

# XML-Relational Mapping

CPS 216  
Advanced Database Systems

---

---

---

---

---

---

---

---

## Announcements (March 18)

2

- ❖ Midterm sample solution available outside my office
- ❖ Course project milestone 2 due March 30
- ❖ Homework #3 due April 6
- ❖ Talk by Amol Deshpande
  - Adaptive Query Processing to Handle Estimation Errors
  - Monday, 11:30am-12:30pm, D106
- ❖ Reading assignment due next Monday
  - Two *VLDB* papers on native XML databases

---

---

---

---

---

---

---

---

## Approaches to XML processing

3

- ❖ Text files (!)
- ❖ Specialized XML DBMS
  - Lore (Stanford), Strudel (AT&T), Tamino/QuiP (Software AG), X-Hive, Timber (Michigan), etc.
  - Still a long way to go
- ❖ Object-oriented DBMS
  - eXcelon (ObjectStore), ozone, etc.
  - Not as mature as relational DBMS
- ❖ Relational (and object-relational) DBMS
  - Middleware and/or object-relational extensions

---

---

---

---

---

---

---

---

## Mapping XML to relational

4

- ❖ Store XML in a CLOB (Character Large Object) column
  - Simple, compact
  - Full-text indexing can help (often provided by DBMS vendors as object-relational “extensions”)
- ❖ Alternatives?
  - Schema-oblivious mapping: well-formed XML → generic relational schema
    - Node/edge-based mapping for graphs
    - Interval-based mapping for trees
    - Path-based mapping for trees
  - Schema-aware mapping: valid XML → special relational schema based on DTD

---

---

---

---

---

---

---

---

---

---

## Node/edge-based: schema

5

- ❖ *Element(eid, tag)*
  - ❖ *Attribute(eid, attrName, attrValue)*
    - Attribute order does not matter
  - ❖ *ElementChild(eid, pos, child)*
    - *pos* specifies the ordering of children
    - *child* references either *Element(eid)* or *Text(tid)*
  - ❖ *Text(tid, value)*
    - *tid* cannot be the same as any *eid*
- ☞ Need to “invent” lots of *id*'s
- ☞ Need indexes for efficiency, e.g., *Element(tag)*, *Text(value)*

---

---

---

---

---

---

---

---

---

---

## Node/edge-based: example

6

```
<bibliography>
<book ISBN="ISBN-10" price="80.00">
<title>Foundations of Databases</title>
<author>Abiteboul</author>
<author>Hull</author>
<author>Vianu</author>
<publisher>Addison Wesley</publisher>
<year>1995</year>
</book>
</bibliography>
```

*Element*

eid	tag
e0	bibliography
e1	book
e2	title
e3	author
e4	author
e5	author
e6	publisher
e7	year

*ElementChild*

eid	pos	child
e0	1	e1
e1	1	e2
e1	2	e3
e1	3	e4
e1	4	e5
e1	5	e6
e1	6	e7
e2	1	t0
e3	1	t1
e4	1	t2
e5	1	t3
e6	1	t4
e7	1	t5

*Attribute*

eid	attrName	attrValue
e1	ISBN	ISBN-10
e1	price	80

*Text*

tid	value
t0	Foundations of Databases
t1	Abiteboul
t2	Hull
t3	Vianu
t4	Addison Wesley
t5	1995

---

---

---

---

---

---

---

---

---

---

## Node/edge-based: simple paths <sup>7</sup>

- ❖ //title
  - SELECT eid FROM Element WHERE tag = 'title';
- ❖ //section/title
  - SELECT e2.eid  
FROM Element e1, ElementChild c, Element e2  
WHERE e1.tag = 'section'  
AND e2.tag = 'title'  
AND e1.eid = c.eid  
AND c.child = e2.eid;

---

---

---

---

---

---

---

---

## Node/edge-based: more complex paths <sup>8</sup>

- ❖ //bibliography/book[author="Abiteboul"]/@price
  - SELECT a.attrValue  
FROM Element e1, ElementChild c1,  
Element e2, Attribute a  
WHERE e1.tag = 'bibliography'  
AND e1.eid = c1.eid AND c1.child = e2.eid  
AND e2.tag = 'book'  
AND EXISTS (SELECT \* FROM ElementChild c2,  
Element e3, ElementChild c3, Text t  
WHERE e2.eid = c2.eid AND c2.child = e3.eid  
AND e3.tag = 'author'  
AND e2.eid = c3.eid AND c3.child = t.tid  
AND t.value = 'Abiteboul')  
AND e2.eid = a.eid  
AND a.attrName = 'price';

---

---

---

---

---

---

---

---

## Node/edge-based: descendent-or-self <sup>9</sup>

- ❖ //book//title

---

---

---

---

---

---

---

---

# Interval-based: schema

- ❖ *Element(left, right, level, tag)*
  - *left* is the start position of the element
  - *right* is the end position of the element
  - *level* is the nesting depth of the element (strictly speaking, unnecessary)
  - Key is *left*
- ❖ *Attribute(left, attrName, attrValue)*
- ❖ *Text(left, level, value)*
- ☞ Where did *ElementChild* go?

---

---

---

---

---

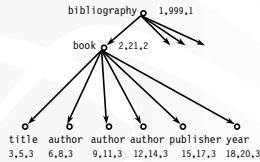
---

---

---

# Interval-based: example

```
1<bibliography>
2<book ISBN="ISBN-10" price="80.00">
3<title>4Foundations of Databases</title>5
6<author>7Abiteboul</author>8
9<author>10Hull</author>11
12<author>13Vianu</author>14
15<publisher>16Addison Wesley</publisher>17
18<year>191995</year>20
</book>21_
</bibliography>999
```




---

---

---

---

---

---

---

---

# Interval-based: queries

- ❖ `//section/title`
  - `SELECT e2.left`  
`FROM Element e1, Element e2`  
`WHERE e1.tag = 'section' AND e2.tag = 'title'`  
`AND e1.left < e2.left AND e2.right < e1.right`  
`AND e1.level = e2.level-1;`
  - ☞ Path expression becomes “containment” joins!
    - Number of joins is proportional to path expression length
- ❖ `//book//title`
  - `SELECT e2.left`  
`FROM Element e1, Element e2`  
`WHERE e1.tag = 'book' AND e2.tag = 'section'`  
`AND e1.left < e2.left AND e2.right < e1.right;`
  - ☞ No recursion!

---

---

---

---

---

---

---

---

# How about XQuery?

DeHaan et al. SIGMOD 2003

- ❖ Evaluating an XQuery expression results in a sequence of environments
  - An environment  $E$  maps each query variable  $v$  to its value: a forest of XML trees (a node-set)  $f_v$
- ❖ Encode using tables with “dynamic intervals”
  - Table  $I$ : increasing sequence of integers, one per environment
  - For each query variable  $v$ , create a table  $T_v(s(tring), l(left), r(right))$  representing the value of  $v$  in all environments
    - Sorted on  $l$  to support efficient processing
    - Different environments form non-overlapping regions

---

---

---

---

---

---

---

---

---

---

# Example $T_v$

$I'$	$T_v$			
$i$	$s$	$l$	$r$	
2	<person>	174	195	
	@id	175	178	
	person0	176	177	
	<name>	179	182	
	Jaak Tempesti	180	181	
	<emailaddress>	183	186	
	mailto:Tempesti@labs.com	184	185	
	<phone>	187	190	
	+0 (873) 14873867	188	189	
	<homepage>	191	194	
	http://www.labs.com/~Tempesti	192	193	
	24	<person>	2088	2109
		@id	2089	2092
person1		2090	2091	
<name>		2093	2096	
Cong Rosca		2094	2095	
<emailaddress>		2097	2100	
mailto:Rosca@washington.edu		2098	2099	
<phone>		2101	2104	
+0 (64) 27711230		2102	2103	
<homepage>		2105	2108	
http://www.washington.edu/~Rosca	2106	2107		
$w_p = 86$				

---

---

---

---

---

---

---

---

---

---

# Translating /

- ❖ Given  $T_v$  for values of  $v$ , compute  $v/\text{name}$

---

---

---

---

---

---

---

---

---

---

## Translating //

16

❖ Given  $T_v$  for values of  $v$ , compute  $v//*$

---

---

---

---

---

---

---

---

## Translating for

17



Figure 6: TRANSLATION OF "for  $x \in e$  do  $e'$ " IN THE ENVIRONMENT FOR  $x_1, \dots, x_m$ .

---

---

---

---

---

---

---

---

## Summary of interval-based mapping

18

- ❖ Path expression steps become containment joins
- ❖ No recursion needed for descendant-or-self
- ❖ Comprehensive XQuery-SQL translation is possible with dynamic interval encoding
  - Looks hairy, but with some special tweaks to the relational engine, it actually performs better than many of the currently available native XQuery products!
- ☞ Set-oriented processing helps!

---

---

---

---

---

---

---

---

# A path-based mapping

## Label-path encoding

❖ *Element(pathid, left, right, value), Path(pathid, path)*

- *path* is a label path starting from the root
- Why are *left* and *right* still needed?

*Element*

<i>pathid</i>	<i>left</i>	<i>right</i>	...
1	1	999	...
2	2	21	...
3	3	5	...
4	6	8	...
4	9	11	...
4	12	14	...
...	...	...	...

*Path*

<i>pathid</i>	<i>path</i>
1	/bibliography
2	/bibliography/book
3	/bibliography/book/title
4	/bibliography/book/author
...	...

---

---

---

---

---

---

---

---

---

---

# Label-path encoding: queries

❖ Simple path expressions with no conditions

`//book//title`

- Perform string matching on *Path*
- Join qualified *pathid*'s with *Element*

❖ Path expression with attached conditions need to be broken down, processed separately, and joined back

`//book[publisher='Prentice Hall']/title`

- Evaluate `//book`
- Evaluate `//book/title`
- Evaluate `//book/publisher[text()='Prentice Hall']`
- Join to ensure `title` and `publisher` belong to the same `book`

---

---

---

---

---

---

---

---

---

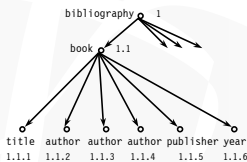
---

# Another path-based mapping

## Dewey-order encoding

❖ Each component of the id represents the order of the child within its parent

- Unlike label-path, this encoding is "lossless"




---

---

---

---

---

---

---

---

---

---

## Dewey-order encoding: queries

22

### ❖ Examples:

```
//title
//section/title
//book//title
//book[publisher='Prentice Hall']/title
```

---

---

---

---

---

---

---

---

## Schema-aware mapping

23

- ❖ Idea: use DTD to design a better schema
- ❖ Basic approach: elements of the same type go into one table
  - Tag name → table name
  - Attributes → columns
    - If one exists, ID attribute → key column; otherwise, need to “invent” a key
    - IDREF attribute → foreign key column
  - Children of the element → foreign key columns
    - Ordering of columns encodes ordering of children

```
<!DOCTYPE bibliography [-
<ELEMENT book (title, -)>           book(ISBN, price, title_id, ...)
<!ATTLIST book ISBN ID #REQUIRED>
<!ATTLIST book price CDATA #IMPLIED> title(id, PCDATA_id)
<ELEMENT title (#PCDATA)->        PCDATA(id, value)
]>
```

---

---

---

---

---

---

---

---

## Handling \* and + in DTD

24

- ❖ What if an element can have any number of children?
- ❖ Example: Book can have multiple authors
  - *book(ISBN, price, title\_id, author\_id, publisher\_id, year\_id)?*
  - ☞ BCNF?
- ❖ Idea: create another table to track such relationships
  - *book(ISBN, price, title\_id, publisher\_id, year\_id)*
  - *book\_author(ISBN, author\_id)*
  - ☞ BCNF decomposition in action!
  - ☞ A further optimization: merge *book\_author* into *author*
- ❖ Need to add position information if ordering is important
  - *book\_author(ISBN, author\_pos, author\_id)*

---

---

---

---

---

---

---

---



## Inlining

- ❖ An author element just has a PCDATA child
- ❖ Instead of using foreign keys
  - *book\_author(ISBN, author\_id)*
  - *author(id, PCDATA\_id)*
  - *PCDATA(id, value)*
- ❖ Why not just “inline” the string value inside *book*?
  - *book\_author(ISBN, author\_PCDATA\_value)*
  - *PCDATA* table no longer stores author values

---

---

---

---

---

---

---

---

## More general inlining

- ❖ As long as we know the structure of an element and its number of children (and recursively for all children), we can inline this element where it appears

```
<book ISBN="...">...
  <publisher>
    <name>...</name><address>...</address>
  </publisher>...
</book>
```

- ❖ With no inlining at all
- ❖ With inlining

<i>book(ISBN, publisher_id)</i>	<i>book(ISBN,</i>
<i>publisher(id, name_id, address_id)</i>	<i>publisher_name_PCDATA_value,</i>
<i>name(id, PCDATA_id)</i>	<i>publisher_address_PCDATA_value)</i>
<i>address(id, PCDATA_id)</i>	

---

---

---

---

---

---

---

---

## Queries

- ❖ *book(ISBN, price, title, publisher, year),*  
*book\_author(ISBN, author), book\_section(ISBN, section\_id),*  
*section(id, title, text), section\_section(id, section\_pos, section\_id)*
- ❖ *//title*
- ❖ *//section/title*
- ❖ *//bibliography/book[author="Abiteboul"]/@price*
- ❖ *//book//title*

These queries only work for the given DTD

---

---

---

---

---

---

---

---

## Pros and cons of inlining

28

- ❖ Not always applicable

---

---

---

---

---

---

---

---

## Result restructuring

29

- ❖ Simple results are fine
  - Each tuple returned by SQL gets converted to an element
- ❖ Simple grouping is fine (e.g., books with multiple authors)
  - Tuples can be returned by SQL in sorted order; adjacent tuples are grouped into an element
- ❖ Complex results are problematic: one SQL query only returns a single table; columns cannot contain sets or structures
  - E.g., books with multiple authors and multiple references
    - Option 1: one table with all combo of authors/references → bad
    - Option 2: two tables, one w/ authors and the other w/ references → join is done as post processing
    - Option 3: sorted "union" of NULL-padded authors and references

---

---

---

---

---

---

---

---

## Comparison of approaches

30

- ❖ Schema-oblivious
  - Flexible and adaptable; no DTD needed
  - Queries are easy to formulate
    - Translation from XPath/XQuery can be easily automated
  - Queries involve lots of join and are expensive
- ❖ Schema-aware
  - Less flexible and adaptable
  - Need to know DTD to design the relational schema
  - Query formulation requires knowing DTD and schema
  - Queries are more efficient
  - XQuery is tougher to formulate because of result restructuring

---

---

---

---

---

---

---

---