

2004

XML

-

2 0 0 4

XML

-

論文 碩士學位 論文 提出

2004 年 6 月

梨花女子大學校 科學技術大學院

學科 崔 恩 慧

# 碩士學位論文 認准

指導教授 \_\_\_\_\_

審查委員 \_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

	-----	v
I .	-----	1
1.1	-----	1
1.2	-----	2
II .	-----	5
2.1	-----	5
2.2	-----	5
2.2.1 DataGuides	-----	6
2.2.2 1-index	-----	7
2.2.3 Index Fabric	-----	7
2.3	-----	8
2.3.1 A(k)-index	-----	9
2.3.2 D(k)-index	-----	9
2.3.3 B+ tree	-----	11
III .	-----	12
3.1 -	-----	12
3.1.1 XML	-----	13
3.1.2 -	-----	15
3.1.3 -	-----	16
3.1.4 -	-----	18
3.2	-----	19
3.2.1 Value B+ tree	-----	20
3.2.2 Name B+ tree	-----	21
IV .	-----	22

4.1		-----	22
4.2		-----	25
4.3	-	-----	26
4.3.1	1	-----	27
4.3.2	2	-----	28
4.3.3	3	-----	29
V.		-----	31
5.1		-----	31
5.2		-----	32
5.3	가	-----	34
VI.		-----	37
		-----	39
		-----	42

[ 3-1]	XML	-----	14
[ 3-2]	XML	-----	15
[ 3-3]	-	-----	17
[ 3-4]		-----	18
[ 3-5]		-----	18
[ 3-6]	-	-----	19
[ 3-7]	Value B+ tree	-----	20
[ 3-8]	Name B+ tree	-----	20
[ 4-1]		-----	24
[ 4-2]		-----	26
[ 5-1]		-----	34
[ 5-2(a)]	(Book )	-----	35
[ 5-2(b)]	(Movie )	-----	35

[ 5-1]		-----	31
[ 5-2]	XML	-----	32
[ 5-3]		-----	32

# I.

## 1.1

XML[1]

가 XML

XML

XML

XPath[2], XQuery[3]

XML

XML

XML

가

[9,15].

XML

[11,12,13,14].

XML

[5,6,7]

XML



/ /

,

가

[4,5,9]. XML

[16]

XML

XML

, XML 가

, 가

XML

가 .

가

## 1.2

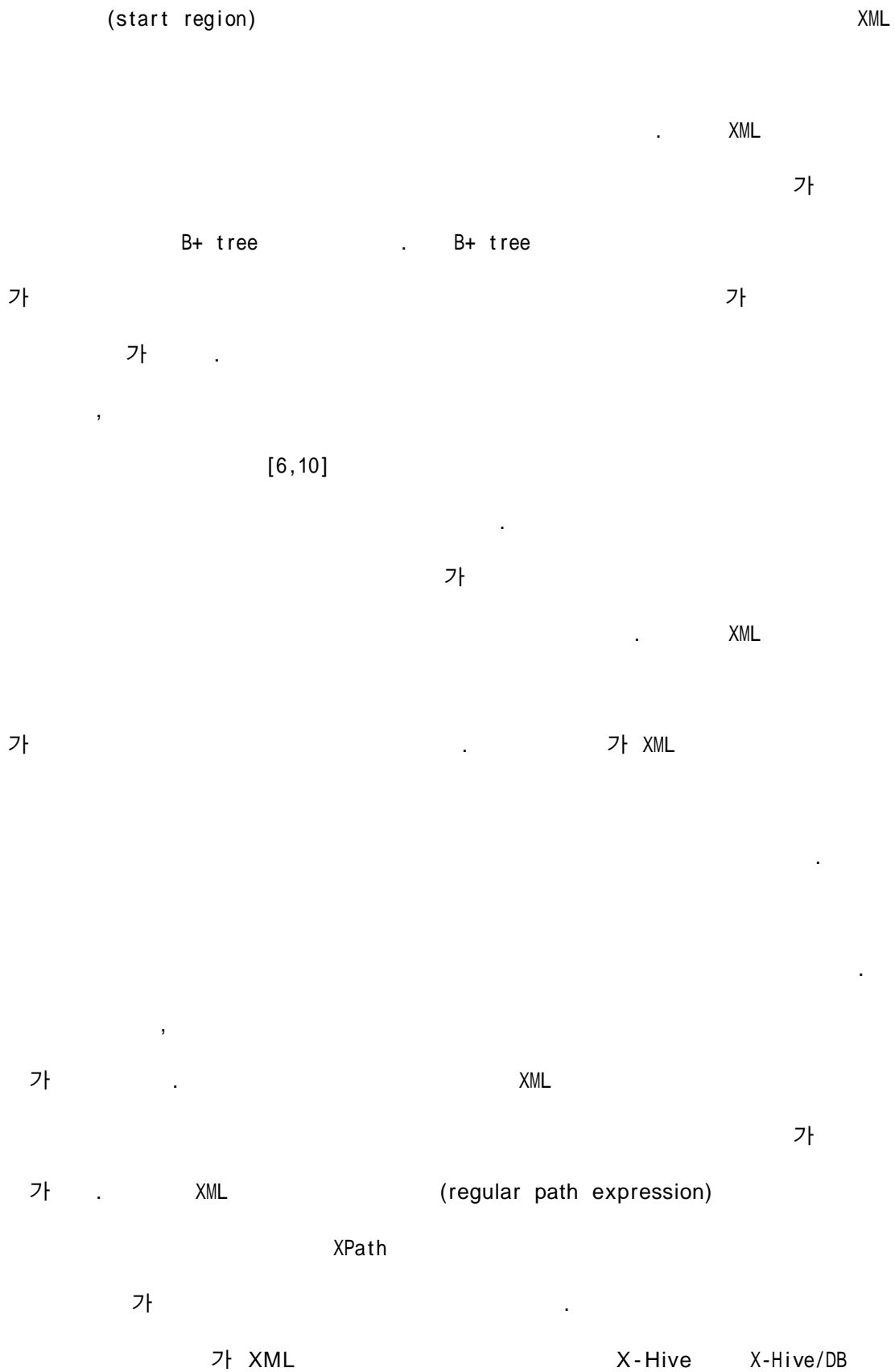
가

XML

XML

(Depth First Search)

XML



1.5

. X-Hive/DB

가

,

가

가

.

.

. 2

XML

, 3

. 4

3

, 5

가

XML

X-Hive/DB

,

6

.

▪  
, XML

## 2.1

XML

/

가

가

XML

가

B+ tree

## 2.2

(path expression)

XML

(Structural

summary)[17,18,19]

[14,17].

XML

DataGuides 1-index 가 ,  
Index Fabric

### 2.2.1 DataGuides

DataGuides[13] Object Exchange Model(OEM)  
(semistructured data)

. DataGuides OEM  
. DataGuides  
가

(regular path expressions)

, OEM 가

OS

가

가

(full path indexing)

가 가

XML

가



가

가

가

XML

- , Index Fabric

, ('//')

가

## 2.2

1-index DataGuides

XML

(regular path expression)

A(k)-index D(k)-index가

B+ tree

XML

(depth-first traversal)

( , )

### 2.2.1 A(k)-index

DataGuides 1-index ,  
 가 가 , 가 XML  
 , 가 . 가  
 K 가  
 A(k)-index[19]가 .  
 K 가 가  
 . A(k)-index K . K

### 2.2.2 D(k)-index

D(k)-index[11] 1-index A(k)-index  
 .  
 가 , 1-index  
 ,  
 .  
 (edge) , 가 1-  
 index, A(k)-index 가 .  
 D(k)-index 가 ,  
 (local similarity) ,  
 (bisimilarity requirement)



가 D(k)-index

가 가 가

가

가

D(k)-index XML

### 2.2.3B+ tree

[5,9] XML

A//B

A/B

가

가 .

[6.7] 가

가) [6]

( ,

)

B+ tree

B+ tree

(containment)

(sibling

pointers) 가

, "Movies//Actor"

Movies Actor

가 .

B+ tree

,

,  
가  
가 ,  
가 .

Movies Actor 가

Movies  
가 .

) XR-Tree [7]가 .

XML

B+ tree , [6]

, (key)

(stab list)

2

XR-tree

, B+ tree

▪

XML

XML

/ , /

XML

XML

, /

가 XML

가

가

XML

/

XML

3.1

XML

### 3.1.1 XML

XML

/ /

. XML

XML

3-1

XML

DTD

. DTD

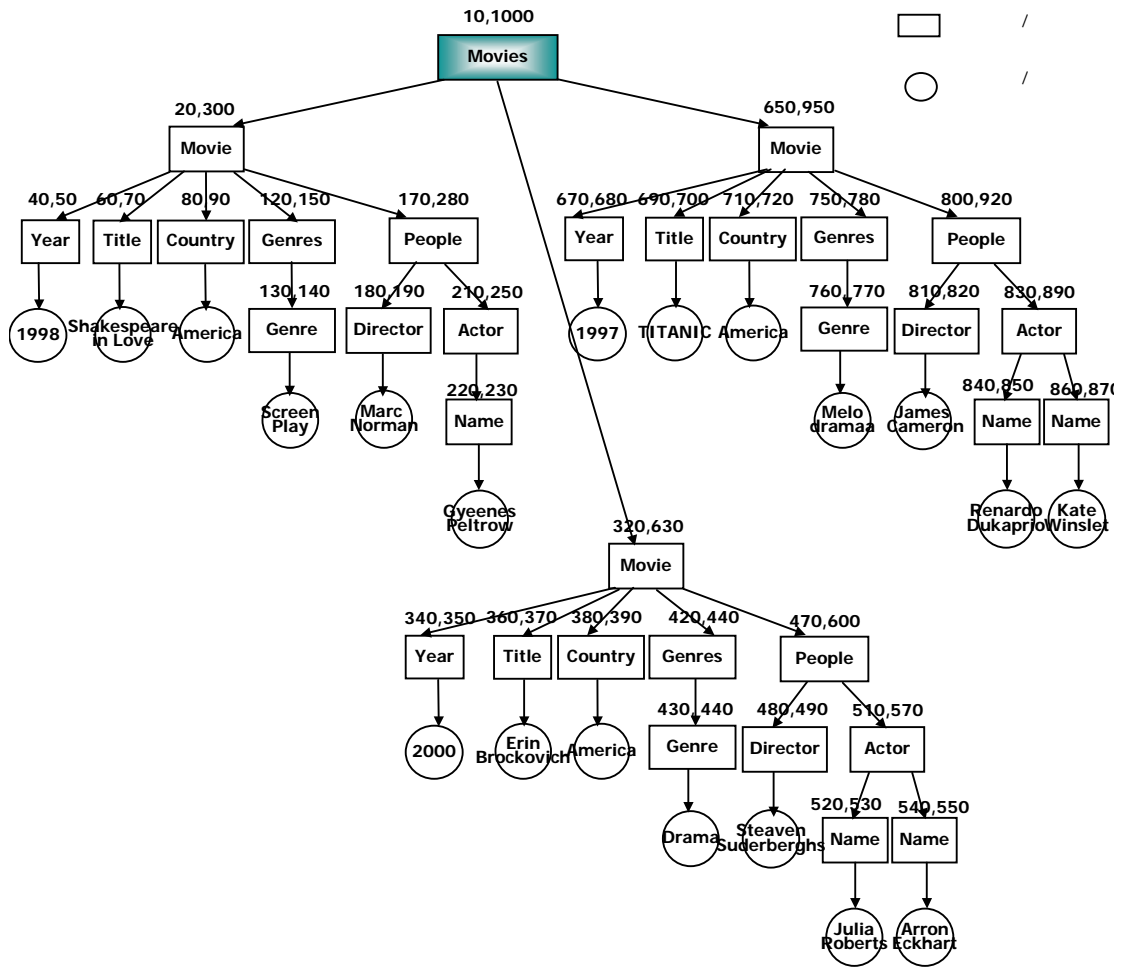
#PCDATA

가

. XML

3-2

<pre> &lt;?xml version ="1.0"?&gt; &lt;!DOCTYPE Movies [ &lt;!ELEMENT Movies (Movie*)&gt; &lt;!ELEMENT Movie ( Title,Country,Genres,People)&gt; &lt;!ATTLIST Movie Year CDATA #REQUIRED&gt; &lt;!ELEMENT Title (#PCDATA)&gt; &lt;!ELEMENT Country (#PCDATA)&gt; &lt;!ELEMENT Genres (Genre*)&gt; &lt;!ELEMENT Genre (#PCDATA)&gt; &lt;!ELEMENT People (Director,Actor)&gt; &lt;!ELEMENT Director (#PCDATA)&gt; &lt;!ELEMENT Actor (Name*)&gt; &lt;!ELEMENT Name (#PCDATA)&gt; ]&gt;  &lt;Movies&gt;   &lt;Movie Year="1998"&gt;     &lt;Title&gt;Shakespeare in Love&lt;/Title&gt;     &lt;Country&gt;America&lt;/Country&gt;     &lt;Genres&gt;       &lt;Genre&gt;ScreenPlay&lt;/Genre&gt;     &lt;/Genres&gt;     &lt;People&gt;       &lt;Director&gt;Mark Norman&lt;/Director&gt;       &lt;Actor&gt;         &lt;Name&gt;Gyeeenes Peltrow&lt;/Name&gt;       &lt;/Actor&gt;     &lt;/People&gt;   &lt;/Movie&gt; </pre>	<pre> &lt;Movie Year="2000" &gt;   &lt;Title&gt;Erin Brockovich&lt;/Title&gt;   &lt;Country&gt;America&lt;/Country&gt;   &lt;Genres&gt;     &lt;Genre&gt;Drama&lt;/Genre&gt;   &lt;/Genres&gt;   &lt;People&gt;     &lt;Director&gt;Steaven Soderberghs&lt;/Director&gt;     &lt;Actor&gt;       &lt;Name&gt;Julia Roberts&lt;/Name&gt;       &lt;Name&gt;Aarron Eckhart&lt;/Name&gt;     &lt;/Actor&gt;   &lt;/People&gt; &lt;/Movie&gt; &lt;Movie Year="2000" &gt;   &lt;Title&gt;Erin Brockovich&lt;/Title&gt;   &lt;Country&gt;America&lt;/Country&gt;   &lt;Genres&gt;     &lt;Genre&gt;Drama&lt;/Genre&gt;   &lt;/Genres&gt;   &lt;People&gt;     &lt;Director&gt;Steaven Soderberghs&lt;/Director&gt;     &lt;Actor&gt;       &lt;Name&gt;Julia Roberts&lt;/Name&gt;       &lt;Name&gt;Aarron Eckhart&lt;/Name&gt;     &lt;/Actor&gt;   &lt;/People&gt; &lt;/Movie&gt; &lt;/Movies&gt; </pre>
---	--



3-2 XML

3.1.2

가

3-2

, Movies

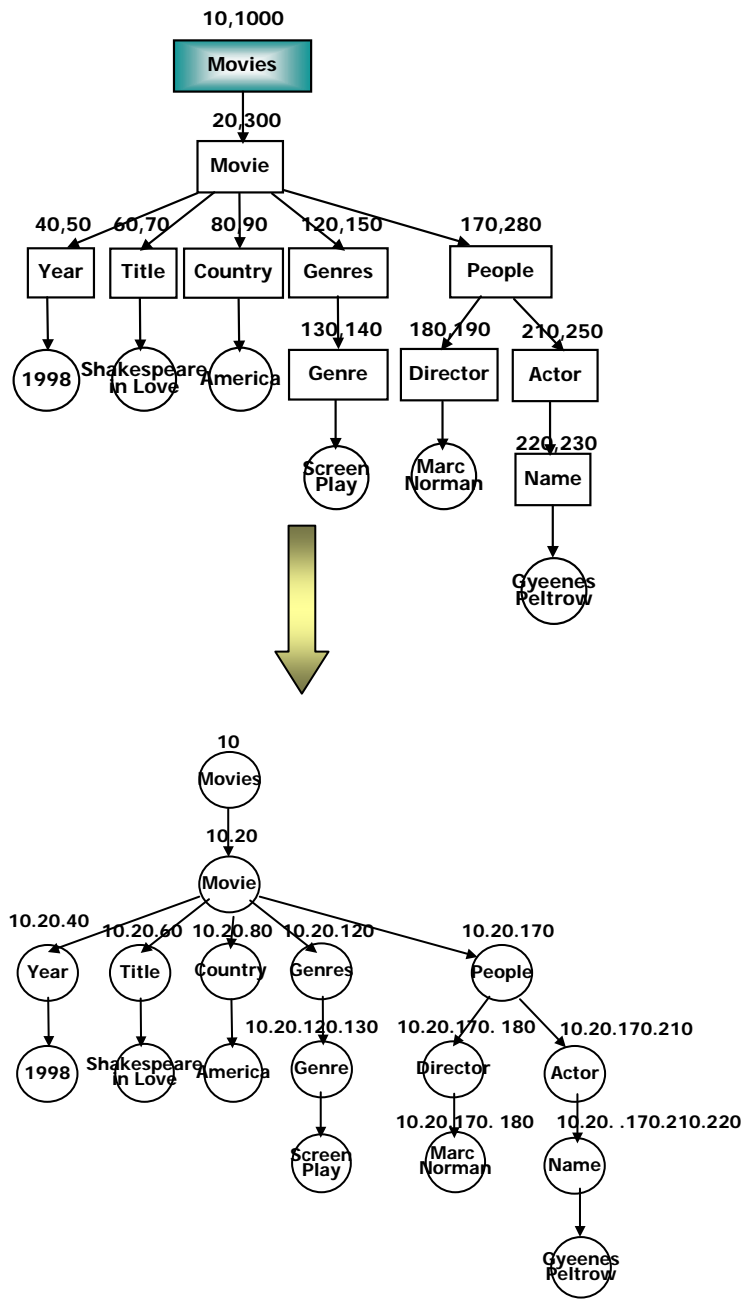
10 1000

Movie 10 20 가 300 .  
 Movie 21 299 .  
 Movie (40,50) 가  
 Year Time 가 (30,35)  
 가 .  
 XML  
 , 3-2 가 1998  
 Year 가 (40,50) 가 .

3.1.3 -

XML -  
 [6,7] ,  
 3-3  
 XML .  
 가 .  
 3-4 Year - 가 10.20.40 ,  
 10 Movies가 Year 10.20 Movie가  
 .  
 , 10.20

Title, Country, Genres, People Year .



3-3

XML

3-3

XML

3-4, 3-5

가 ,



Doc_ID	Start_Region	Node_Name	Path_Info
1	10	MovieList	10
1	20	Movie	10.20
1	30	ID	10.20.30
1	40	Year	10.20.40
1	50	Title	10.20.50
1	60	Country	10.20.60
1	70	Staff	10.20.70
1	80	Scenarist	10.20.70.80
1	90	Directed_By	10.20.70.90
1	100	Director	10.20.70.90.100
1	110	Link	10.20.70.90.110
1	120	Create Year	10.20.70.90.110.120
...	...	...	...

3-4

Doc_ID	Start_Region	Text_Value	Path_Info
1	30	Sw	10.20.30
1	40	1977	10.20.40
1	50	Star Wars	10.20.50
1	60	America	10.20.60
1	80	Jim Andrew	10.20.70.80
1	100	George Lucas	10.20.70.90.100
1	120	1980	10.20.70.90.110.120
1	140	Peter Morison	10.20.70.140
1	170	Adventure	10.20.140.170
1	180	Action	10.20.140.180
1	190	Sci-Fi	10.20.140.190
...	...	...	...

3-5

3.1.4

가

3-2

Movie

(10.20), (10.320), (10.650)

가

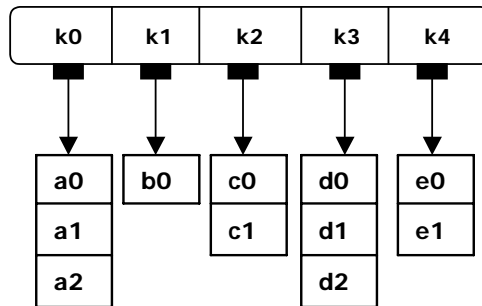
가

3-6 XML

가 5

k 5 가

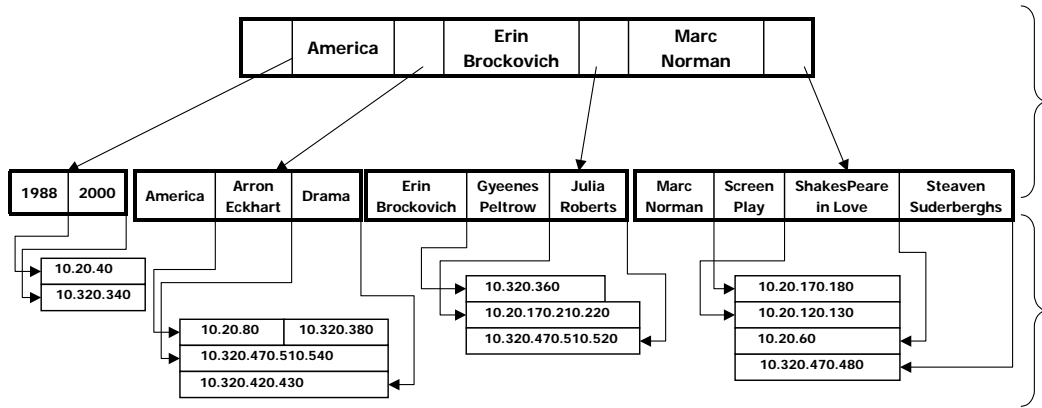
/ 가  
 - , k<sub>0</sub> XML  
 - a<sub>0</sub> , a<sub>1</sub> , a<sub>2</sub> . - XML  
 /



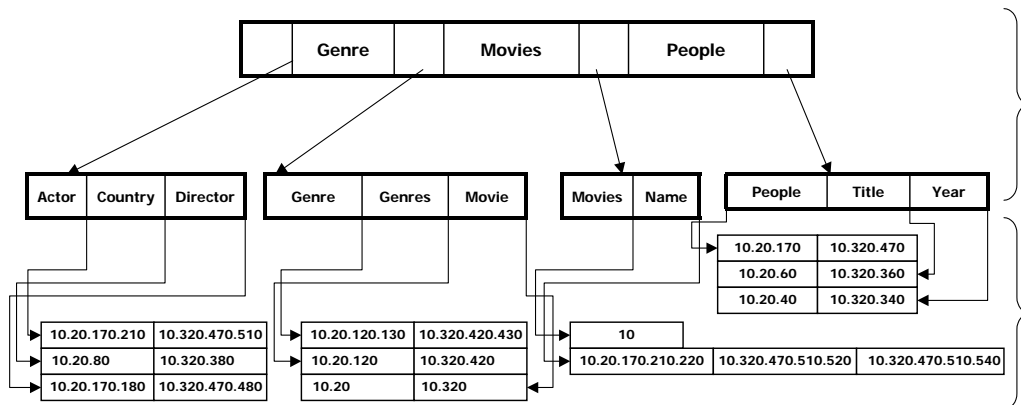
3-6 -

### 3.2

B+ tree . B+ tree  
 가  
 가 XML  
 B+ tree . B+ tree  
 XML /  
 B+ tree , -  
 /  
 - Value B+ tree , /  
 - Name B+ tree .



3-7 Value B+ tree



3-8 Name B+ tree

3.2.1 Value B+ tree

3-7 Value B+ tree

(key entry)

-6

가 , Value B+ tree

가

XML

[Year="1998"] 가 XPath ,  
 3-7 1998 Value B+ tree  
 - 1998 (10.20.40)  
 . [Year="1998" AND Country="America"]  
 가 가

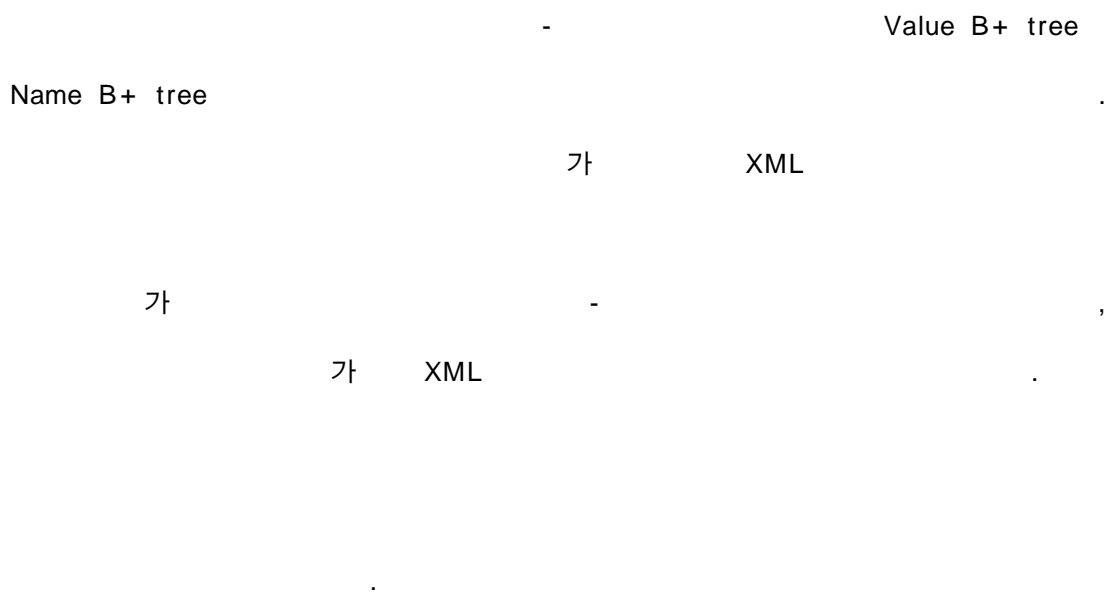
### 3.2.2 Name B+ tree

3-8 Name B+ tree

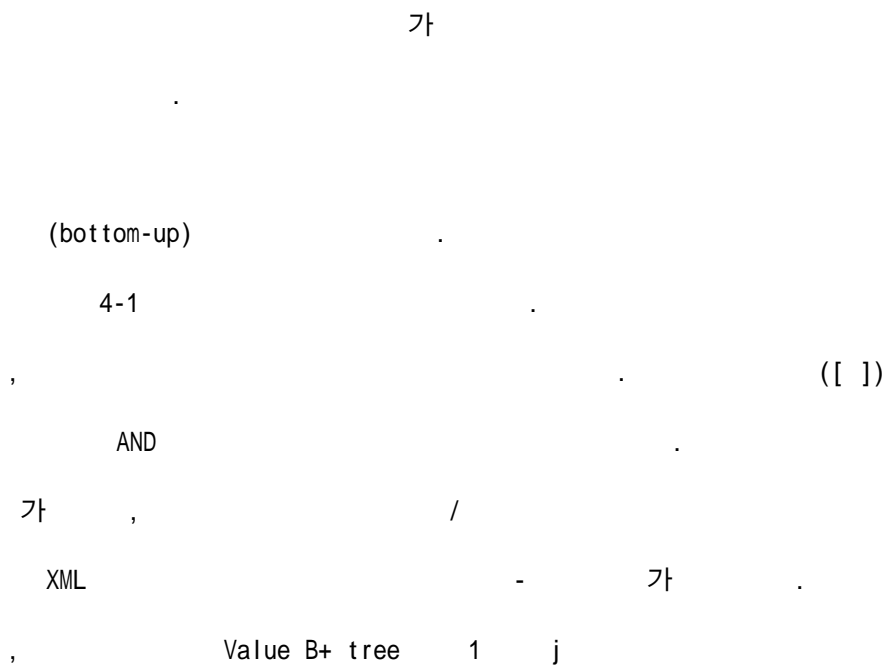
, XML -  
 가 . , Name B+ tree -

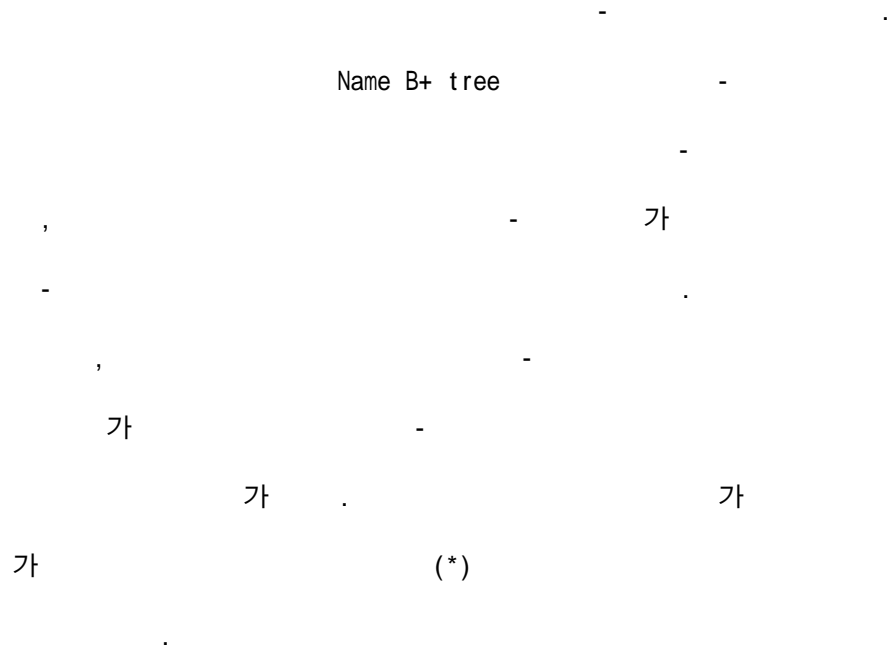
Movies//Movie[Title="Erin Brockovich"] XPath  
 Movies Movie Title , Name B+ tree  
 , -  
 . , Movies -  
 (10) , Movie (10.20),(10.320),(10.650) Movies .  
 Title - 가 (10.20).(10.320) (10.650) Movie  
 가  
 (10.20.60),(10.320.360),(10.650,690). Name B+ tree XML

# IV.



## 4.1





---

**1**


---

**Definition:**

$-pn_1, pn_2, \dots, pn_j$ ;  
 $-pv_1, pv_2, \dots, pv_j$ ;  
 $-qn_1, qn_2, \dots, qn_x$ ;  
 $-qv_1, qv_2, \dots, qv_y$ ;  
 $-conA_1, conA_2, \dots, conA_n$ ;

**Description:**

가

**Input:**  $Q = \text{[/*[}pn_1 = "pv_1" \text{ AND } pn_2 = "pv_2" \text{ ... AND } pn_j = "pv_j" \text{ ]]}$

Value B+-tree,

Name B+-tree

**Output:**

**Function** *SearchCommonAncestor()* {

**while**(1    j    ) **do**

$qv_y \leftarrow$  Value B+-tree     $pv_j$     -    .;

$qn_x \leftarrow$  Name B+-tree     $pn_j$     -    .;

**end**

  /\* ( $qn_x, qv_y$ )    -    가    \*/

**while** (1    x    ) **do**

**for** (1    y    ) **do**

$comA_n \leftarrow$   $qn_x$      $qv_y$         $qn_x$     .

      ,     $qn_x$     .;

**end**

**end**

**while**(1    n    ) **do**    /\* $conA_n$     self join    \*/

**for** (1    n    ) **do**

$comPrefix \leftarrow$     -

**end**

**end**

**return**  $comPrefix$  ; //

}

---

## 4.2

가

-

.

4-2

.

( ' / ' )

( ' / / ' )

2

.

Anc

Des

,

Name B+ tree

Anc

Des

가

-

-

.

Anc

가

Des

,

.

Anc

Des

-

Anc

-

.

, 4.1

Anc

-

Des

.

Des

, Anc

-

Dec

-

, 4.1

Dec

-

.

Dec

-

.



---

2

---

**Definition:**

-Anc :  
 -Des :  
 - $ap_1, ap_2, \dots, ap_j$  : -  
 - $dp_1, dp_2, \dots, dp_j$  : -  
 - $cap_1, cap_2, \dots, cap_n$  :

**Description:**

가

**Input:**  $Q = //Anc[pn_1="pv_1" \text{ AND } pn_2="pv_2" \dots \text{ AND } pn_j="pv_j"]//des$   
 Name B+-tree,  
 Value B+-tree

**Output:**

**Function SearchDescendant() {**

```

/*
    Name B+-tree      anc      -      .      ;
    Value B+-tree     des      -      .;
*/
cp = SearchCommonAncestor();

while (1 < x ) do
    cap_n ← ap_x cp      ap_x      .
end

/*
return cap_n      dp_j;
*/
}
    
```

---

4-2

4.3

-

XML

XPath

([])

. XML

' / ' ' // '

, - 가 ,

- .

‘ \* ’ .

3-1 XML XPath

### 4.3.1 1

Country가 ‘ America ’ 가 Title

XPath .

**Query = `//*[Country="America"]//Title`**

`[Country="America"]`

) `*/Title` - .

가. ‘ America ’

. , Value B+ tree ‘ America ’

- America : (10.20.80),(10.320.380)

‘ America ’ ‘ Country ’ . , Name

B+ tree ‘ Country ’ - .

- Country : (10.20.80),(10.320.380),(10.690.730)

, Country America가 - America

- Country : (10.20.80), (10.320.380)

‘ Title ’ Name B+ tree

‘ Title ’

- Title : (10.20.60), (10.320.360)

, 가) Country - Title -

- Country : (10.20.80), (10.320.380)

- Title : (10.20.60), (10.320.360)

10.20	10.320	가	Country	Title
		가		
(10.20.60)	(10.320.360)	가		
‘Shakespeare in Love	‘Erin Brockovich’	가		Title

4.3.2 2

1) 가

가 Movie

**Query = //Movie[Country="America"]/Title**

가. 1) 가) ‘ America ’ 가

Country -

- Country : (10.20.80), (10.320.380)

. Movie Name B+ tree ‘ Movie ’

- Movie : (10.20),(10.320)

, 가) Country - Movie -  
가 - . Country -  
가 Movie -

. Name B+ tree 'Title'

- Title : (10.20.60), (10.320.360)

, Movie Title )  
Movie - Title .

- Movie : (10.20),(10.320)

- Title : (10.20.60), (10.320.360)

. (10.20.60) (10.320.360) 'Shakespeare in Love' 'Erin  
Brockovich'가 Title .

### 4.3.3 3

1,2 2 가

가 XPath

1-2 XML , Genres Genre 가

. Genres 가 .

Country가 'America' Genres가 'Drama' Title

XPath .

Query = //Movie[Country="America" AND Genres="Drama"]//Title

가. 2 가

America Drama . , Value B+ tree America Drama

- America : (10.20.80),(10.320.380)

- Drama : (10.320.420.430)

. , Country Genres . ,

Name B+ tree Country Genres - .

- Country : (10.20.80),(10.320.380)

- Genres : (10.320.420)

Country Genres -

가 - (10.320)

가

. Movie

Title Movie//Title

Name B+ tree Movie Title - .

-Movie : (10.20),(10.320)

-Title : (10.20.60),(10.320.360)

(10.320) Movie(10.320)

, Title . (10.320.43)

. (10.20.60) (10.320.360) 'Shakespeare in Love' 'Erin

Brockovich' .

3 XML 가

XML X-Hive X-Hive/DB

5.1

X-Hive/DB XML (Native XML Database) XML  
 ( ) XML  
 DOM XML

	Windows 2000 Advanced Server
<b>CPU/RAM</b>	Pentium IV processor 2.40MHz/512M RAM
	JAVA(JDK 1.4.2)
<b>DBMS</b>	Microsoft SQL Server
	Borland JBuilder 9.0

5-1

5-1

Intel Pentium 4 CPU 2.40GHz Windows

2000 Advanced Server

512MB . RDBMS MS SQL Server MS SQL Server  
 XML , JBuilder 9.0 JAVA

5.1.1

5-2 XML  
 가 XML  
 Movie , XML Book

	Movie	Book
	100	19
	155KB	570KB
	4776	5445
	5	9
	2908	3630
	477600	103455

5-2 XML

5.1.2

Movie	Q1	<code>//*[Country="America"]//Title</code>
	Q2	<code>//Movie[Country="America"]//Title</code>
	Q3	<code>/MovieList/Movie[Country="America"]//Title</code>
	Q4	<code>//*[*="America"]//Title</code>
	Q5	<code>//*[Country="America" AND Year="1998"]//Title</code>
	Q6	<code>//Movie[Country="America" AND Year="1998"]//Title</code>





가

### 5.1.3 가

2

XML

6

5-1

XML

5-3

X-Hive

Book

Movie

6

5-2

msec

5-2

2

XML

가

가

X-Hive

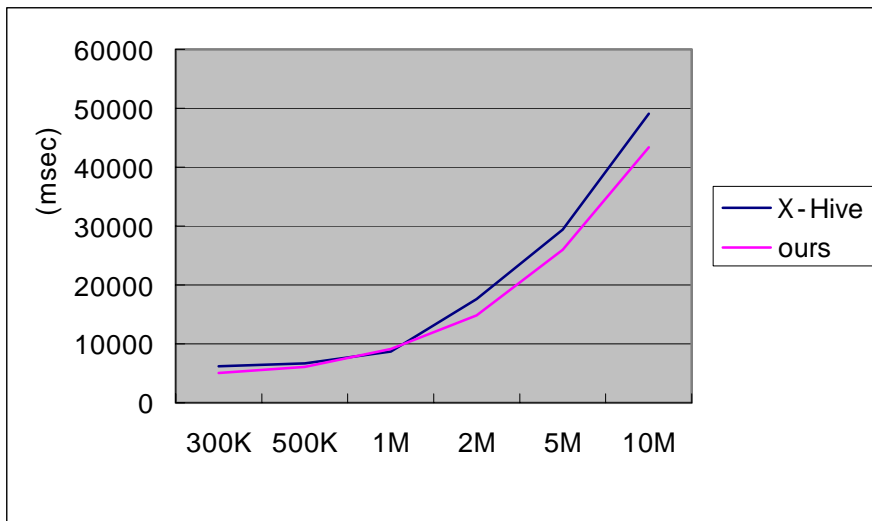
1.5

X-Hive

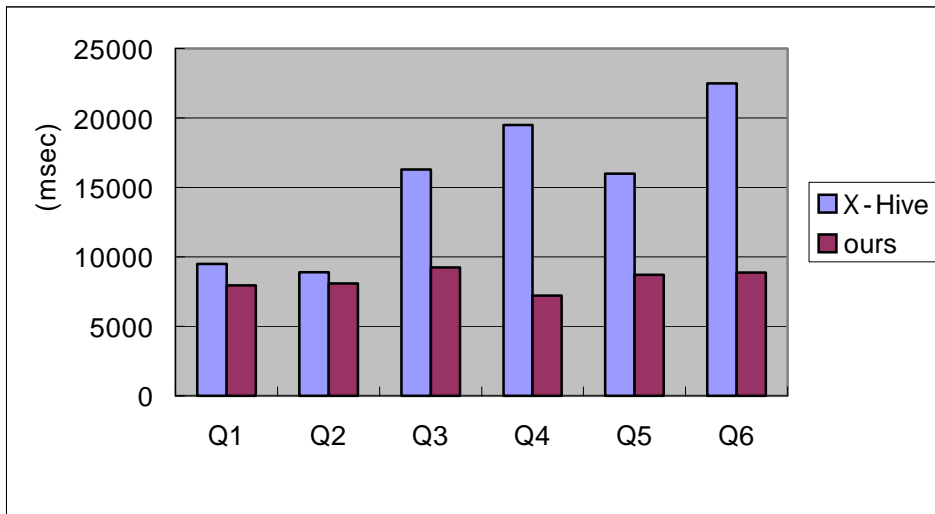
가

가

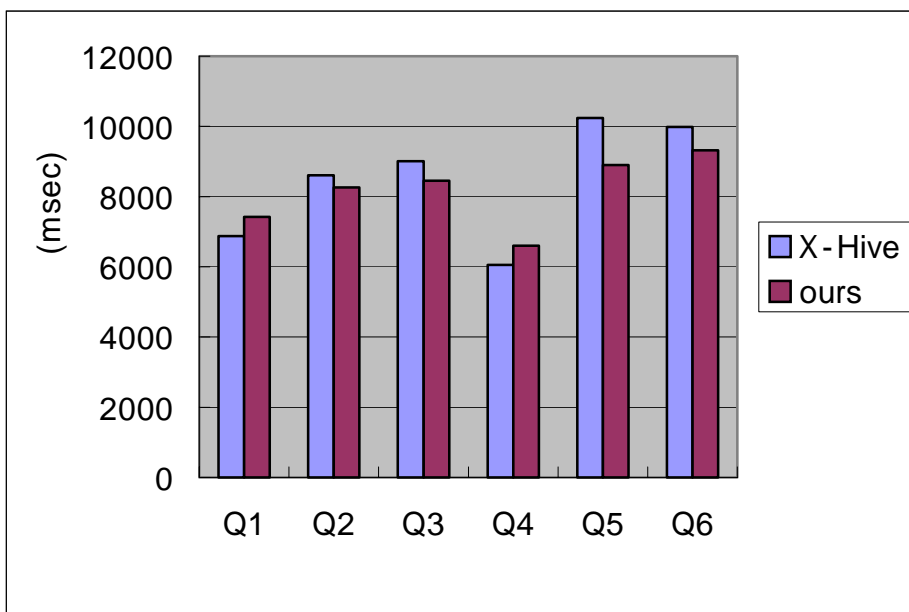
가



5-1



5-2 (a) (Book )



5-2 (b) (Movie )

Q1, Q2, Q3

Q1

가

Q2, Q3

가

.  
가 . Q4  
. , Q1 가  
. 2 가 Q5,Q6 1 Q1,Q2  
가 . 가 가  
. .

# VI.

XML

가

가

.

가

XML

.

가

.

,

XML

‘ - ’

XML

.

가

XML

.

,

가

.

.

,

-

XML

,

XML

. XML

가 XML

가 XML  
XPath XML  
가 XML  
가 XML

- [1] W3C Consortium, XML 1.0 (Second Edition), W3C Recommendation 06 Oct. 2000, available at <http://www.w3.org/TR/REC-xml>.
  
- [2] W3C Consortium, XML Path Language (XPath), version 2.0, W3C Recommendation Nov 12, 2003. <http://www.w3.org/TR/xpath20.html>.
  
- [3] S.Boag, D. Chamberlin, M. F. Fernandez, D. Florescu, J. Robie, and J. Simeon, XQuery 1.0: An XML Query Language, W3C Working Draft Nov 12, 2003. <http://www.w3.org/TR/xquery/>.
  
- [4] D. Srivastava, S. Al-Khalifa, H. V. Jagadish, N. Koudas, J. M. Patel, and Y. Wu. Structural joins: A primitive for efficient XML query pattern matching. Proc. Of Int. Conf. On Data Engineering, pages 141 - 151, 2002
  
- [5] C. Zhang, J. F. Naughton, D. J. DeWitt, Q. Luo, and G. M. Lohman. On supporting containment queries in relational database management systems. Proc. of the ACM SIGMOD Int. Conf. on Management of Data, pages 425 - 436, 2001.
  
- [6] S.-Y. Chien, Z. Vagena, D. Zhang, V. Tsotras, and C. Zaniolo. Efficient Structural Joins on Indexed XML documents. Proc. Of the VLDB, pages 263 - 274, 2002.
  
- [7] H. Jiang, H. Lu, and W. Wang. XR-Tree: Indexing XML Data for Efficient

- Structural Joins. Proc. Of International Conference On Data Engineering, pages 253-264, 2003.
- [8] B. F. Cooper, N. Sample, M. J. Franklin, G. R. Hjaltason, and M. Shadmon. A fast index for semistructured data. Proc. Of the VLDB, pages 341-350, 2001
- [9] Q. Li and B. Moon. Indexing and Querying XML Data for Regular Path Expressions, Proc. Of the VLDB, pages 361-370, 2001.
- [10] A. M. Flavio Rizzolo. Indexing XML Data with Toxin. In WebDB, pages 49--54, 2001.
- [11] Q. Chen, A. Lim, and K. W. Ong, D(k)-Index: An adaptive structural summary for graph-structured data, Proc. of the ACM SIGMOD Int. Conf. on Management of Data, 2003
- [12] H. Jiang, H. Lu, and W. Wang. XR-Tree: Indexing XML Data for Efficient Structural Joins. Proc. Of International Conference On Data Engineering, pages 253-264, 2003.
- [13] Roy Goldman , Jennifer Widom, DataGuides: Enabling Query Formulation and Optimization in Semistructured Databases, Proc. Of the VLDB, pages 436-445, August 25-29, 1997.
- [14] T. Milo and D. Suciu, Index Structures for path expressions, Proc. of the IEEE Int.

Conf. On Data Theory, pages 277-295, 1999

- [15] Albrecht Schmidt, Martin Kersten, Menzo Windhouwe. Querying XML Documents made Easy : Nearest Concept Queries, Proc.Of International Conference on Data Engineering, 2001
  
- [16] S.W.Kim, et al. Indexing and Retrieval of XML-encoded Structured Documents in Dynamic Environment. Lecture Notes in Computer Science (LNCS) Vol.2480, 2002
  
- [17] Peter Buneman, Susan Davidson, Gerd Hillebrand, and Dan Suciu. A query language and optimization techniques for unstructured data. Proc of the ACM SIGMOD International Conference on the Management of Data, pages505-516.
  
- [18] Svetlozar Nestorov, Jeffrey Ullman, Janet Wiener, and Sudarshan Chawathe. Representative objects: concise representations of semistructured, hierarchical data. Proc. of the IEEE International Conference on Data Engineering, 1997:79-90.
  
- [19] R. Kaushik et al., "Exploiting Local Similarity for Indexing Paths in Graph-Structured Data", Proc.Of International Conference On Data Engineering, 2002.





## ABSTRACT

# A Region-Path Indexing Technique for Content Based Retrieval of XML documents

*Department of Computer Science & Engineering*

*Ewha Institute of Science and Technology*

*Choi Eun Hye*

XML is a markup language for documents containing contents and structured information. Especially data on the internet are usually represented and transferred as XML. The XML data is represented as a tree and therefore, indexing techniques are needed to efficiently support hierarchical structural properties. XML queries are represented as regular path expressions and evaluated by traversing each object of the tree. Several indexes are proposed to fast evaluate regular path expressions.

However, in some cases they may not cover all possible paths because they require a great amount of disk space. In order to efficiently evaluate the queries in such cases, we propose an optimized traversing which uses the region-path index. Minimizing traverse of not only data object but also index object, we can execute the XPath expressions fast. That is, we propose a region-path index scheme as a new index scheme to efficiently process queries with extended path expressions. Our proposed index scheme allocates a unique path identifier for every possible single

path in as extended path expression and provides functionalities of both single path indexing and multiple path indexing through the composition of index key and path identifier while using only a index structure. The proposed index scheme provides better performance than single-path index schemes, and is practical since it can be implemented by little modification of leaf record of a B+ tree index.

We conducted diverse experiments to show that our indexing techniques achieves a better performance than the other approach.