

# XML

XML

가

가

가

DTD

XML

XML

## Abstract

XML is becoming the *de facto* standard for data exchange over the Internet as a semistructured data which properties are irregular and incomplete. Therefore, to handle these data efficiently, what we use storage devices and storage techniques are main factors. In this paper, we developed storage techniques, which take the virtues of an object-relational database and support various query types needed for XML query languages without regard to the DTD. The techniques are capable of connecting naturally with conventional data and reducing overheads caused by the characteristics of an XML data model.

Keyword: storage technique, XML, Schema mapping

## 1. Introduction

XML[1]

가

가

가

XML

, XML

XML

XML 가 [2,3,4].

XML

XML

, 가 , ,

[5,6,7]

가[8-15]

XML

XML

가

가

, XML 가

XML

[16-21]

XML

가

가

XML

1 XML

2

, 가 XML DTD가  
 , DTD가 XML DTD 가  
 , XML XML  
 - , XML  
 , - , - ,  
 , XML - -  
 (mapping rule)  
 XML DTD가  
 - XML  
 가 XML ,  
 XML 가  
 2 , 3 XML  
 4 - XML 5  
 가 6 .

## 2. Related Work

가 XML

[5-9,13,14,22,23].

STORED[8]

(semistructured)

, STORED

/ /

, XML 가

가

[9]

가

XML

8가

XML

가

XML

[13] XML

XML

DTD

XML

XML

가

XML

가

DTD가

[8,9,13,23]

XML

XML

[14] 가 XML  
가  
DTD가  
XML DBMS[26] XML  
XML  
2  
[27] DTD 가 가 , DTD [13]  
XML XADT  
ORDBMS  
Oracle XML DB[5] SQL 2000 [6] XML  
가 XML 가  
, SQL XML  
가  
DTD  
XML

### 3. Structure Extraction for an XML Document

XML -  
3 , DTD가 XML  
, -

. , XML

가

### 3.1 Structure Tree

XML

**3.1** T = (L, V, E, O)

- L

- V

- E (structure symbol)

가 XML DTD

XML

가

가

\*, +, ?가 \* 0

+

, ?

가

가

- O

가

가

DTD 가 DTD  
, DTD가

3.2 Structure Tree for Valid-Documents

DTD가 3.1  
DTD

**Rule 3.1** DTD L

**Rule 3.2** V 가  
가 PCDATA, CDATA, EMPTY, ANY  
가 CDATA, ID, IDREF, NMTOKEN, enumerated,

IDREFS, NMTOKENS

**Rule 3.3** E DTD  
가  
, \*, +, ?

**Rule 3.4** DTD 가  
O

**Sub-Rule 3.1** XML DTD ENTITY, ENTITIES

Sub-Rule 3.2 XML DTD

[13]

DTD

a , b  
c 가 가

$(a   b) \rightarrow a?, b?$	$a^{**} \rightarrow a^*$	$\dots, a^*, \dots, a^* \rightarrow a^*$
$(a, b)^* \rightarrow a^*, b^*$	$a^*? \rightarrow a^*$	$\dots, a^*, \dots, a^? \rightarrow a^*$
$(a, b)? \rightarrow a?, b?$	$a^{??} \rightarrow a?$	$\dots, a?, \dots, a? \rightarrow a^*$

Sub-Rule 3.3 DTD

Example 3.1

XML DTD

가 가

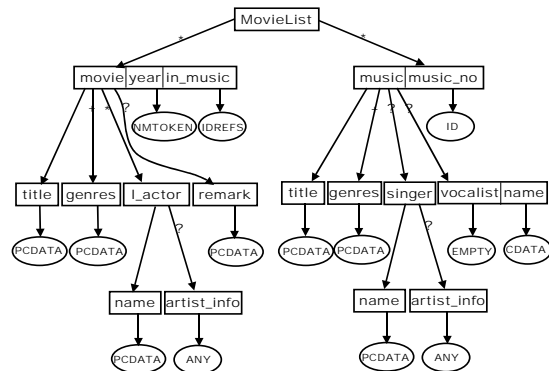
가

DTD

XML DTD

```

<!ELEMENT MovieList (movie, music)*>
<!ELEMENT movie (title, genres+, leading_actor*, remark?)>
<!ATTLIST movie year ID #IMPLIED>
<!ATTLIST movie insertion_music IDREFS #REQUIRED>
<!ELEMENT title (#PCDATA)>
<!ELEMENT genres (#PCDATA)>
<!ELEMENT leading_actor (name, artist_info?)>
<!ELEMENT name (#PCDATA)>
<!ELEMENT artist_info ANY>
<!ELEMENT remark (#PCDATA)>
<!ELEMENT music (title, genres+, (singer | vocalist))>
<!ATTLIST music music_no ID #REQUIRED>
<!ELEMENT singer (name, artist_info)>
<!ELEMENT vocalist EMPTY>
<!ATTLIST vocalist name CDATA #IMPLIED>
    
```



3-1 XML DTD

XML

DTD

movie

가

XML

3.1



3-1

가

singer

vocalist

DTD

( | )

3.2

?

### 3.3 Structure Tree for Well-formed Documents

XML DTD가

가

XML

가

(directed graph)

IDREF, IDREFS, XLink, URI

가

DTD가

XML

#### Rule 3.5 XML

0

$n$

$i(0 \leq i < n)$

$i+1$

가

가

$i$

$i+1$

가

•  $i$

$i+1$

가

0

\*

- 1 + .
- 가 .
- 0 1 가 ? .

3.2 XML 가 3-2 3.1 DTD

DTD가 .

.

1 가

movie A, B가 , music C, D, F가

movie(A, B) music(C, D, F) .

가 1 + . music(C, D, F)

singer 가 C F singer

가 D singer 가 ?

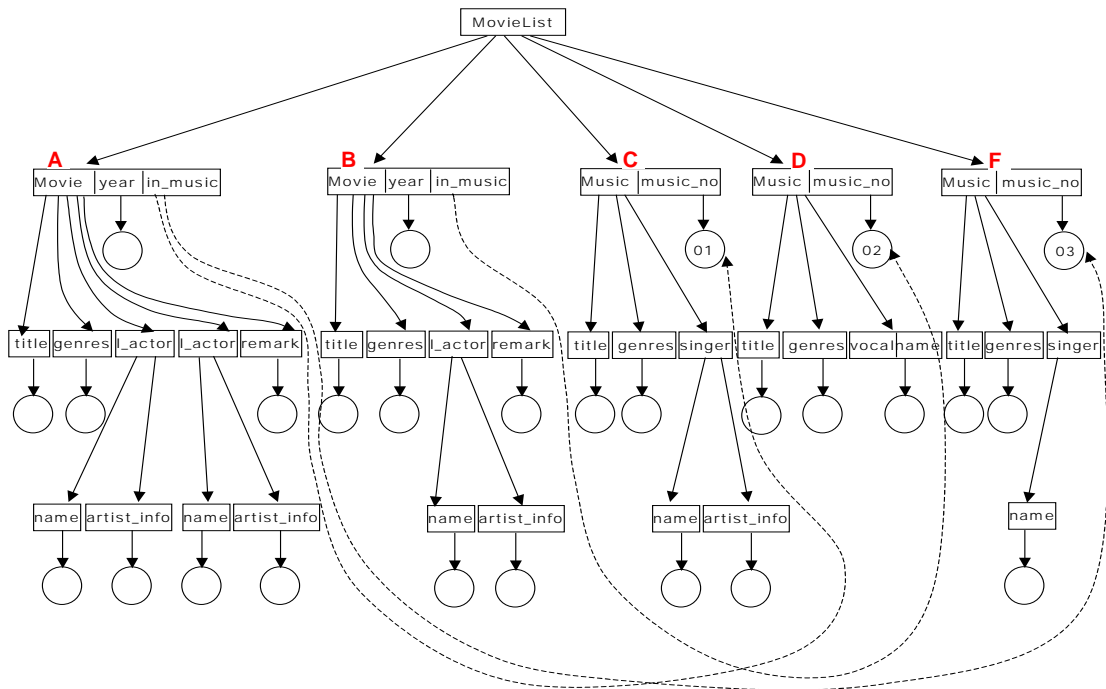
. Vocalist 가 . movie(A, B) title 가

title 가

. l\_actor A , B 가 +

. 3-2 3.5 3.6 3-3

.



3-2 XML

**Rule 3.6**

가

CDATA

,

가

가

NMTOKENS

. 가

EMPTY

.

가

CDATA ,

IDREFS

**Example 3.3**

3-2

CDATA

,

in\_music

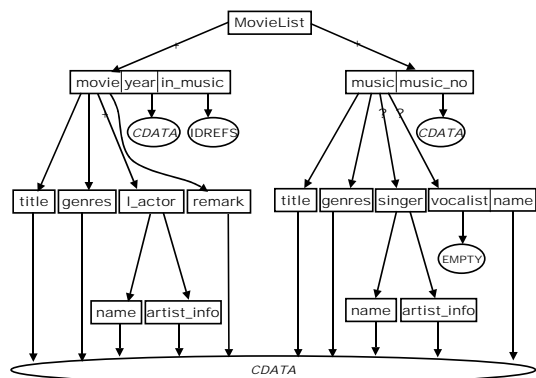
가

가

IDREFS

vocalist

가



3-3 3-2

EMPTY

3.5 3.6 3-2 XML

3-3 DTD

3-1

#### 4. Storing XML Documents into an Object-Relational Database

①

, ②XML 가

XML

③

가

#### 4.1 Rules for Mapping an XML Structure Tree into an Object-Relational Database

XML

##### Rule 4.1

##### Rule 4.2

A 가

A

A

A 가

\*

+

4.3

?

4.4

**Rule 4.3**

\* +

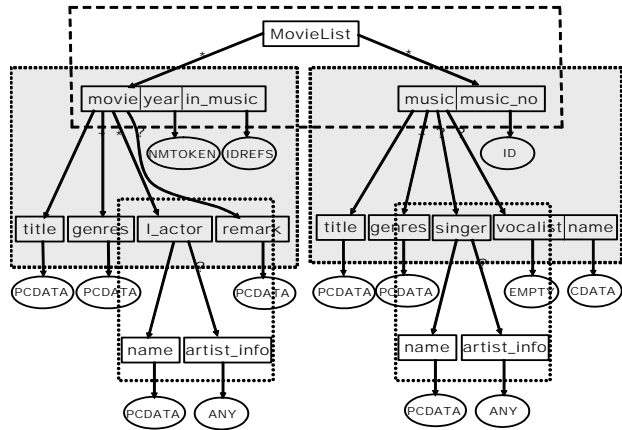
가 . ,

2

**Rule 4.4**

?

가



4-2

**Rule 4.1**

3-1

3-2

4.1-4.4

-

4-2

5 가

, , Movielist

movie music , movie year, in\_music,

title, genres, l\_actor, remark

**Rule 4.5**

PCDATA, CDATA, NMTOKEN, ANY

가

ID IDREF

**Rule 4.6**

EMPTY

**Rule 4.7**

IDREFS NMTOKENS

4.3

가

DTD가

3-3

#### 4.2 Assigning Identifiers to XML Data

4.1-4.7

, XML

XML

4.8

**Rule 4.8** XML

( $P_a, P_b$ )

.  $P_a, P_b$

{0,...,9}

,  $P_a$

,

,  $P_b$

,

$\alpha$

$P_a, P_b$

$P_a + P_b$ 가

$\beta$   $P_a$ 가

$P_b$

$\beta$

$P_a P_b$

**Example 4.2**

3-2

4-1

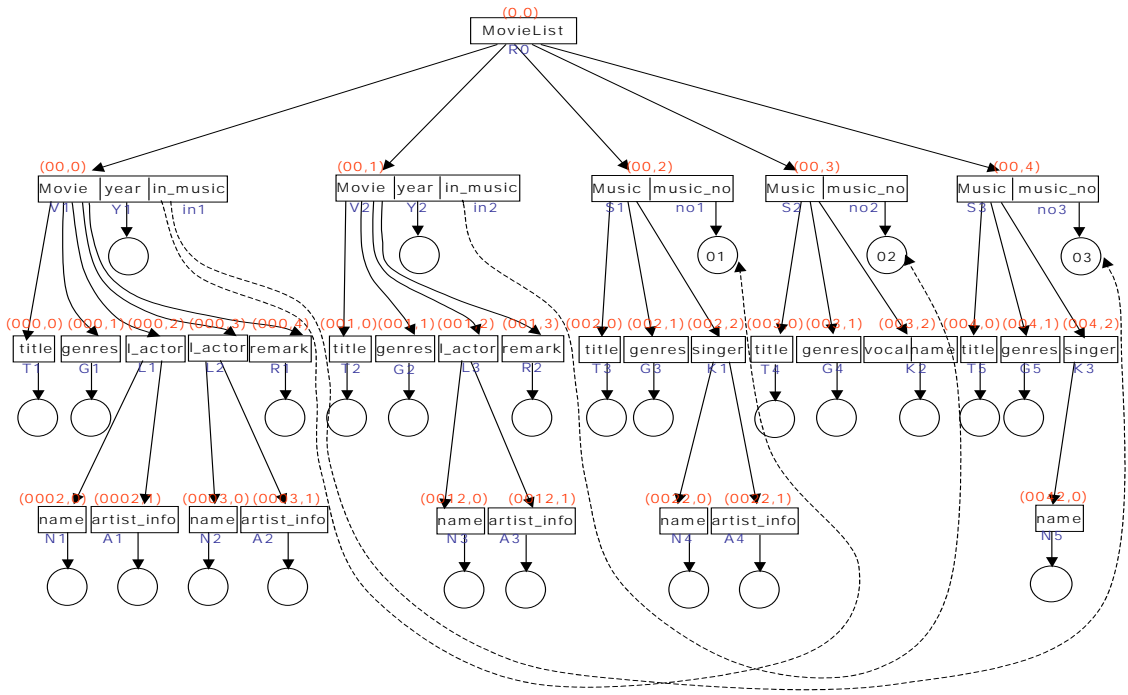
가

$P_a P_b$

.  $P_a P_b$

XML

[28].



4-1 XML

### 4.3 Storage Technique

XML

#### 4.3.1 Structure-Text Mixed Storage Technique

XML

가

가

ID(identification)

XML 가

4.2

ID

4.3

4.3

ID

( | ) ( a | b ) a?,

b?

a , b

ID

a 가 가 b

, a 가

b

3.1

3.3

( - i: - , ..., - n: - )

**Movilist** (MovieList\_ID:string, movie:reference, music:reference)

**Movie** (movie\_ID:string, year\_ID:string, year:string, in\_music(in\_music\_ID:string, in\_music:reference)  
title\_ID:string, title:string, genres(genres\_ID:string, genres:string), leading\_actor\_ID:string,  
leading\_actor:reference, remark\_ID:string, remark:string)

**Music** (music\_ID:string, music\_no\_ID:string, title\_ID:string, genres(genres\_ID:string, genres:string),  
SingerVocalist\_ID:string, singer:ref, vocalist\_name:string)

**L\_actor** (l\_actorID:string, name\_ID:string, actor\_info\_ID:string, actor\_info:string)

**Singer** (singerID:string, name\_ID:string, name:string, artist\_info\_ID:string, artist\_info:string)



XML

music

music ID	Music_noID	Music_no	Title ID	title	genres		singer	Vocal_noID	Vocal_no
					genre ID	genres	singer		
00,2	00,2	01	002,0	Scabro ough fair	002,1	pop/balla d	Ref(K1)	Null	Null
00,3	00,3	02	003,0	Two of us	003,1	pop/soft rock	Null	003,2	The Beatles
00,4	00,4	03	004,0	The sound of silence	004,1	pop/folk	Ref(K3)	null	null

4-2

가

4-2 -

가

가

+ \*

XML

XML

가

### 4.3.2 Structure-Text Separated Storage Technique

4.1

4.3

XML

(internal node)

(internal node)

**Movelist** (movelistID:string, movie:reference, music:reference)  
**Movie** (movieID:string, year:reference, in\_music(in\_music:reference), title:reference, genres(genres:reference), l\_actor(l\_actor:reference), remark:reference)  
**Music** (musicID:string, music\_no:reference, title:reference, genres(genres:reference), singer:reference, vocalist\_name:reference)  
**l-actor** (l\_actorID:string, name:reference, artist\_info:reference)  
**singer** (singerID:string, name:reference, artist\_info:reference)

**Year** (yearID:string, year:string)  
**In\_music** (in\_musicID:string, in\_music:string)  
**Music\_no** (music\_noID:string, music\_no:string)  
**Vocalist\_name** (vocal\_nameID:string, name:string)  
**Title** (titleID:string, title:string)  
**Remark** (remarkID:string, remark:string)  
**Genres** (genresID:string, genres:string)  
**name** (nameID:string, name:string)  
**Artist\_info** (artist\_info:string, artist\_info:string)

music					
musicID	music_no	title	genres	singer	Vocal_no
00,2	Ref(00,2)	Ref(003,0)	Ref(002,1)	Ref(002,2)	Null
00,3	Ref(00,3)	Ref(004,0)	Ref(003,1)	Null	Ref(003,2)
00,4	Ref(00,4)	Ref(005,0)	Ref(004,1)	Ref(004,2)	Null

music_no		genres	
music_noID	year	genresID	genres
00,2	01	000,1	Drama
00,3	02	001,1	Drama
00,4	03	002,1	pop/ballad
		003,1	pop/soft rock
		004,1	pop/folk

title	
titleID	title
000,0	The Graduate
001,0	I am Sam
002,0	Scaborough fair
003,0	Two of us
004,0	The sound of silence

vocal_no	
vocalID	name
003,2	The Beatles

4-3 -

가 . ,

4-3

genres

select \* from genres

가

, XML

XML

XML

, XML

가

### 4.3.3 Name-based Storage Technique

XML

가

EMPTY

가

가

<b>MovieList</b>	<b>(MovieListID:string, MovieList:string)</b>
<b>Movie</b>	<b>(movieID:string, movie:string)</b>
<b>Music</b>	<b>(musicID:string, music:string)</b>
<b>Title</b>	<b>(titleID:string, title:string)</b>
<b>Genres</b>	<b>(genresID:string, genres:string)</b>
<b>leading_actor</b>	<b>(l_actorID:string, l_actor:string)</b>
<b>remark</b>	<b>(remarkID:string, remark:string)</b>
<b>name</b>	<b>(nameID:string, name:string)</b>
<b>artist_info</b>	<b>(a_infolD:string, a_info:string)</b>
<b>singer</b>	<b>(singerID:string, singer:string)</b>
<b>vocalist_name</b>	<b>(v_nameID:string, v_name:string)</b>
<b>year</b>	<b>(yearID:string, year:string)</b>
<b>in_music</b>	<b>(in_musicID:string, in_music:string)</b>
<b>music_no</b>	<b>(music_noID:string, music_no:string)</b>

4-4 XML

가

XML

가

가

movie	
MovieID	Movie
V1(00,0)	<pre>&lt;movie year='1968' insertion_music='01 03'&gt; &lt;title&gt;The Graduate&lt;/title&gt; &lt;genres&gt;drama&lt;/genres&gt; &lt;leading_actor&gt; &lt;name&gt;Dustin Hoffman&lt;/name&gt; &lt;artist_info&gt;other_movies: Midnight Cowboy, Straw Dogs, Kramer Vs. Kramer, Rain man awards: 1979(Kramer Vs. Kramer), 1988(Rain man) &lt;/artist_info&gt; &lt;/leading_actor&gt; &lt;leading_actor&gt; &lt;name&gt;Katharine Ross&lt;/name&gt; &lt;artist_info&gt;other_movies: Butch Cassidy and The Sundance Kid, The Swarm&lt;/artist_info&gt; &lt;/leading_actor&gt; &lt;remark&gt;Torn apart by race riots in Detroit the preceding summer, and still reeling from the deaths of Martin Luther King and President Kennedy...&lt;/remark&gt; &lt;/movie&gt;</pre>
V2(00,1)	<pre>&lt;movie year='2001' insertion_music='02'&gt; &lt;title&gt;I am Sam&lt;/title&gt; &lt;genres&gt;drama&lt;/genres&gt; &lt;leading_actor&gt; &lt;name&gt; Sean Penn &lt;/name&gt; &lt;artist_info&gt;other_movies: Hurlyburly, She's so Lovely &lt;/artist_info&gt; &lt;remark&gt;I Am Sam is the compelling story of Sam Dawson a mentally-challenged father...&lt;/remark&gt; &lt;/movie&gt;</pre>

title		L-actor	
titleID	title	L-actorID	Leading_actor
000,0	<title>The Graduate</title>	L1(000,2)	<pre>&lt;leading_actor&gt; &lt;name&gt;Dustin Hoffman&lt;/name&gt; &lt;artist_info&gt;other_movies: Midnight Cowboy, Straw Dogs, Kramer Vs. Kramer, Rain man awards: 1979(Kramer Vs. Kramer), 1988(Rain man)... &lt;/artist_info&gt; &lt;/leading_actor&gt;</pre>
001,0	<title>I am Sam</title>	L2(000,3)	<pre>&lt;leading_actor&gt; &lt;name&gt;Katharine Ross&lt;/name&gt; &lt;artist_info&gt;other_movies: Butch Cassidy and The Sundance Kid, The Swarm &lt;/artist_info&gt; &lt;/leading_actor&gt;</pre>
002,0	<title>Scaborough fair</title>	L3(001,2)	<pre>&lt;leading_actor&gt; &lt;name&gt; Sean Penn &lt;/name&gt; &lt;artist_info&gt;other_movies: Hurlyburly, She's so Lovely... &lt;/artist_info&gt; &lt;/leading_actor&gt;</pre>
003,0	<title>Two of us</title>		
004,0	<title>The sound of silence</title>		

genres	
genresID	genres
G1(000,1)	<genres>drama</genres>
G2(001,1)	<genres>drama</genres>
G3(002,1)	<genres>pop/ballad</genres>
G4(003,1)	<genres>pop/soft rock</genres>
G5(004,1)	<genres>pop/folk</genres>

4-4

### 5. Performance Evaluation

XML

[29],

[28].

XML

가

가

가

가

XML

XML

, XML

XML 가 가 가

가

STORED[8], [9] 8가 가

가 (attribute table with inlining), XML-DBMS [26], [13]

가 (hybrid inlining), 가

5-1

**Table 5-1 Comparison of Storage Techniques**

Consideration Items		STORED	Attribute table with inlining	XML-DBMS	Hybrid inlining	Structure-text Mixed ST	Structure-text separate ST	Name-based ST
Preserve XML properties?	Hierarchical information	x						
	Order information	x			x			
Process completely irregular and incomplete structures?	Optional elements							
	*, ?, + type elements							
	ANY type							
	Multi-value attributes				x			
Minimize data redundancies?								x
Minimize null values or overhead data?			x	x	x			
Minimize table fragmentations?			x				x	x
Is reusable as relational databases seamlessly?			x					x
Minimize post-processing to return XML documents			x					

× , ,  
 .  
 XML 가  
 SQL  
 × , 가  
 .  
 XML 가 XML  
 .  
 , × .  
 XML 10%~30%  
 , × .  
 XML 가  
 가  
 ,  
 20%~40% ,  
 × .  
 XML 가  
 . , XML  
 가  
 가 .  
 ×,  
 가 가 XML 가 가

가

가 . , , ,

2 x,

XML

가

가

가 (degree)

x, 가 가

가

가

6. Conclusion

XML 가 DTD 가 가

XML 가 가

XML

, XML 가

References

[1] W3C Consortium, XML1.0.(Second Edition), W3C Recommendation, 6 Oct. 2000, available at <http://www.w3.org/TR/2000/WD-xml-2e-20000814>

- [2] V. Aguilera, S. Cluet, P. Veltri, D. Vodislav, and F. Watez, Querying XML Documents in Xyleme, SIGIR, 2000
- [3] Z. G. Ives, A. Y. Levy, and D. S. Weld, Efficient Evaluation of Regular path Expressions on Streaming XML Data, Technical report, 2000
- [4] J. McHugh, S. Abiteboul, R. Goldman, D. Quass, and J. Widom. Lore: A Database Management System for Semistructured Data. SIGMOD Record, 26(3):54-66, September 1997
- [5] S. Banerjee, Oracle XML DB, Oracle Corporation Technical White Paper Release 9.2, Jan. 2002
- [6] S. Howlett and D. Jennings, SQL Server 2000 and XML: Developing XML-Enabled Data Solutions for the Web, MSDN magazine, Jan. 2002, available at <http://msdn.microsoft.com/library/default.asp?url=/msdnmag/issues/0800/sql2000/toc.asp>
- [7] IBM Corporation, DB2 XML Extender, IBM Corporation, 2000, available at <http://www-4.ibm.com/>
- [8] A. Deutsch, M. Fernandez, and D. Suciu, Storing Semistructured Data with STORED, SIGMOD, Philadelphia, PN, 1999
- [9] D. Florescu and D. Kossman, A Performance Evaluation of Alternative Mapping Schemes for Storing XML Data in a Relational Database, INRIA, Rocquencourt, France, 1999
- [10] Q. Li and B. Moon, Indexing and Querying XML Data for Regular Path Expressions, VLDB, 2001
- [11] F. Rizzolo and A. Mendlzon, Indexing XML Data with ToXin, 4<sup>th</sup> Int. Workshop on the Web and Database, 2001
- [12] J. Shanmugasundaram, E. Shekita, R. Barr, M. Carey, B. R. B. Lindsay, and H. Pirahesh, Efficiently publishing Relational Data as XML Documents, VLDB, 2000
- [13] J. Shanmugasundaram, K. Tuffe, G. He, C. Zhang, D. DeWitt, and J. Naughton, Relational Databases for Querying XML Documents: Limitations and Opportunities, VLDB, 1999
- [14] I. Tatarinov, S. D. Viglas, K. Beyer, J. Shanmugasundaram, E. Shekita, and C. Zhang, Storing and Querying Ordered XML Using a Relational Database System, SIGMOD, 2002
- [15] C. Zhang, J. Naughton, D. DeWitt, Q. Luo, and G. Lohman, On Supporting Containment Queries in Relational Database Management Systems, SIGMOD, 2001
- [16] S. Adler, A. Berglund, J. Caruso, S. Deach, T. Graham, P. Grosso, E. Gutentag, A. Milowski, S. Parnell, J. Richman, S. Zilles, Extensible Stylesheet Language (XSL) Version 1.0, W3C Proposed Recommendation Aug. 2001, available at <http://www.w3.org/TR/xsl/>
- [17] S. Abiteboul, D. Quass, J. McHugh, J. Widom, and J. L. Wiener, The Lorel Query Language for Semistructured Data, International Journal on Digital Libraries, Apr. 1997
- [18] S. Boag, D. Chamberlin, M. F. Fernandez, D. Florescu, J. Robie, J. Siméon, XQuery 1.0: An XML Query Language, W3C Working Draft 16 Aug. 2002, available at <http://www.w3.org/TR/xquery/>
- [19] A. Deutsch, M. Fernandez, D. Florescu, A. Levy, and D. Suciu, XML-QL: A Query Language for XML, Submitted to the W3C 19 Aug. 1998, available at <http://www.w3.org/TR/1998/NOTE-xml-ql-19980819>



- [20] J. Robie, XQL (XML Query Language), Aug. 1999, available at <http://www.ibiblio.org/xql/xql-proposal.html>
- [21] W3C Consortium, XML Path Language (XPath) Version 1.0, W3C Recommendation 16 Nov. 1999, available at <http://www.w3.org/TR/xpath.html>
- [22] M. Fernandez, W. -C. Tan, and D. Suciu, SilkRoute: Trading between relational and XML In Proc. of the WWW9, 2000
- [23] T. Shimura, M. Yoshikawa, and S. Uemura, Storage and Retrieval of XML Documents Using Object-Relational Databases, DEXA, 1999.
- [24] M. J. Carey, D. Florescu, Z. G. Lves, Y. Lu, and J. Shanmugasundaram, E. Shekita, and S. Subramannian, XPERANTO: Publishing Object-Relational Data as XML, In Proc. of the Int. Workshop on Web and Databases, 2000
- [25] I. Tatarinov, Z. G. Ives, A. Y. Halevy, and D. S. Weld, Updating XML, SIGMOD, 2001
- [26] R. Bourret, XML-DBMS: Middleware for Transferring Data between XML Documents and Relational Databases, available at <http://www.rpbouret.com/xmldbms/>
- [27] K. Runapongsa and J. M. Patel, Storing and querying XML data in ORDBMSs, EDBT workshop, 2002
- [28] , , XML , ,
- [29] J. Kim, W. Lee, K. Lee. The Cost Model for XML Documents, In Proc. of ACS/IEEE International Conference on Computer Systems and Applications, Beirut, Lebanon, Jun. 2001.
- [30] R. Murthy and S. Banerjee, XML Schema in Oracle XML DB, In Proc. Of the 29<sup>th</sup> VLDB Conference, Berlin, Germany, 2003