \dots e_t)) \subseteq h(\text{hash}(e'_1t'_1 \dots e'_t)),$$

and $[N = N' \text{ or } N \in N']$.

where $\text{hash}(e_1t_1e_2t_2 \dots e_t)$ and $\text{hash}(e'_1t'_1e'_2t'_2 \dots e'_t)$ denote the hash bit-vector signatures of the path structures of the query and the database matching instance p , respectively. For instance, given two path structures q and p , where $\text{hash}(q) = 100001$ and $\text{hash}(p) = 101101$, then p is called a matching instance of q because $h(\text{hash}(q)) = \{1, 6\} \subseteq h(\text{hash}(p)) = \{1, 3, 4, 6\}$.

Consider searching for matching instances to the NaP q_2 of Example 1. The NaP expression

$$p = \text{/management//employee/salary/72,000}$$

is the only form of matching instance that can result from an exact path matching scheme. However, in constructing the path structure signature, a fuzzy choice of matching functions f and $\text{hash}()$ from the family of *Locality Sensitive Hashing* schemes [4] would not only help identifying the *exact match* instances such as p but also the *approximate matches* simply *resembling* the query, such as

$$\begin{aligned} p' &= \text{/management//employee//salary//72,000, and} \\ p'' &= \text{//manager/employee//salary/value/72,000.} \end{aligned}$$

NOTATION 1. (Path Components). *Any NaP path expression $p = e_1t_1e_2t_2 \dots e_kN$ consists of two main components, a **path expression component** denoted by $Q^p = e_1t_1e_2t_2 \dots e_k$ and a **numerical predicate (sentinel) component** $Q^n = N$.*

For instance, the NaP expression

$$P = \text{/management//employee[65,000≤/salary≤73,580]}$$

consists of two components: the *path expression component* $P^p = \text{/management//employee/salary}$ and the *numerical predicate component* $P^n = 65,000 \leq \text{/salary} \leq 73,580$. Given an NaP query Q , FLUX proceeds in two different phases, (i) finding the regions in the database satisfying the numerical predicate component Q^n of the query (**range matching**), and (ii) matching the query path component Q^p against the *range-matched* instances of the database (**path matching**). Numerical matching is the initial

step and the results of this stage are passed to the path matching phase for structure matching and refinement of the answers. The following sections provide the details of the range and path matching procedures.

3. RANGE MATCHING

Any range query may benefit from efficient indexing mechanisms to quickly locate and retrieve the interesting portions of the database satisfying the range predicate. Popular indexing techniques such as B⁺-trees and R-trees have been extensively applied to alleviate such problems in the general context of numerical predicate queries. The *range matching phase* of FLUX employs an indexing technique based on B⁺-trees on the *numerical predicate component* Q^n of the query for the effective reduction of the search space.

An offline procedure constructs a B⁺-tree index on the *indexable attributes* (e.g, numerical, textual, date, ...) of the XML document dataset. Figure 1[2] depicts a portion of one such an index tree, constructed on the **age** attribute of a typical XML employee database. Given the constructed B⁺-tree of height L , level 0 of the index represents the range index for the root node and the corresponding range index information for the leaves are present in level $L - 1$. Level L stores the leaf buckets of the actual contents of the designated numerical attribute. For instance, the last bucket stores the **age** content information for two existing **age** values 65 and 66 in the database. Each instance (e.g. **age** = 66) also holds the bit-vector signature of the actual path component leading to this node (details provided in the next section), and its corresponding **ELEM-ID** information. The **ELEM-ID** is the *preorder traversal rank* of the corresponding node in the actual XML document. For instance, the node instance with **age** = 66 has preorder rank of 72, which is shown in the document tree of Figure 1, named as the node ⁷²66. Note that, each individual occurrence of an internal or leaf node has a unique preorder value.

Various *numbering schemes* have been proposed [1, 6, 11, 16] which, instead, associate interval/regional encoding with every node, based on the document order. For instance, each label may consist of (*start, end, level*) values for each node, acquired from the *preorder* traversal of the document, which is used to (i) help identify PC or AD relationships, and (ii) impose a logical document order among the nodes. We argue that, it is enough to use the preorder ranks of the nodes to impose the document order. Moreover, each node is associated with a parent pointer in order to locate its parent node. Given a leaf instance node n_i , the parent pointer $\text{parent}(n_i)$ is used to construct the complete leaf-to-root path originating from n_i . This complete path structure is constructed in the last stage of the path matching phase as the final round of path comparison.

4. PATH MATCHING

Given a NaP query Q and the range-matched instances p_i of the database, the *path matching phase* performs the necessary steps to identify those path structure instances p_i whose path component p'_i matches the path expression component Q^p . In the offline phase, each path expression of the database is hashed-mapped and summarized by a compact bit-vector signature through collectively applying a family of hash functions on the element contents of each path. The following section describes the notion of Bloom filters [4]

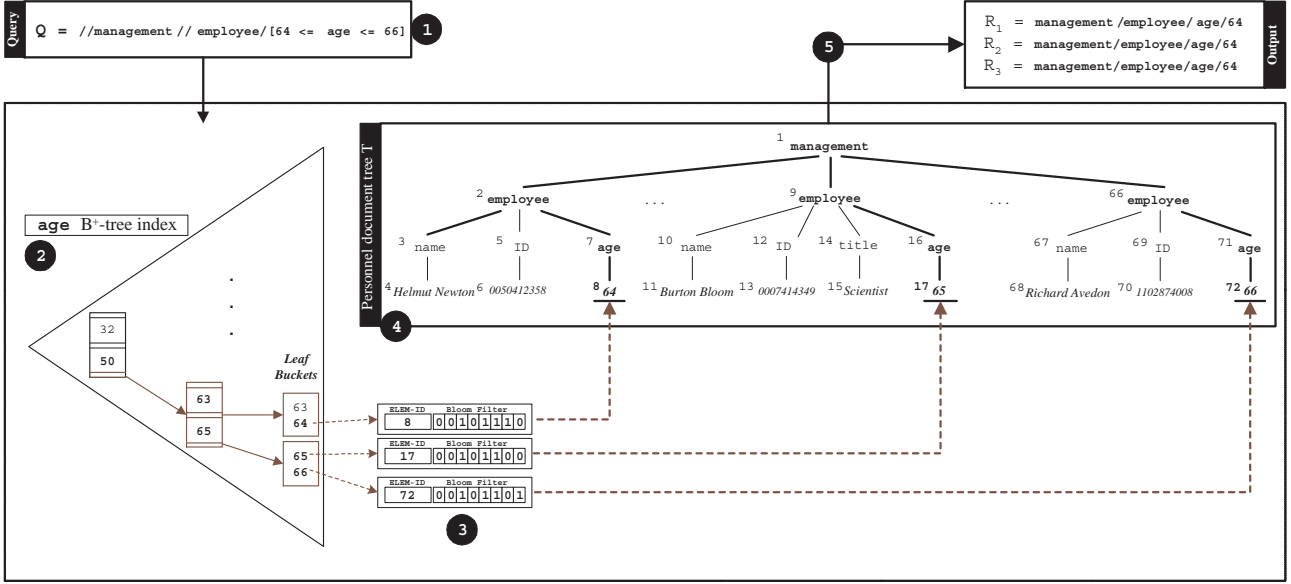


Figure 1: FLUX Search model.

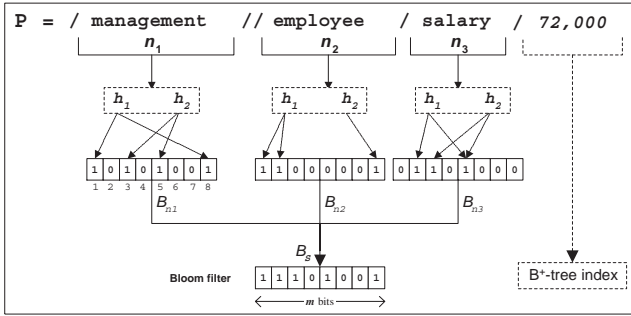


Figure 2: A Bloom filter example.

which is used to map the dataset’s path components into bit-vector signatures. It introduces the theory of bloom filters and the motivations behind incorporating them.

4.1 Bloom Filter: Preliminaries

Bloom filter is a space-efficient data structure to probabilistically represent a set and its elements to support highly accurate *set membership* queries [4]. The bloom filter B consists of a bit vector of length m , and a family of k independent hash functions. Given a set $S = \{n_1, n_2, \dots, n_{|S|}\}$, a family of hash functions is used to construct a bit-vector signature for S . Figure 2 depicts the construction of a bloom filter bit-vector signature using $k = 2$ independent $m = 8$ -bit hash functions h_1 and h_2 , on the path set $S = \{\text{management}, \text{employee}, \text{salary}\}$, where $n = |S| = 3$, from an employee database.

In general, given each element $n_i \in S$, the family of hash functions h_j ($1 \leq j \leq k$) are used to map n_i into a bit-vector. All the entries of the bit vector are initially set to zero. In order to construct the desired bloom bit-vector B_{n_i} , all the k hash functions h_j are applied to n_i . The application of each

h_j on n_i results in “setting” some entries of B_{n_i} to 1. For instance in Figure 2, the application of hash function h_1 on n_1 , $h_1(n_1 = \text{management})$ sets the 1st and 8th bits of the corresponding bit-vector B_{n_1} . Similarly, $h_2(\text{management})$ sets the 3rd and 5th bits of B_{n_1} . To construct the bloom bit-vector for the whole set $S = \{n_1, n_2, \dots, n_{|S|}\}$, the resulting bit-vectors B_{n_i} are combined to form the bloom bit-vector B_S . The combination of the bit vectors B_{n_i} may be performed through a simple logical OR operation. That is, the bit vectors resulting from the application of h_1 and h_2 on the path element *management*, *employee*, and *salary* of Figure 2, are combined using a logical OR function to construct the bloom bit-vector signature $B_{P^\rho} (=B_S)$ for the path component $P^\rho = \text{/management//employee/salary}$. The i^{th} entry of B_{P^ρ} is set to 1 if and only if the i^{th} bit vector entry of at least one of the path components B_{n_1} , B_{n_2} or B_{n_3} has been set to 1. For instance, in Figure 2 the 8th-bit of the final bloom filter B_S is set because the 8th-bit of B_{n_1} (or similarly B_{n_2}) is set. Note that, such application of bloom filter relaxes the edge requirement as imposed by the query. This feature helps to additionally identify and report those instances whose path structure components are very similar to the query, yet having different edge structure.

Subsequently, to test whether the query’s path component Q^ρ is similar to an instance path component B_{P^ρ} of the database, the same set of hash functions are applied to B_{Q^ρ} and all the corresponding bit-vector entries are set to 1. If all the “set” entries of B_{Q^ρ} match with their counterpart in B_{P^ρ} (that is $h(B_{Q^\rho}) \subseteq h(B_{P^\rho})$), it implies that the database path component B_{P^ρ} is identical to B_{Q^ρ} with some probability. The set of all such path structure instances is a superset of the actual (exactly-matched) answer set.

However, there is a chance of B_{P^ρ} and B_{Q^ρ} being identical while the actual path components P^ρ and Q^ρ are different (e.g. by-chance collisions/similarity of the “set” entries of Q^ρ and P^ρ). In such a case, a filter error (*false positive*) is

Algorithm RangeQuery

Input: 1. XPath Range Query Q of the form Q^p/Q^n or $Q^p//Q^n$ where Q^p = Internal path component and, Q^n = Range predicate component on the *RangeAttribute* in the form: $RangeAttribute = value$ or $LB \leq RangeAttribute \leq UB$.

2. The B'-Tree index on the values of *RangeAttribute* in the database along with the corresponding bloom filter bit-vector signature of the path from root to the *RangeAttribute*.

```

01: Procedure CandidateSelection
02: /* search the RangeAttribute's B'Tree to find the qualifying
03:    candidate bucket nodes satisfying the range predicate */
04: candidateSet = RangeQuery( B'Tree T,
05:                            double LB, double UB );
06: queryBloomFilter = ConstructBloomFilter(Qp);
07: /* initialize result set to empty */
08: resultSet = EMPTY;
09:
10: while ( candidateSet.hasMoreElement ) {
11:   curElement = candidateSet.getNext();
12:
13:   /* filter irrelevant paths by using bloom bit-vector */
14:   if ( BloomFiltration( curElement.bloomFilter,
15:                        querybloomFilter, t ) ) then
16:
17:     /* load the actual path ending at the document tree
18:        node whose preorder rank is equal to
19:        curElement.ID, for comparison */
20:     candidPath = loadPathFromDocTree( docTreeFile,
21:                                       curElement.ID );
22:     if (compare( Qp, candidPath ) then
23:       resultSet.add(candidPath);
24:     end if
25:   end if
26: end while
27: return resultSet;
28:
29:
30: Procedure BloomFiltration( BitVector u, BitVector v,
31:                           Threshold t )
32: /* Find out the similarity of the bloom filter u and v */
33: if ( u == v OR d(u,v) >= C )
34:   return true;
35: else
36:   return false;

```

Figure 3: FLUX Algorithm.

said to have occurred. The performance of the hash functions of the bloom filter depends on the *filter error ratio*, which is proven by B.H. Bloom [4] to be as follows. Let n be the number of nodes (or elements) in the set S (or path component P^p), m the size of the bit vector and k the total number of hash functions. Filter error ratio (α) is defined:

$$\alpha = \left(1 - \left(1 - \frac{1}{m} \right)^{nk} \right)^k \approx \left(1 - e^{-\frac{kn}{m}} \right)^k.$$

For instance, the formula suggests that the filter error for a set of $n = 3$ nodes and a Bloom filter of $m = 8$ bits with $k = 2$ hash functions, calculates to $\alpha = 0.028$. That is a filter error would occur with a probability of 2.8%. Moreover, one of the most interesting features of the bloom filter is that it guarantees not to incur any *false negatives* while being highly accurate and very space-efficient.

One of the shortcomings of this approach is the lack of support for updates. To tackle this problem, Fan et al. [10] propose *Counting Bloom Filters*, where each entry i of the signature B_{P^p} of path components in the database is associated with a *counter* c_i where the counter of the i^{th} entry of B_{P^p} is equal to c_i , which means that there are exactly c_i

element nodes in P^p which set the i^{th} bit of the bloom signature vector B_{P^p} . Hence, deletion of a node e_t in the path component (e.g., $e_t = \text{management}$) would be reflected on the bloom signature of the path component P^p , by simply decrementing all those corresponding counters which were incremented by 1 when applying the hash functions h_i on the node e_t . This procedure removes the contribution of the node e_t ($\in P^p$) to the bit-vector signature of the path component P^p . These counters may also be used to estimate the *selectivity* of the element contents along each path.

4.2 Usage of Bloom Filter in FLUX

Given an XML document, in the *offline* phase, all the root-to-leaf path structures of the document tree are extracted. Next, the bloom bit-vector signature of the path component of each such path structure is constructed and stored in an offline profile for each given document tree. The bloom representation of each path structure facilitates an efficient mechanism to compare each path component of the document tree against their counterpart in the query. We now introduce the overall procedure of the FLUX algorithm which combines the features of range matching and path matching schemes.

5. FLUX ALGORITHM

Given the document tree T , the offline phase starts by performing a *preorder traversal* on T and assigns preorder ranks (ELEM-ID) to each node of T (the number on the top-left of each node in Figure 1). These preorder ranks create a virtual document order. Figure 3 depicts the algorithmic details of the FLUX procedure which consists of five individual phases, described in the following.

1) Offline Index Creation. The FLUX offline manager constructs a B^+ -tree index structure on the *indexable attributes* of the XML document collection (e.g. *age*, *salary*, *year*, and *date*). The leaves of each such B^+ -tree store the attribute content (e.g., *age* value), ELEM-ID, and the bloom bit-vector signature of the root-to-leaf path structures of the corresponding nodes. For instance, the node corresponding to *age* = 64 at the leaf bucket level of B^+ -tree of Figure 1[2] stores the *preorder rank* ELEM-ID (e.g. 8 in this instance) of the actual node of the document tree whose *age* attribute has the value 64. Moreover, it stores the bloom bit-vector signature of the root-to-leaf path structure ending at that particular node. For instance, for the node *age* = 64 located at the B^+ -tree leaf bucket of Figure 1[2], the bit-vector 00101110 represents the bloom signature of the root-to-leaf path structure $/^1\text{management}/^2\text{employee}/^7\text{age}$ of the node 864 of the document tree in Figure 1[4], where the numbers 1,2 and 7 denote the ELEM-IDs of the element tag instances of *management*, *employee* and *age* element nodes, respectively.

2) Query Segmentation. This phase segments the query expression $Q = //\text{management}//\text{employee}/[64 \leq \text{age} \leq 66]$ into the path component $Q^p = //\text{management}//\text{employee}/\text{age}$ and the numerical predicate component $Q^n = [64 \leq \text{age} \leq 66]$.

3) Range Lookup. The search part of this phase corresponds to the lines 2-5 of the algorithm in Figure 3. The corresponding B^+ -tree of the *age* range attribute is searched for potential candidate bucket nodes matching the predicate in Q^p (e.g. nodes 64, 65 and 66 in Figure 1[2] for $[64 \leq \text{age} \leq 66]$).

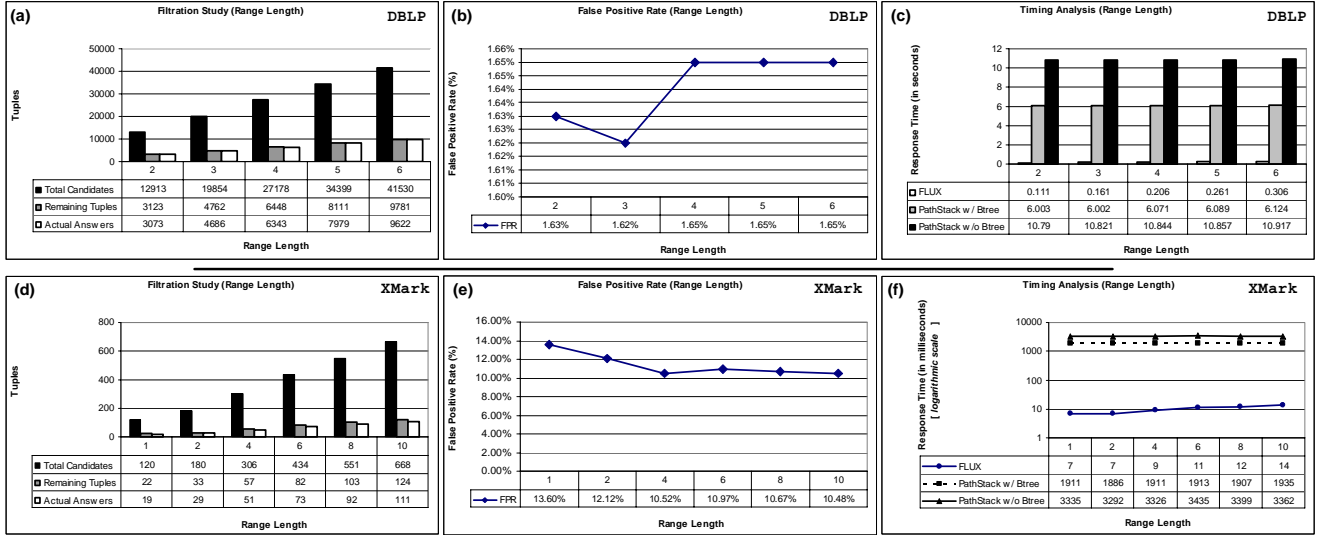


Figure 5: Effect of range length variation on the filtration, FPR and response time.

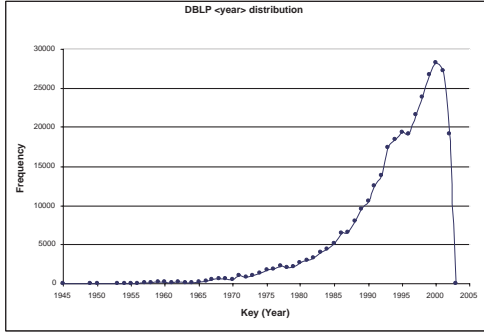


Figure 4: The $\langle year \rangle$ element distribution of dblp dataset.

4) **Path Matching and Filtration.** Let B_{p_1}, \dots, B_{p_k} denote the bloom signatures of each of the k matches of the database (e.g., the bloom filter of the path component `/1management/2employee/7age` which end at node ⁸64), whose numerical contents have already been matched with the query’s numerical predicate Q^n . This stage is responsible for matching the path component of the query B_{Q^p} against the path components of the range-matched instances B_{p_1}, \dots, B_{p_k} . It ranks each matching instance B_{p_i} based on its similarity to B_{Q^p} . The path matching procedure corresponds to the invocation of the `BloomFiltration()` function at lines 14-15 of Figure 3 where its definition is provided at lines 30-36. After filtering out the false positives, the `candidPath` holds the results of matches to Q in the database. The function d in line 33 is the feature of the algorithm which facilitates the *fuzzy (approximate)* matching of the bloom filters. More specifically, the function $d(u, v)$ returns the number of “*set*” bit entries shared by u and v . For instance, in Figure 2 $d(B_{n_1}, B_{n_2}) = 3$ because these signature are different in three bit entries, that is the 2^{nd} , 3^{rd} , and 5^{th} bits. Finally, the actual path structures of the non-filtered matches are constructed (using the node pointers from leaf-

to-root), compared against the query and reported to the user. Note that, the set of additional matching instances reported through function d and the set of false positives constitute the set of approximate answers which are reported to the user at no extra cost. However, our experimental evaluations are only focusing on the analysis of exact matching.

6. EXPERIMENTAL EVALUATIONS

We implemented the FLUX system using *Java 1.4.2* and ran our experiments on a *Pentium M-2GHz* processor with *2GB* of main memory, using a page size of *1KB* (determine the number of indexed data items which a leaf node can have and the number of key/pointers which an internal node can have for the B^+ -tree.), cache size of *100KB*, and LRU cache replacement policy. We compared our proposed technique with PathStack [6] which is the best in the literature for simple XPath queries. Also PathStack is implemented by using *Java 1.4.2*. Two variations of PathStack are used when retrieving the XML document elements residing in the range specified in the query for the structural join: one variation is to use B^+ -tree index and the other variation is not to use B^+ -tree index.

The experimental evaluations were performed on a set of both synthetic (XMark [23] containing information about an auction site) and real (dblp³) XML datasets. The dblp dataset (sized of 127MB) constituted 3,332,130 element nodes with an average and maximum depth of 2.9 and 6, respectively. We generated a set of synthetic XMark datasets with scaling factor ranging from 0.1 to 1.2 for the experimental evaluation. The average depth for the XMark datasets is 5. Moreover, different amount of random noise was imposed on the dblp and XMark datasets to create path structure variation at the element names. The number of hash functions used for constructing the bloom filter is 4. For each element

³Acquired from the University of Washington’s XML Data Repository accessible through <http://www.cs.washington.edu/research/xmldatasets/>

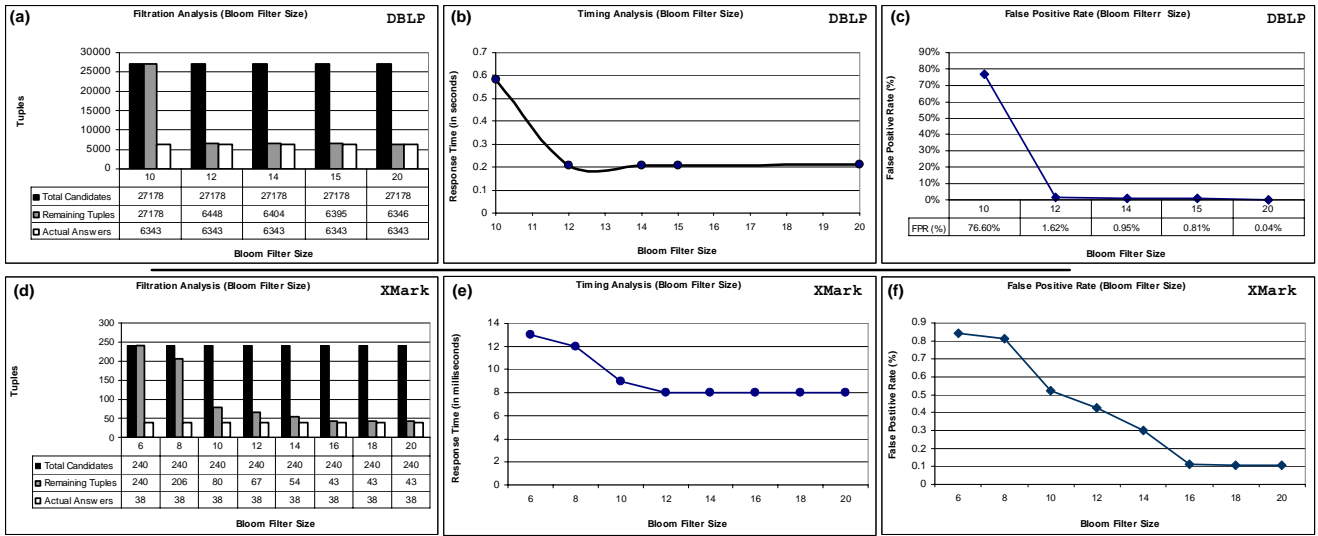


Figure 6: Effect of the size of bloom filter signature (in *bits*) on filtration, FPR, and response time.

of along the path which leads to an instance of the range attribute, its MD5 digest (a 128-bit cryptographic message) is computed. This 128-bit message is evenly divided into 4 groups. Each 32-bit group is further transformed into an integer ranging from 0 to the bloom filter size - 1. Unless otherwise stated, the bloom filter size was chosen to be 14 bits for dblp dataset and 16 bits for XMark datasets which will be explained later in this section.

The results presented in this section were generated by averaging the results from running a workload of 100 random queries on dblp and XMark datasets. The dblp query template was chosen as $Q_D = \text{dblp/article}[\$LB \leq \text{year} \leq \$UB]$, for different random values of $\$LB$ and $\$UB$. Similarly, the XMark query template was selected as $Q_X = \text{regions//item//mail}[\$LB \leq \text{date} \leq \$UB]$. The range values $[\$LB, \$UB]$ were chosen randomly from the $\langle \text{year} \rangle$ and $\langle \text{date} \rangle$ domain space (e.g., Figure 4) in the year range 1945 to 2003 and date range 01/01/1998 to 12/28/2001. Figure 4 depicts the frequency distribution of the occurrence of the $\langle \text{year} \rangle$ element in the incorporated dblp dataset. The dblp dataset includes 328,831 path instances leading to the year element, which is the reason behind using Q_D as the query template for dblp dataset while it provides a large candidate set. The richness of the path structure which lead to the $\langle \text{date} \rangle$ element is the reason behind to choose Q_X as the query template for XMark dataset (more structural variations on Q_X can be applied for structure effect study). Following are some notations used in the upcoming figures:

- **Total Candidates:** Number of all the possible year instances (dblp) and date instances (XMark) in the database for the inspected range resulting from the range query search on the B^+ -tree index structure lookup phase.

- **Remaining Tuples:** The number of candidates left for further inspection after pruning the intermediate results by comparing their bloom filter signature against the bloom filter signature of the query.

- **Actual Answer:** The number of actual answers in the database to the query.

- **False Positive Rate (FPR):** The FPR is calculated as $(\text{RemainingTuples} - \text{ActualAnswers}) / \text{RemainingTuples}$, which indicates how close the filtration gets to the actual answer set.

Figures 5-9 analyze the effect of *range length*, *bloom filter size*, the imposed *noise*, and the *scalability analysis* on the *Filtration*, *False Positive Rate (FPR)* and *Response Time* effectiveness of FLUX, on the dblp and XMark datasets, respectively.

6.1 Effect of Range Length

Figure 5 depicts the effect of the range length $r = |\$UB - \$LB|$ on the performance of FLUX on dblp and XMark (scaling factor = 1, size $\approx 113\text{MB}$ and noise = 30%) datasets. The query's range length/extent is varied from 2 (narrow) to 6 (moderately wide), and 1 (narrow) to 10 (wide) on the dblp⁴ and XMark datasets, respectively. FLUX succeeds in pruning a substantial fraction of the candidate result set in the bloom filter comparison phase. For instance, in Figure 5(a), the column pertaining to $r = 3$ indicates that the application of bloom filtration reduces the number of total candidates from 19854 tuples to 4762 tuples, or in other words, to 24% of the total candidate result set. A consistent invariance to the range length on the performance of the filtration is observed on both dblp and XMark datasets (Figure 5(a) and 5(d)). Figures 5(c) and 5(f) depict the total response time of performing the designated operations, as a function of range length compared with PathStack [6] (with and without using B^+ -tree index structure). The running time of FLUX consistently outperforms PathStack on both dblp and XMark datasets. For instance, in Figure 5(f) FLUX performs 100-times faster on average when compared to PathStack (with B^+ -tree index structure). Figures

⁴e.g. $1999 \leq \text{year} \leq 2003$ has the range extent of $r = |2003 - 1999| = 4$

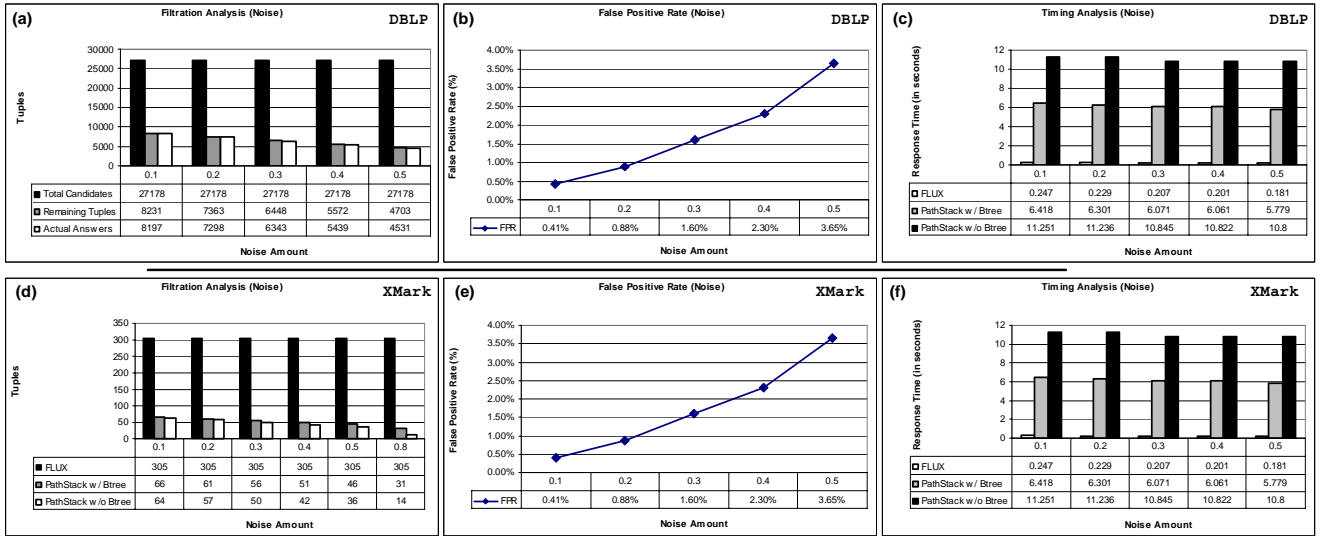


Figure 7: The results of applying random noise with various intensity on the element names.

| Query | FLUX | PSB | PS |
|--|------|------|------|
| $Q_1 = \text{regions//mail/date}$ | 7.9 | 1521 | 2937 |
| $Q_2 = \text{regions//item//mail/date}$ | 8 | 1901 | 3323 |
| $Q_3 = \text{regions//item/mailbox/mail/date}$ | 8.1 | 2307 | 3708 |

Table 1: Response time (in milliseconds) comparison of FLUX v.s. PathStack on XMark dataset varying the query structure. PSB = “PathStack with Btree” and PS = “PathStack without Btree”.

8(a)-(b) further demonstrate the speedup achieved by using FLUX against PathStack (with B^+ -tree index structure) as the range length varies on both dblp and XMark datasets. The speedup ranges from 54 times (at $r = 2$) to 19 times (at $r = 6$) on the dblp dataset, and consistently around 100 times faster on the XMark dataset. Despite the declination of the speedup⁵ as the range length increases, the observed speedup depicts a 20-times faster response time in the worst case. Figures 5(b) and 5(e) depict the stability of False Positive Rate (FPR), which stays within 2% of the remaining tuples as the range length varies for dblp dataset and 14% for the XMark dataset.

6.2 Effect of Bloom Filter Size

Figure 6 analyzes the effect of bloom filter size (in bits) as it varies from 10 to 20 bits and 6 to 20 bits on dblp and XMark datasets. The XMark dataset of this section was generated with a scaling factor of 1, with about 113MB in size and 30% imposed noise at the path element names. Figures 6(a) and 6(d) validate the intuitive expectation that the larger choice of the bloom signature length should result in more effective filtration. Figures 6(b) and 6(e) depict the response time analysis of FLUX when varying the bloom bit-vector size in answering the same set of 100 random queries on each respective dataset. The filtration (Figures 6(a) and 6(d)) and response time (Figures 6(b) and 6(e)) performance

⁵This artifact is because the larger range lengths/extents need more bloom filter calculations.

of FLUX improves consistently as the size of the bloom bit-vector increases from 10 to 14 bits for the dblp dataset and 10 to 16 bits for the XMark dataset. This is due to the fact that, the chance of bloom signature collision⁶ reduces as the size of the bloom signatures increases. When the bloom bit-vector increases from 14 to 20 for the dblp dataset and 16 to 20 for the XMark dataset, the filtration effectiveness still increases while the query response time does not due to the fact that larger size of bloom filter will incur more time to retrieve the corresponding data. Hence we choose 14 bits and 16 bits for the dblp dataset and 16 bits for the XMark datasets for constructing bloom filters in a timely manner. Moreover, Figures 6(c) and 6(f) demonstrate the filtration effectiveness of FLUX which is shown in the reduction of FPR when increasing the size of bloom filter. This again coincides with the intuitive expectation on the bloom filter size as mentioned above.

6.3 Effect of Noise in Data

For this set of experiments, we introduced random noise at the element names, varying from 1% to 5% on dblp dataset and 1% to 8% on XMark dataset, respectively. Figure 7 depicts the effect of the imposed noise ratio on the overall performance of FLUX. As expected, the introduction of more noise results in larger FPR as shown in Figures 7(b) and 7(e). However, despite the introduction of noise, FLUX performs very efficiently in filtration ratio and response time as observed in Figures 7(a) and 7(d), and Figures 7(c) and 7(f), respectively. FLUX substantially outperforms PathStack regardless of the amount of noise imposed on the data as shown in Figures 7(c) and 7(f). Figures 8(c) and 8(d) show the amount of speedup achieved by using FLUX versus PathStack (with B^+ -tree index structure) which ranges from 26 times (at $\text{noise} = 0.1$) to 32 times (at $\text{noise} = 0.5$) on dblp dataset. Similarly, FLUX consistently outperforms PathStack on XMark dataset with an average of 100 times faster response time. Relative to PathStack, FLUX per-

⁶The probability bloom hash functions assign an identical bloom signature to two different path structures.

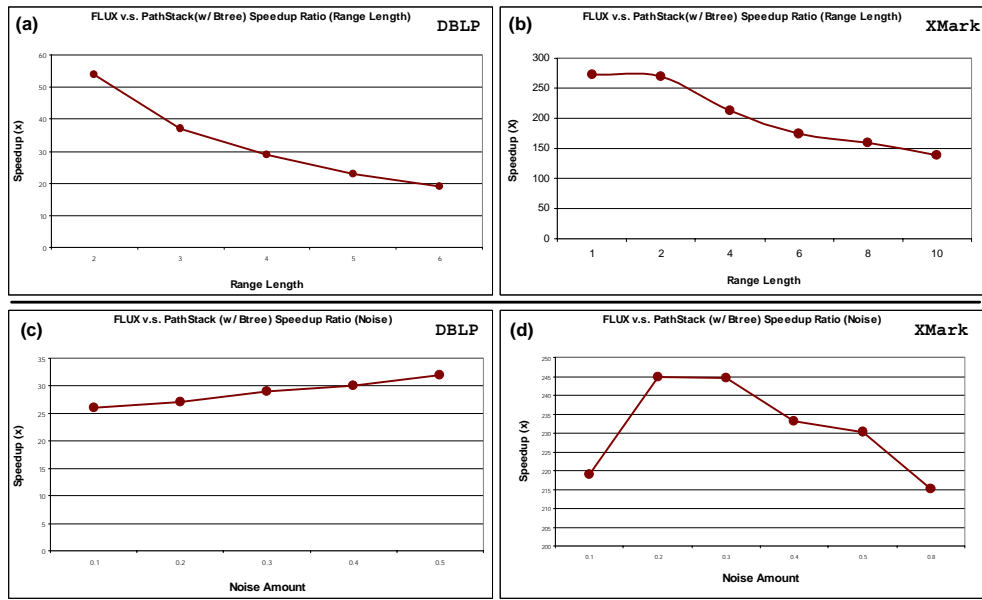


Figure 8: The speedups of FLUX v.s. PathStack (w/ B⁺-tree) when the range length and the imposed noise amount varies.

forms even better when more noise is inherent in the dataset, which is a very desirable feature when the query is posed on datasets with variations in their representation or not necessarily conforming to a unified schema or Document Type Definition (DTD).

6.4 Query Structure Variation

Table 6.1 depicts the response time analysis when varying the query structure in FLUX and PathStack (with and without B⁺-tree index). From type Q_1 to Q_3 , more structures are imposed on top of the range attribute. The results were acquired by averaging the running time of 100 random range queries of type Q_i (of Table 6.1). The range domain was selected in the 01/01/1998 to 12/28/2001 date range and each random range query has the length of 4. The bloom filter size was selected to be 16 bits. The incorporated XMark dataset was generated using a scaling factor of 1 with 30% noise. In all the observed cases, FLUX consistently outperformed PathStack. The performance of FLUX is slightly affected when the path structure of the query tends to get more complicated due to the bottom-up computation approach. The set of the remaining tuples for each type of query is the same after using the bloom filtration. Thus the cost of retrieving the corresponding paths for the remaining tuples for further inspection against the query is approximately the same. However, for PathStack, more structures with the query will incur more document elements retrieved from the disk for the structural join to produce the matching instances of the query. Hence, the performance of PathStack will decrease when more path structures are imposed on the same range attribute.

6.5 Scalability Analysis

In this set of experiments, we generated a set of XMark datasets with scaling factors ranging from 0.1 to 1.2 to study the effects of document size on the effectiveness of FLUX.

Figure 9 depicts the filtration efficiency and response time analysis of FLUX versus PathStack resulted from running a set of the same 100 random range queries selected in the 01/01/1998 to 12/28/2001 date range. The performance of both FLUX and PathStack suffers as the size of the dataset increases, however, FLUX experiences from 98 times to 215 times slower performance degradation rate compared with PathStack with B⁺-tree index structure. The comparison with PathStack without B⁺-tree index structure, is even more dramatic.

7. CONCLUSION

This paper proposed an efficient technique, named FLUX, for answering complex range queries in a database of XML documents. FLUX incorporated a B⁺-tree based index structure on the contents of range attributes. It uses the notion of Bloom filters to associate a structure signature to each range attribute instance. The filtration performed by the bloom signatures of FLUX reduced the search space to a minor fraction of the intermediate result set. Experimental results demonstrate that the filtration, response time, false positive rate, speedup and scalability of FLUX consistently outperforms PathStack [6] on both real and synthetic datasets.

8. REFERENCES

- [1] S. Al-Khalifa, H.V. Jagadish, J.M. Patel, Y. Wu, N. Koudas and D. Srivastava, Structural Joins: A Primitive for Efficient XML Query Pattern Matching. ICDE, 141–152 (2002).
- [2] S. Al-Khalifa et al., Querying Structured Text in an XML Database. SIGMOD, 4–15 (2003).
- [3] S. Amer-Yahia, L.V.S. Lakshmanan and S. Pandit, FleXPath: Flexible Structure and Full-Text Querying for XML. SIGMOD, 83–94 (2004).

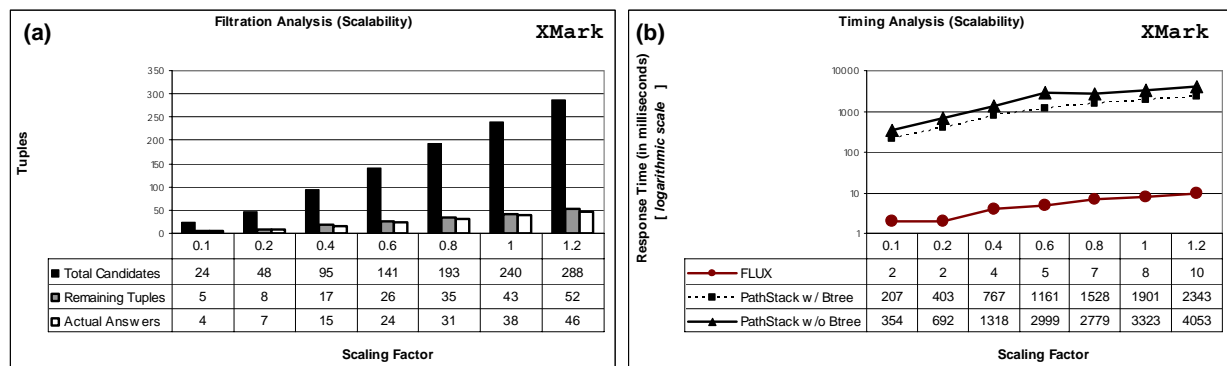


Figure 9: The scalability analysis on XMark datasets

- [4] B.H. Bloom, Space/Time Trade-offs in Hash Coding with Allowable Errors. *Communications of the ACM* **13(7)**, 422–426 (1970).
- [5] C. Botev, J. Shanmugasundaram and S. Amer-Yahia, A TeXQuery-Based XML Full-Text Search Engine. *SIGMOD*, 943–944 (2004).
- [6] N. Bruno, N. Koudas and D. Srivastava, Holistic twig joins: optimal XML pattern matching. *SIGMOD*, 310–321 (2002).
- [7] D. Chamberlin, Daniela Florescu, Jonathan Robie, Jérôme Siméon and Mugur Stefanescu, XQuery: A Query Language for XML. W3C Working Draft, <http://www.w3.org/TR/xquery> (2001).
- [8] S. Chien, Z. Vagena, D. Zhang, V.J. Tsotras and C. Zaniolo, Efficient Structural Joins on Indexed XML Documents. *VLDB*, 263–274 (2002).
- [9] CiteSeer Scientific Literature Digital Library and Search Engine, <http://citeseer.ist.psu.edu>
- [10] L. Fan, P. Cao, J.M. Almeida and A.Z. Broder, Summary Cache: A Scalable Wide-Area Web Cache Sharing Protocol. *SIGCOMM*, 254–265 (1998).
- [11] T. Grust, Accelerating XPath location steps. *SIGMOD*, 109–120 (2002).
- [12] L. Guo, J. Shanmugasundaram, K.S. Beyer and E.J. Shekita, Efficient Inverted Lists and Query Algorithms for Structured Value Ranking in Update-Intensive Relational Databases. *ICDE*, (2005).
- [13] H. Jiang, W. Wang, H. Lu and J. Xu Yu, Holistic Twig Joins on Indexed XML Documents. *VLDB*, 273–284 (2003).
- [14] R. Kaushik, P. Shenoy, P. Bohannon and E. Gudes, Exploiting Local Similarity for Indexing Paths in Graph-Structured Data. *ICDE*, 129–140 (2002).
- [15] DBLP Bibliography Server, <http://dblp.uni-trier.de/>
- [16] Q. Li and B. Moon, Indexing and Querying XML Data for Regular Path Expressions. *VLDB*, 361–370 (2001).
- [17] J. Lu, T. Chen and T.W. Ling, Efficient Processing of XML Twig Patterns with Parent Child Edges: A Look-ahead Approach. *CIKM*, 533–542 (2004).
- [18] J. Lu, T.W. Ling, C.Y. Chan and T. Chen, From Region Encoding to Extended Dewey: On Efficient Processing of XML Twig Pattern Matching. *VLDB*, 193–204 (2005).
- [19] A. Marian, S. Amer-Yahia, N. Koudas and D. Srivastava, Adaptive Processing of Top-K Queries in XML. *ICDE*, (2005).
- [20] P. Rao, B. Moon, PRIX: Indexing And Querying XML Using Prüfer Sequences. *ICDE*, 288–300 (2004).
- [21] F. Weigel, H. Meuss, K.U. Schulz and F. Bry, Content and Structure in Indexing and Ranking XML. *WebDB*, 67–72 (2004).
- [22] H. Wang, S. Park, W. Fan and P.S. Yu, ViST: A Dynamic Index Method for Querying XML Data by Tree Structures. *SIGMOD*, 110–121 (2003).
- [23] A. R. Schmidt et al., The XML Benchmark Project. Technical Report INS-R0103, CWI (2001).