

PIX: XML 문서 검색을 위한 분할 색인 기법

PIX: Partitioned Index for Keyword Search over XML Documents

초록

XML 문서 검색을 위한 분할 색인 기법. 이 논문에서는 XML 문서에 대한 키워드 검색을 위한 새로운 색인 기법을 제안한다. 제안된 기법은 XML 문서의 구조를 고려하여 데이터를 분할하고, 이를 효율적으로 색인화한다. 이를 통해 기존 방법보다 검색 성능을 향상시키고, 대용량 XML 문서에 대한 검색을 지원한다. 제안된 기법은 XML 문서의 구조를 고려하여 데이터를 분할하고, 이를 효율적으로 색인화한다. 이를 통해 기존 방법보다 검색 성능을 향상시키고, 대용량 XML 문서에 대한 검색을 지원한다. 제안된 기법은 XML 문서의 구조를 고려하여 데이터를 분할하고, 이를 효율적으로 색인화한다. 이를 통해 기존 방법보다 검색 성능을 향상시키고, 대용량 XML 문서에 대한 검색을 지원한다.

1.

XML(Extensible Markup Language) W3C

[1]. 가 SGML

XML 1998

. XML

,

XML

, , , ,

가 XML

. HTML

XML

가

XML

가

XML

XML

. XML

가

RDBMS

SQL

XML

가

XQL[19], XML-QL[20]

가

XPath[21],

XQuery[22] W3C

가

가

가

가

가

가

,

가 XML

가

가

XML

가

가 (passage retrieval) [23,24,25] 가 XML

1 DBLP XML 가
Jagadish가 XML 가

```

<inproceedings key="conf/sigmod/Jagadish90">
  <author>H. V. Jagadish</author>
  <title>Linear Clustering of Objects with Multiple Atributes.</title>
</inproceedings>
<inproceedings key="conf/sigmod/ZhangRL96">
  <author>Tian Zhang</author>
  <author>Raghu Ramakrishnan</author><author>Miron Livny</author>
  <title>BIRCH: An Efficient Data Clustering Method for Very Large Databases.</title>
</inproceedings>
<inproceedings key="conf/dbpl/JagadishLST01">
  <author>H. V. Jagadish</author> <author>Laks V. S. Lakshmanan</author>
  <author>Divesh Srivastava</author> <author>Keith Thompson</author>
  <title>TAX: A Tree Algebra for XML.</title>
</inproceedings>

```

1 DBLP XML

가 Jagadish XML
1
가 Jagadish가 TAX
XML 가 가
XML

[6,7,10,11,12,13].

[2]

XML

가 ,

가

,

가

[3, 5],

[4]

가

XML

가

XML

가

(partition)

2

XML

3

4

5

가

6

2.

XML

XML

가

XML

XQL, XML-QL

XPath가

XQuery

XQuery

“XML

”

(information retrieval)

.1) XML

가

가

1) XQuery Requirements

XQuery FullText

IR

Use Cases

가

가 가

(ranking)가 . XIRQL[6] XQL 가

(weighting) , XML , 가

. XXL [7] (path algebra) [8]

가 (similarity join)

. XML [9]가

. XKeyword[10]

XML 가

XKeyword 가

XML 가

XML

XML

(IR)

. XRANK[11],

XSEarch[12], [13] . XRANK

PageRank[14] 가

(least common ancestor)

DIL , Top-k

RDIL, RDIL

가 HDIL . XSEarch

XML

XML

relationship) . XSEarch (interconnection
 XRANK 가
 XML 가
 가 가 . [13]
 [15] 가
 가 가 .
 XML 가 XRANK
 XML

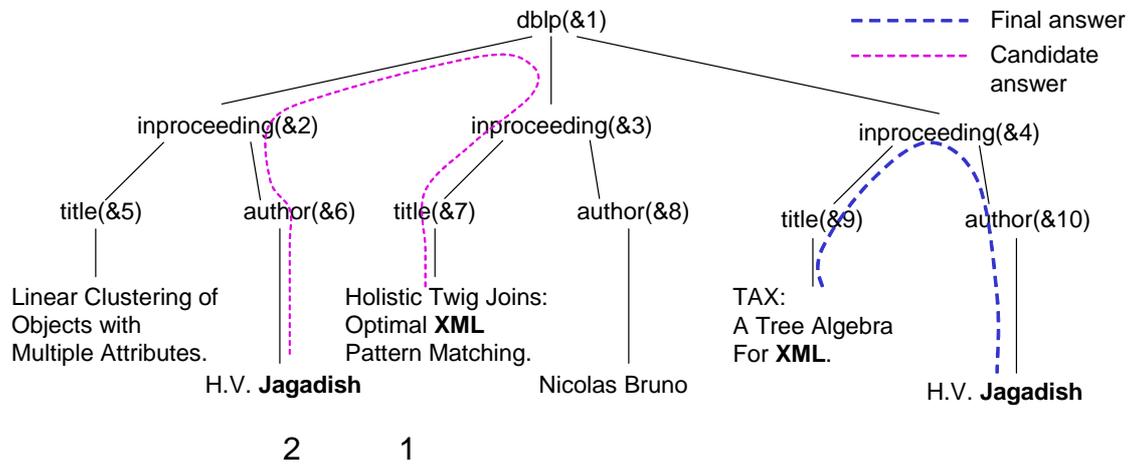
3.

3.1. XML

XML 가 가 ,
 (ordered node labeled tree)
 (document node)가 0
 가 ' ' ' '
 가 가
 가 가 30 ~ 50%

가

[2].



2 1 XML

가

가

3.2.

가

가 AND

가 OR

가

XML

가

XML

(LCA, least common ancestor)

DTD

가

가

XML

DTD

가

가

가

.

가

가

[10,11].

XML

“Jagadish가

XML

”

("Jagadish", "XML")

.

가

2

&4

.

"Jagadish", "XML“

.

&1

가

(&4)

&4

.

4. PIX (Partitioned Index for Keyword Search over XML Documents)

4.1.

(inverted index)

가

.

()

가

.

,

,

가

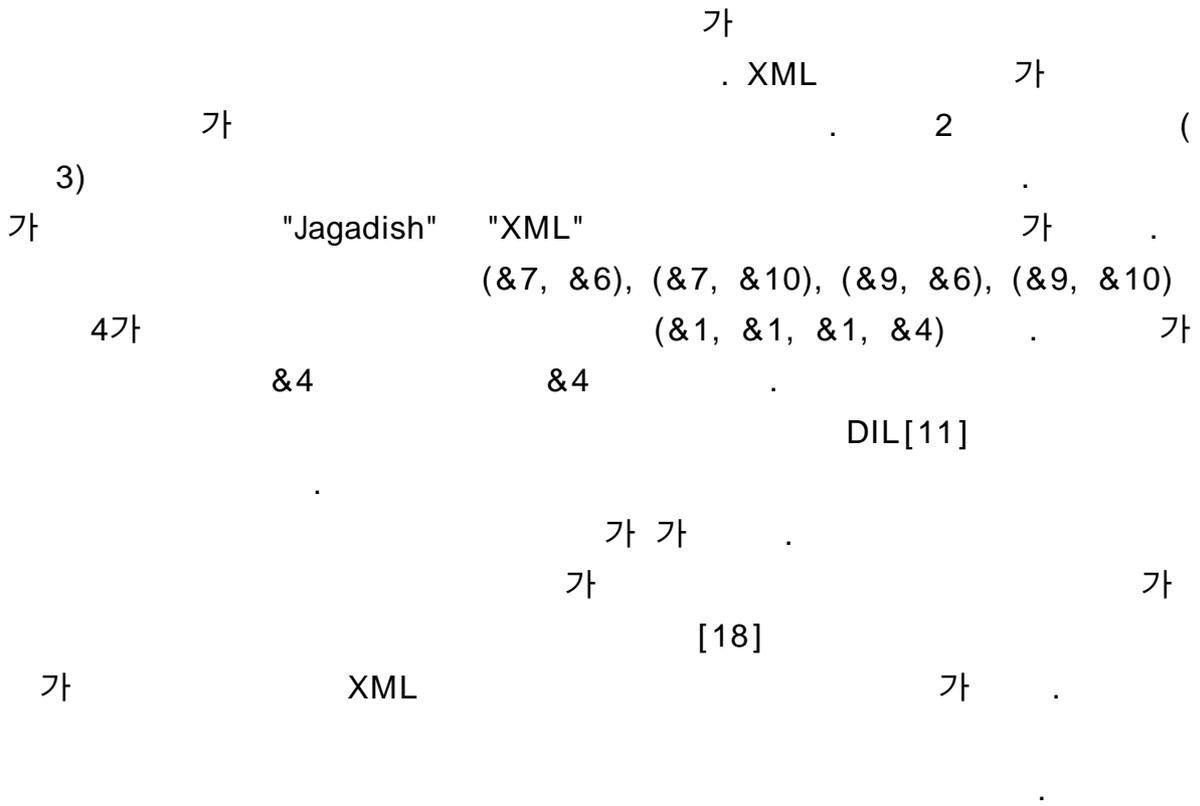
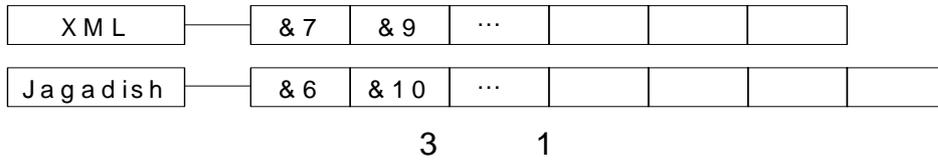
. XML

가

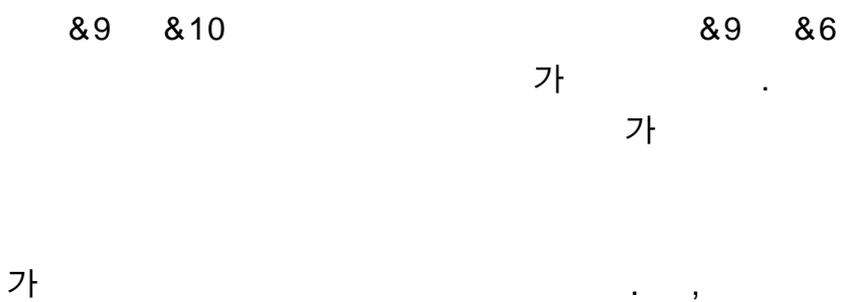
가

3

.



4.2.



(base level)

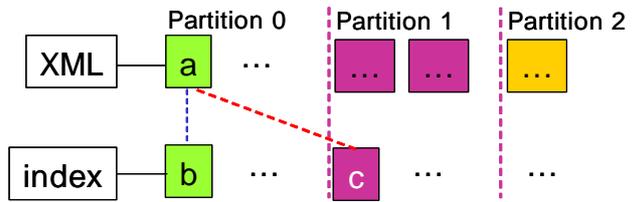
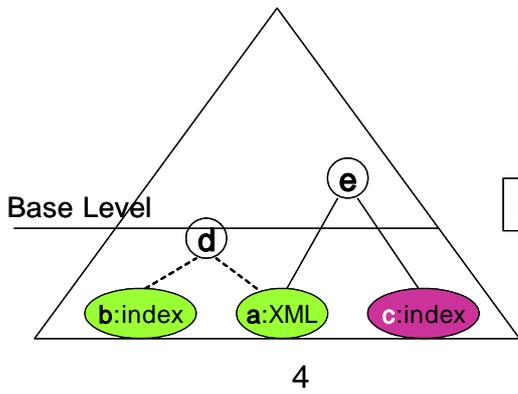
가

Top-k

k

가

가



4

a b

가

a c

가

a c

(partition)

a b

4

a c

PIX

가

가

가

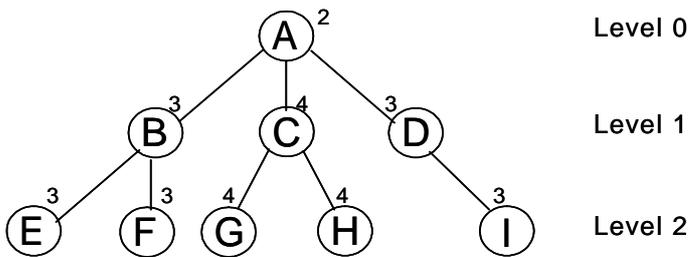
PIX

PF

가 2)

$1 (PF)$ Value \times Node Value
 $PF_i(n)$, i PF_i PF n $PF(i, n)$
 T n_a T , n_b T
 n_a n_b T i 가
 $PF_i(n_a) = PF_i(n_b)$

2



가 $f(B) = f(D) = f(E) = f(F) = f(I) = 3$, $f(C) =$
 $f(G) = f(H) = 4$
 $i = 1$ PF f PF_1 E F
 $f(E) = f(F)$ G I $f(G)$
 $f(I)$ 1 0
 $f(E) = f(I)$ E I 1
 PF_i i PF_i

\cong_i
 PF 가

2) 가 가 가

2 (\cong_i) i
 PF (binary relation)
 $n_a \cong_i n_b$ $n_a \cong_i n_b$
 $j \leq i, PF_j(n_a) = PF_j(n_b)$
 가 가 가
 가 i 가 i 가 i
 가
 1 (equivalence relation)
 가 $n_a \cong_i n_a$
 (reflexive relation) n_a, n_b n_b, n_a
 (symmetric relation) $n_a \cong_i n_b$ $n_a \cong_i n_b$
 i 가 $n_b \cong_i n_c$ $n_b \cong_i n_c$ i
 n_c $n_a \cong_i n_b$ n_x, n_b
 n_y $n_x \cong_i n_y$ $n_a \cong_i n_c$
 가 (transitive relation)

P

PF PF
 가 가

3 (**P**) i

PF_i P-value P-value가 v i
 $P_{i,v}$

(index level)

PIX

PIX

P

PF

i

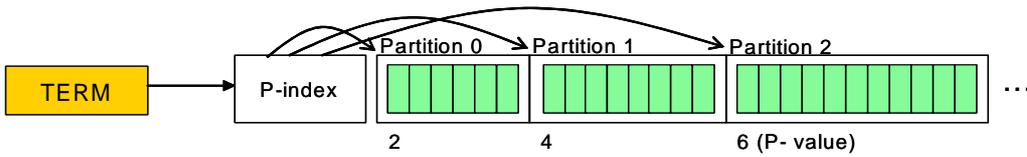
PF_i

PIX

PIX

P-index 가

5



5 PIX

3 XML

6 A

“index” A, B

“XML”

B

1 가 0

(Partition 0), 1

(Partiton 1), 2

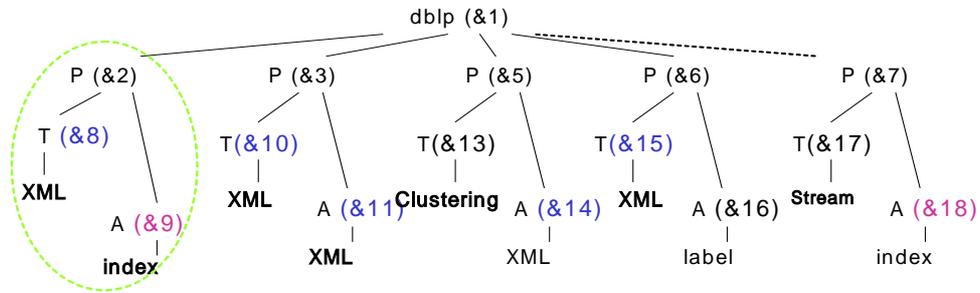
(Partitoin 2)

(&10, &11)

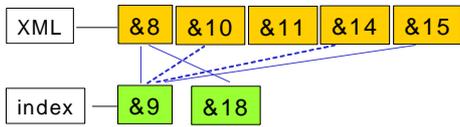
(&9)

가

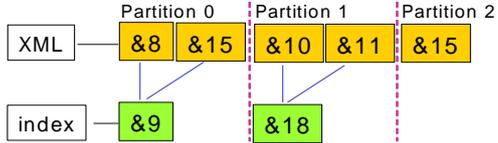
가 1 가 ,
 (&8, &15) (&19) 가 . &8 &9 &2
 &15 &9 .



A. Conventional Inverted Index



B. Partitioned Inverted Index



6

XML

) 가

, PIX

, XRANK[11] DIL

. PIX

P-value

P-value

P-value

가

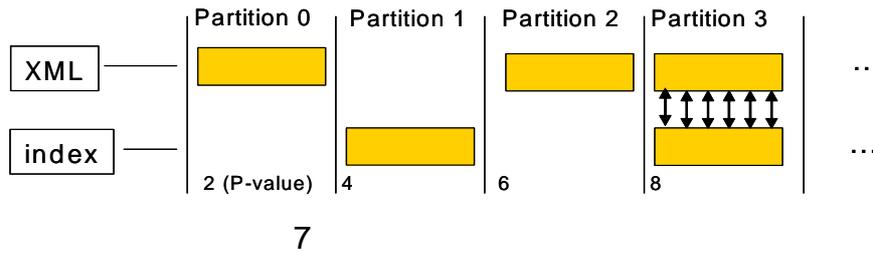
가

가

1. P-index

2. P-value가

3.



XML index . XML P-value가
 2, 6, 8 index P-value가 4, 8
 index P-value 8 . XML
 . XML P-value 2, 6 index P-value
 4

가 . PIX
 가

가 .

가

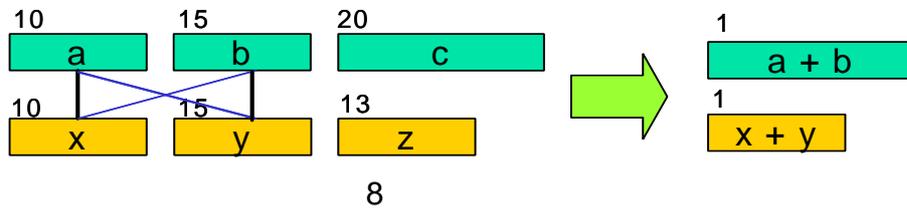
가 .

(recall)

가 P-value 가 P-value P-value

- 1.
2. P-value P-value

가 Top-k (k)



4 XML index XML index P_{4,10}³⁾, P_{4,15}, P_{4,20} , index P_{4,10}, P_{4,13}, P_{4,15} 가 4 가 4 가 3 P-value 1 가 , 15 1, 20 2 가 P_{4,10} P_{4,15} 3) 4 P-value가 10 3

P-value가 3 P-value가 1
 P_{4,10} P_{4,15} P_{3,1} XML P_{3,1}
 index P_{3,1}
 P_{3,1} 4 가
 가 P_{4,10} 4 가 P_{4,15}

가 , 가

PIX

PIX

PIXPF

PIXMERGE

PF

4 PIXPF PIXMERGE

$$PIXPF_i(n) = \sum_{k=1}^i (o_k(n) \bmod 2) 2^k$$

i : i

n :

$o_k(n)$: n 가 2^k 로 나눈 나머지가 0인 경우

PIXMERGE (a, b)

i 번째 비트까지 비교하여 $P_{i,a}$ 와 $P_{i,b}$ 를 구한다.

$$a \bmod 2^i = b \bmod 2^i \quad P_{i,a} = P_{i,b}$$

PIXPF

i 번째 비트까지 비교하여 (sibling) 이 나오면

0 , 1 가 나오면 i 번째 비트를 무시하고 다음 비트를 비교한다.

가 나오면 2^k 를 곱하여 P-value를 구한다.

$$0 \times 2^0 + 1 \times 2^1 = 2 \quad \text{P-value} \quad \text{P-value}$$

가 i 번째 비트까지 비교하여 2^i 를 곱하여 P-value를 구한다.

4 PIXPF PF n_a n_b 가 i 번째 비트까지 비교하여

T n_a n_b 가 i 번째 비트까지 비교하여

n_a n_b 가 i 번째 비트까지 비교하여

n_a n_b 가 i 번째 비트까지 비교하여

n_a n_b $1 \leq k \leq i$ k o_k 가

$$\sum_{k=1}^i (o_k(n_a) \bmod 2) 2^k = \sum_{k=1}^i (o_k(n_b) \bmod 2) 2^k \quad \text{PIXPF}_i(n_a) = \text{PIXPF}_i(n_b)$$

PIXPF PF

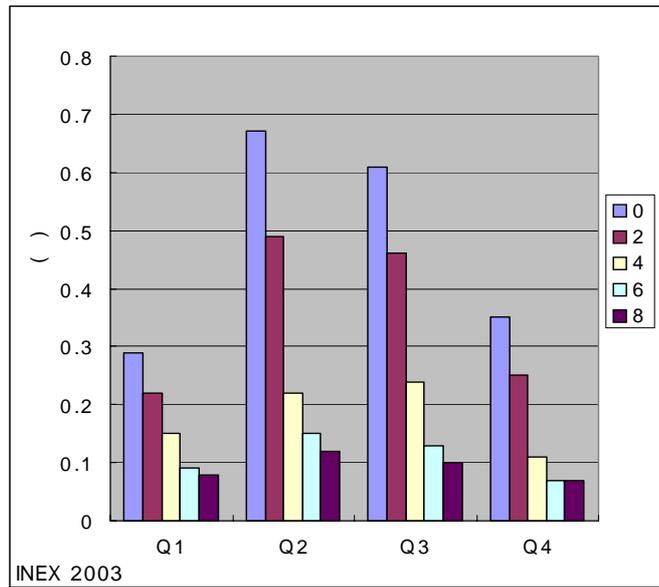
5.

PIX C++ BerkeleyDB[16]

. 512MB 가 PentiumIII 993 MHz PC .
INEX 2003[17] Shakespeare[26] .

가
가 ,
가 가
. PIXPF i 2^i

. INEX 2003 가 . INEX 2003
XML 1995 ~ 2001
XML 500M . INEX
CO(content only) CAS(content and structure)
CO CAS
. PIX
CO
가 Q1 ~ Q4 3,4
Q5 ~ Q7 가



9

1

9 Q1 ~ Q4 INEX

4 가

가

12 XML

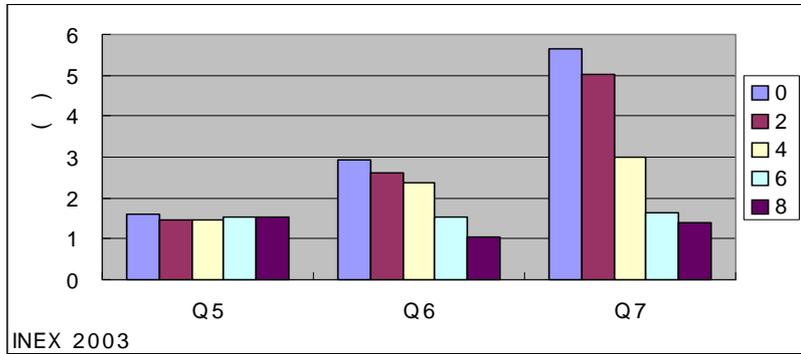
4)

가

가 가

4)

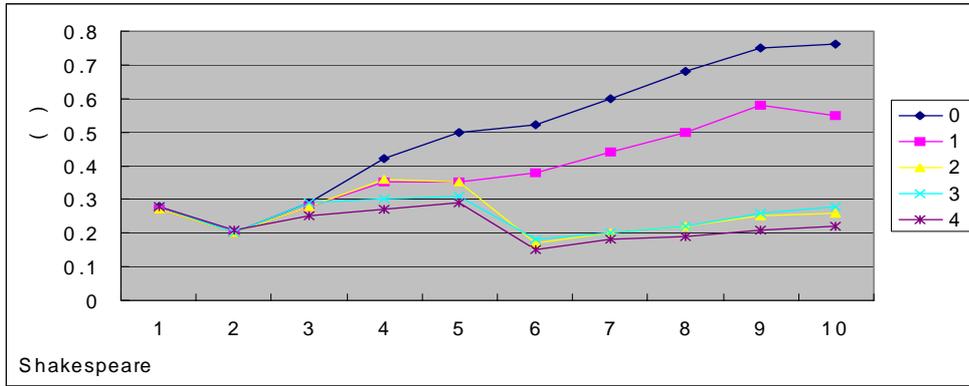
가 가



10

2

INEX 10 가 가 Q5 ~ Q7
 가 9
 Q5 가
 4 가 15 , 4
 가 16 가 Q7
 가 4 16 4
 PIX 가 PF
 가 PIX PIXPF P-value
 , INEX DTD 가 가
 가 가 가



11

11 Shakespeare

가

가

가

가 가

가

가 가

가

가

가
가

가

INEX

2

12

2

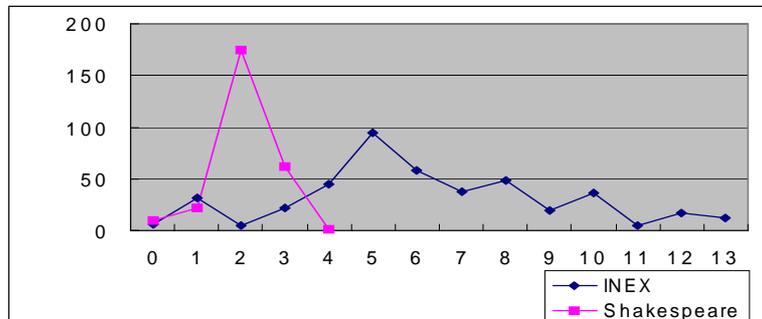
가

가

가 가

INEX

가



12

6.

가 XML 가 가
가 . PIX XML 가
가 , .
가 가 Top-k
가 가
가 XML 가 가
가 XML 가

[1] <http://www.w3.org/XML/>

[2] L. Mignet, D. Barbosa, P. Veltri. The XML Web: a First Study. WWW 2003

D Florescu, et al. Integrating Keyword Search into XML Query Processing.

WWW '99

- [3] S. Putz. Using a Relational Database for an Inverted Text Index. XEROX Technical Report '91
- [4] V. N. Anh, O. Krester, A. Moffat. Vector-Space Ranking with Effective Early Termination. SIGIR '01
- [5] D. Cutting, J. Pedersen. Optimizations for Dynamic Inverted Index Maintenance. SIGIR Conf. on Research & Development in IR '90
- [6] N. Fuhr, K. Grojohann. XIRQL: A Query Language for Information Retrieval in XML Documents. SIGIR '01
- [7] A. Theobald, G. Weikum. The Index-based XXL Search Engine for Querying XML Data with Relevance Ranking. EDBT '02
- [8] A. Theobald, G. Weikum. Adding Relevance to XML. WebDB '00
- [9] D. Florescu, et al. Integrating Keyword Search into XML Query Processing. WWW '99
- [10] V. Hritidis, Y. Papakonstantinou, A. Balmin. Keyword Proximity Search on XML Graph. ICDE 2003
- [11] L. Guo, et al. XRANK: Ranked Keyword Search over XML Documents. SIGMOD '03
- [12] S. Cohen, J. Mamou, Y. Kanza, Y. Sagiv. XSEarch: A Semantic Search Engine for XML. VLDB 2003
- [13] D. Carmel, Y. S. Maarek, M. Mandelbrod, Y. Mass, A. Soffer. Searching XML Documents via XML Fragments. SIGIR 2003
- [14] S. Brin, L. Page, The Anatomy of a Large-Scale Hypertextual Web Search Engine. WWW7 '98
- [15] G. Salton and M.J. McGill, "Introduction to Modern Information Retrieval". McGraw-Hill, New York, 1983.
- [16] <http://www.sleepycat.com>
- [17] Initiative for the evaluation of XML retrieval
- [18] A. Moffat, J. Zobel. Self-Indexing Inverted Files for Fast Text Retrieval. TODS Vol. 14, No. 4, 1996.

- [19] J. Robie, J. Lapp, and D.S Schach. XML query language(XQL). The Query Languages Workshop. W3c, Dec. 1998.
<http://www.w3.org/TrandS/QL/QL98/pp/xql.html>
- [20] A. Deutsch, M. Fernandez, D. Florescu, A. Levy, and D. Suciu. XML-QL: A query language for XML. The Query Languages Workshop. W3c, Dec. 1998. <http://www.w3.org/TR/1998/NOTE-xml-ql-19980819/>.
- [21] XQuery: A query language for XML, Feb. 2001.
<http://www.w3.org/XML/Query>
- [22] XPath: XML Path language, Nov. 1999. <http://www.w3.org/TR/xpath>
- [23] J.P. Callan. Passage-Level Evidence in Document Retrieval. SIGIR 1994.
- [24] R. Wilkinson. Effective retrieval of structured documents. SIGIR 1994.
- [25] J. Zobel, A. Moffat, R. Wilkinson, and R. Sacks-Davis. Effientent retireval of partial documents. Information Processing and Management, 31(3):361-377, 1995
- [26] <http://www.ibiblio.org/xml/examples/shakespeare/>