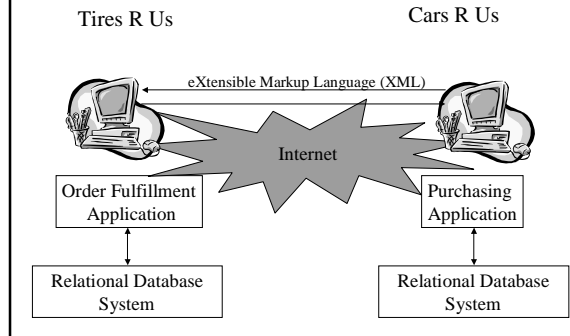


Bridging Relational Technology and XML

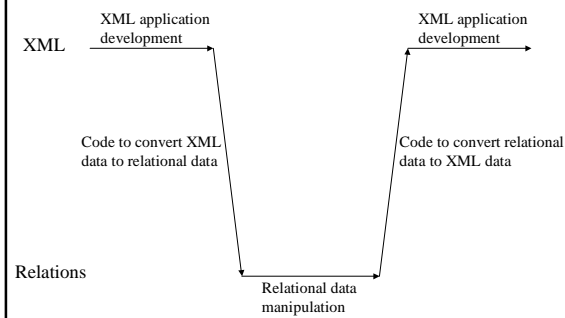
Jayavel Shanmugasundaram

Cornell University
(Work done at University of Wisconsin & IBM Almaden Research Center)

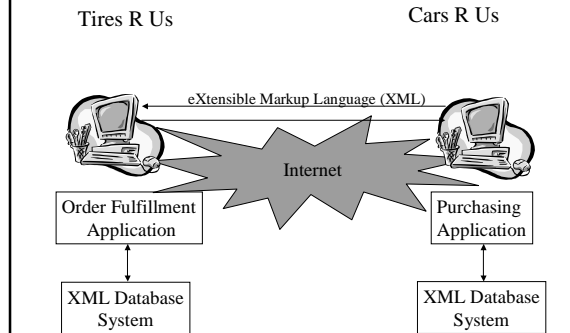
Business to Business Interactions



Shift in Application Developers' Conceptual Data Model



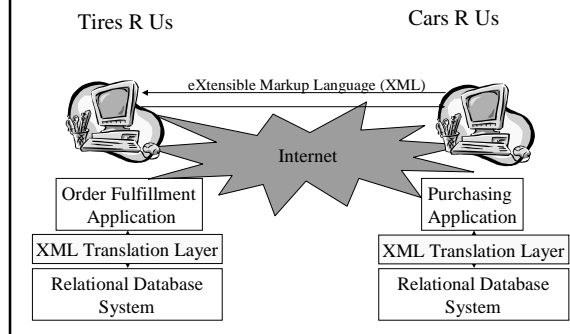
Are XML Database Systems the Answer?



Why use Relational Database Systems?

- Highly reliable, scalable, optimized for performance, advanced functionality
 - Result of 30+ years of Research & Development
 - XML database systems are not “industrial strength” ... and not expected to be in the foreseeable future
- Existing data and applications
 - XML applications have to inter-operate with existing relational data and applications
 - Not enough incentive to move all existing business applications to XML database systems
 - Remember object-oriented database systems?

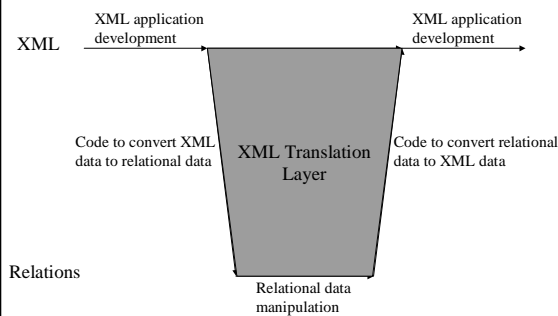
A Solution



XML Translation Layer (Contributions)

- Store and query XML documents
 - Harnesses relational database technology for this purpose [VLDB'99]
- Publish existing relational data as XML documents
 - Allows relational data to be viewed in XML terms [VLDB'00]

Bridging Relational Technology and XML



Outline

- Motivation & High-level Solution
- Background (Relations, XML)
- Storing and Querying XML Documents
- Publishing Relational Data as XML Documents
- Conclusion

Relational Data

PurchaseOrder

Id	Customer	Day	Month	Year
2001	Cars R Us	10	June	1999
3001	Bikes R Us	null	July	1999

Item

Pid	Name	Quantity	Cost
2001	Firestone Tire	50	2000.00
3001	Schwinn Tire	100	2500.00
3001	Trek Tire	20	400.00
2001	Goodyear Tire	200	8000.00

Payment

Pid	Installment	Percentage
2001	1	40%
2001	2	60%
3001	1	100%

SQL Query

Find all the items bought by "Cars R Us" in the year 1999

```

Select it.name
From PurchaseOrder po, Item it
Where po.customer = "Cars R Us" and
      po.year = 1999 and
      po.id = it.pid

```

} Predicates
} Join

XML Document

Self-describing tags

```

<PurchaseOrder id="2001" customer="Cars R Us">
  <Date>
    <Day> 10 </Day>
    <Month> June </Month>
    <Year> 1999 </Year>
  </Date>
  <Item name="Firestone Tire" cost="2000.00">
    <Quantity> 50 </Quantity>
  </Item>
  <Item name="Goodyear Tire" cost="8000.00">
    <Quantity> 200 </Quantity>
  </Item>
  <Payment> 40% </Payment>
  <Payment> 60% </Payment>
</PurchaseOrder>

```

XML Document

Self-describing tags `<PurchaseOrder id="2001" customer="Cars R Us">`
`<Date>`
`<Day> 10 </Day>`
`<Month> June </Month>`
`<Year> 1999 </Year>`
`</Date>`
 Nested structure
`<Item name="Firestone Tire" cost="2000.00">`
`<Quantity> 50 </Quantity>`
`</Item>`
`<Item name="Goodyear Tire" cost="8000.00">`
`<Quantity> 200 </Quantity>`
`</Item>`
`<Payment> 40% </Payment>`
`<Payment> 60% </Payment>`
`</PurchaseOrder>`

XML Document

Self-describing tags `<PurchaseOrder id="2001" customer="Cars R Us">`
`<Date>`
`<Day> 10 </Day>`
`<Month> June </Month>`
`<Year> 1999 </Year>`
`</Date>`
 Nested structure
`<Item name="Firestone Tire" cost="2000.00">`
`<Quantity> 50 </Quantity>`
`</Item>`
 Nested sets
`<Item name="Goodyear Tire" cost="8000.00">`
`<Quantity> 200 </Quantity>`
`</Item>`
`<Payment> 40% </Payment>`
`<Payment> 60% </Payment>`
`</PurchaseOrder>`

XML Document

Self-describing tags `<PurchaseOrder id="2001" customer="Cars R Us">`
`<Date>`
`<Day> 10 </Day>`
`<Month> June </Month>`
`<Year> 1999 </Year>`
`</Date>`
 Nested structure
`<Item name="Firestone Tire" cost="2000.00">`
`<Quantity> 50 </Quantity>`
`</Item>`
 Nested sets
`<Item name="Goodyear Tire" cost="8000.00">`
`<Quantity> 200 </Quantity>`
`</Item>`
 Order
`<Payment> 40% </Payment>`
`<Payment> 60% </Payment>`
`</PurchaseOrder>`

XML Schema

PurchaseOrder → `<PurchaseOrder id={integer} customer={string}>`
`Date (Item)* (Payment)*`
`</PurchaseOrder>`

Date → `<Date>`
`Day? Month Year`
`</Date>`

Day → `<Day> {integer} </Day>`

Month → `<Month> {string} </Month>`

Year → `<Year> {integer} </Year>`

Item → `<Item name={string} cost={float}>`
`Quantity`
`</Item>`

... and so on

XML Schema (contd.)

PurchaseOrder → `<PurchaseOrder id={integer} customer={string}>`
`Date? (Item | Payment)*`
`</PurchaseOrder>`

PurchaseOrder → `<PurchaseOrder id={integer} customer={string}>`
`(Date | Payment*) (Item (Item Item)* Payment)*`
`</PurchaseOrder>`

PurchaseOrder → `<PurchaseOrder id={integer} customer={string}>`
`Date Item (PurchaseOrder)* Payment`
`</PurchaseOrder>`

XML Query

Find all the items bought by "Cars R Us" in 1999

For \$po in /PurchaseOrder
 Where \$po/@customer = "Cars R Us" and \$po/date/year = 1999
 Return \$po/Item

XML Query (contd.)

//Item

//Item[5]

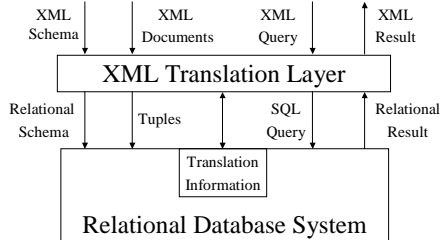
//Item Before //Payment

/(Item/(Item/Payment)*/(Payment | Item))*/Date

Outline

- Motivation & High-level Solution
- Background (Relations, XML)
- Storing and Querying XML Documents
- Publishing Relational Data as XML Documents
- Conclusion

Storing and Querying XML Documents [Shanmugasundaram et. al., VLDB'99]



Outline

- Motivation & High-level Solution
- Background (Relations, XML)
- Storing and Querying XML Documents
 - Relational Schema Design and XML Storage
 - Query Mapping and Result Construction
- Publishing Relational Data as XML Documents
- Conclusion

XML Schema

PurchaseOrder → <PurchaseOrder id={integer} customer={string}>
(Date | (Payment)* (Item (Item Item)* Payment))*
</PurchaseOrder>

Desired Properties of Generated Relational Schema \mathcal{R}

- All XML documents conforming to XML schema should be “mappable” to tuples in \mathcal{R}
- All queries over XML documents should be “mappable” to SQL queries over \mathcal{R}
- Not Required: Ability to re-generate XML schema from \mathcal{R}

Simplifying XML Schemas

- XML schemas can be “simplified” for translation purposes

PurchaseOrder \rightarrow \langle PurchaseOrder id={integer} customer={string}>
(Date | (Payment)* (Item (Item Item)* Payment))*
 \langle /PurchaseOrder



PurchaseOrder \rightarrow \langle PurchaseOrder id={integer} customer={string}>
Date? (Item)* (Payment)*
 \langle /PurchaseOrder

- All without undermining storage and query functionality!

Why is Simplification Possible?

- Structure in XML schemas can be captured:
 - Partly in relational schema
 - Partly as *data values*

PurchaseOrder \rightarrow \langle PurchaseOrder id={integer} customer={string}>
Date? (Item)* (Payment)*
 \langle /PurchaseOrder

- Order field to capture order among siblings
- Sufficient to answer ordered XML queries
 - PurchaseOrder/Item[5]
 - PurchaseOrder/Item AFTER PurchaseOrder/Payment
- Sufficient to reconstruct XML document

Simplification Desiderata

- Simplify structure, but preserve differences that matter in relational model
 - Single occurrence (attribute)
 - Zero or one occurrences (nullable attribute)
 - Zero or more occurrences (relation)

PurchaseOrder \rightarrow \langle PurchaseOrder id={integer} customer={string}>
(Date | (Payment)* (Item (Item Item)* Payment))*
 \langle /PurchaseOrder



PurchaseOrder \rightarrow \langle PurchaseOrder id={integer} customer={string}>
Date? (Item)* (Payment)*
 \langle /PurchaseOrder

Translation Normal Form

- An XML schema production is either of the form:

$P \rightarrow \langle P \text{ attr}_1=\{\text{type}_1\} \dots \text{attr}_m=\{\text{type}_m\} \rangle$
 $a_1 \dots a_p a_{p+1}^? \dots a_q^? a_{q+1}^* \dots a_r^*$
 $\langle /P \rangle$

where $a_i \neq a_j$

- ... or of the form:

$P \rightarrow \langle P \text{ attr}_1=\{\text{type}_1\} \dots \text{attr}_m=\{\text{type}_m\} \rangle$
 $\{\text{type}\}$
 $\langle /P \rangle$

Example Simplification Rules

$(e_1 | e_2) \rightarrow e_1^? e_2^?$

$(\text{Date} | (\text{Payment})^*) (\text{Item} (\text{Item Item})^* \text{Payment})^*$



$\text{Date}^? (\text{Item})^* (\text{Item} (\text{Item Item})^* \text{Payment})^*$

$e^* \rightarrow e^*$

$\text{Date}^? (\text{Payment})^* (\text{Item} (\text{Item Item})^* \text{Payment})^*$



$\text{Date}^? (\text{Item})^* (\text{Item} (\text{Item Item})^* \text{Payment})^*$

Simplified XML Schema

PurchaseOrder \rightarrow \langle PurchaseOrder id={integer} customer={string}>
Date (Item)* (Payment)*
 \langle /PurchaseOrder

Date \rightarrow \langle Date>
Day? Month Year
 \langle /Date

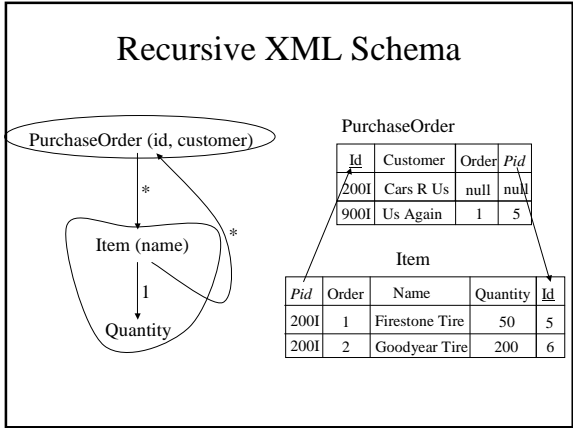
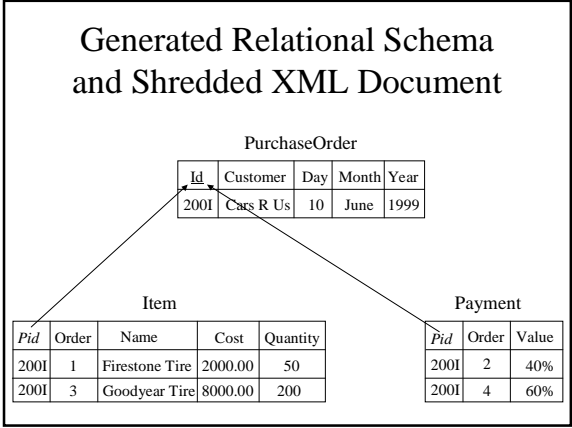
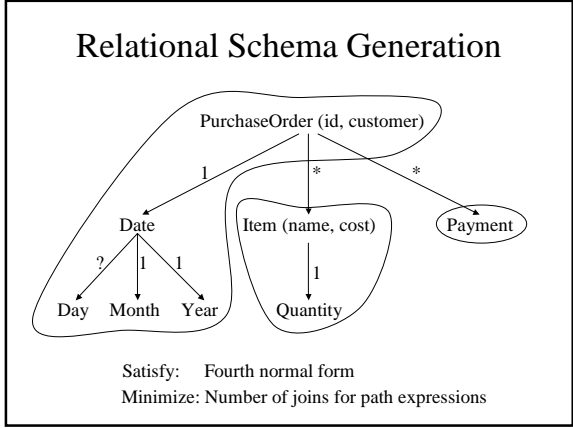
Day \rightarrow \langle Day> {integer} \langle /Day

Month \rightarrow \langle Month> {string} \langle /Month

Year \rightarrow \langle Year> {integer} \langle /Year

Item \rightarrow \langle Item name={string} cost={float}>
Quantity
 \langle /Item

... and so on



- ### Relational Schema Generation and XML Document Shredding (Completeness and Optimality)
- Any XML Schema X can be mapped to a relational schema R , and ...
 - Any XML document XD conforming to X can be converted to tuples in R
 - Further, XD can be recovered from the tuples in R
 - Also minimizes the number of joins for path expressions (given fourth normal form)

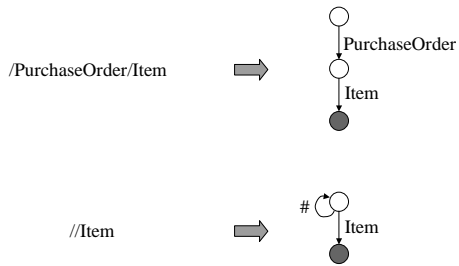
- ### Outline
- Motivation & High-level Solution
 - Background (Relations, XML)
 - Storing and Querying XML Documents
 - Relational Schema Design and XML Storage
 - Query Mapping and Result Construction
 - Publishing Relational Data as XML Documents
 - Conclusion

XML Query

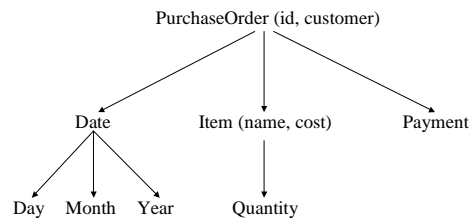
Find all the items bought by "Cars R Us" in 1999

For \$po in /PurchaseOrder
Where \$po/@customer = "Cars R Us" and \$po/date/year = 1999
Return \$po/Item

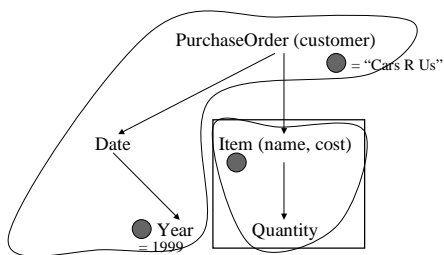
Path Expression Automata (Moore Machines)



XML Schema Automaton (Mealy Machine)



Intersected Automaton



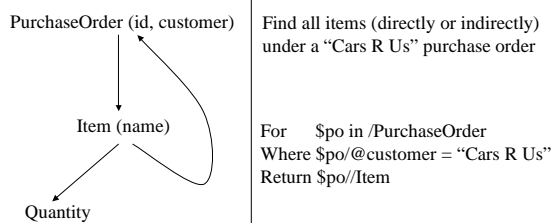
Generated SQL Query

```

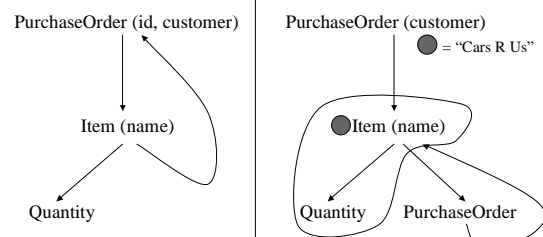
Select i.name, i.cost, i.quantity
From PurchaseOrder p, Item i
Where p.customer = "Cars R Us" and
      p.year = 1999 and
      p.id = i.pid
  
```

} Predicates
 } Join condition

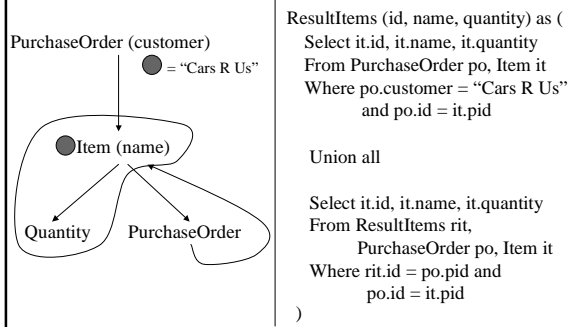
Recursive XML Query



Recursive Automata Intersection



Recursive SQL Generation



SQL Generation for Path Expressions (Completeness)

- (Almost) all path expressions can be translated to SQL
- SQL does not support
 - Nested recursion
 - Meta-data querying
- Meta-data query capability provided in the XML translation layer

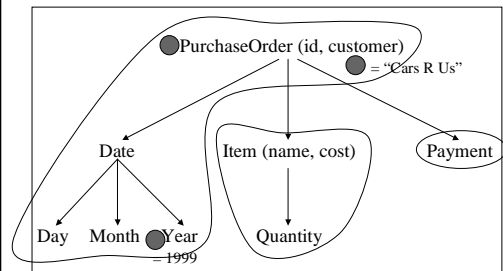
Constructing XML Results

```
<item name = "Firestone Tire" cost="2000.00">
  <quantity> 50 </quantity>
</item>
<item name="Goodyear Tire" cost="8000.00">
  <quantity> 200 </quantity>
</item>
```



("Firestone Tire", 2000.00, 50)
 ("Goodyear Tire", 8000.00, 200)

Complex XML Construction



Outline

- Motivation & High-level Solution
- Background (Relations, XML)
- Storing and Querying XML Documents
- Publishing Relational Data as XML Documents
- Conclusion

Relational Schema and Data

PurchaseOrder

<u>Id</u>	Customer	Day	Month	Year
2001	Cars R Us	10	June	1999

Item

<i>Pid</i>	Name	Quantity	Cost
2001	Firestone Tire	50	2000.00
2001	Goodyear Tire	200	8000.00

Payment

<i>Pid</i>	Installment	Percentage
2001	1	40%
2001	2	60%

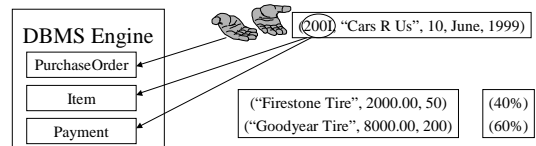
XML Document

```

<PurchaseOrder id="2001" customer="Cars R Us">
  <Date>
    <Day> 10 </Day>
    <Month> June </Month>
    <Year> 1999 </Year>
  </Date>
  <Item name="Firestone Tire" cost="2000.00">
    <Quantity> 50 </Quantity>
  </Item>
  <Item name="Goodyear Tire" cost="8000.00">
    <Quantity> 200 </Quantity>
  </Item>
  <Payment> 40% </Payment>
  <Payment> 60% </Payment>
</PurchaseOrder>
  
```

Naïve Approach

- Issue many SQL queries that mirror the structure of the XML document to be constructed
- Tag nested structures as they are produced



Problem 1: Too many SQL queries
 Problem 2: Fixed (nested loop) join strategy

Relations to XML: Issues

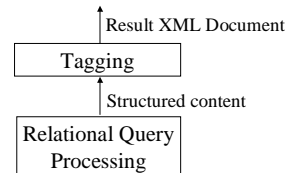
[Shanmugasundaram et. al., VLDB'00]

- Two main differences:
 - Ordered nested structures
 - Self-describing tags
- Space of alternatives:

	Early Tagging	Late Tagging
Early Structuring		✓
Late Structuring	■	

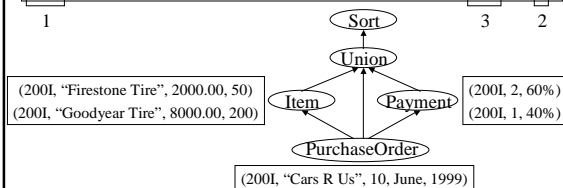
Late Tagging, Early Structuring

- *Structured* XML document content produced
 - In document order
 - “Sorted Outer Union” approach
- Tagger just adds tags
 - In constant space



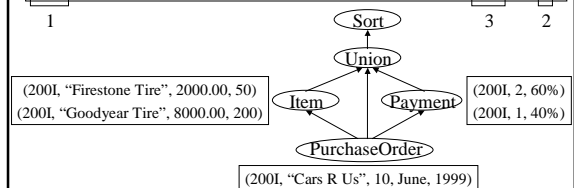
Sorted Outer Union Approach

(2001	null	, null, null, null,	null	, null, null,	1,	40%	2)
(2001	"Cars R Us", 10, June, 1999,	null	, null, null,	null, null,	0)		
(2001	null	, null, null, null,	null	, null, null,	2,	60%	2)
(2001	null	, null, null, null,	"Firestone Tire", 2000.00, 50	null, null,	1)		
(2001	null	, null, null, null,	"Goodyear Tire", 8000.00, 200	null, null,	1)		



Sorted Outer Union Approach

(2001	"Cars R Us", 10, June, 1999,	null	, null, null,	null, null,	0)		
(2001	null	, null, null, null,	"Firestone Tire", 2000.00, 50	null, null,	1)		
(2001	null	, null, null, null,	"Goodyear Tire", 8000.00, 200	null, null,	1)		
(2001	null	, null, null, null,	null	, null, null,	1,	40%	2)
(2001	null	, null, null, null,	null	, null, null,	2,	60%	2)



XML Document Construction (Completeness and Performance)

- Any nested XML document can be constructed using “sorted outer union” approach
- 9x faster than previous approaches [VLDB’00]
 - 10 MB of data
 - 17 seconds for sorted outer union approach
 - 160 seconds for “naïve XML application developer” approach

Outline

- Motivation & High-level Solution
- Background (Relations, XML)
- Storing and Querying XML Documents
- Publishing Relational Data as XML Documents
- Conclusion

Conclusion

- XML has emerged as the Internet data format
- But relational database systems will continue to be used for data management tasks
- Internet application developers currently have to explicitly bridge this “data model gap”
- Can we design a system that automatically bridges this gap for application developers?

For \$cust in /Customer
Where \$cust/name = “Jack”
Return \$cust

```

// First prepare all the SQL statements to be executed and create cursors for them
Exec SQL Prepare CustStat From "select cust_id, cust_name from Customer cust where cust_name = 'Jack'"
Exec SQL Declare CustCursor Cursor For CustStat
Exec SQL Prepare AcctStat From "select acct_id, acct_account from Account acct where acct.custid = ?"
Exec SQL Declare AcctCursor Cursor For AcctStat
Exec SQL Prepare ProductStat From "select product_id, product_name, product_desc from ProductOrder product
where product.custid = ?"
Exec SQL Declare ProductCursor Cursor For ProductStat
Exec SQL Prepare ItemStat From "select item_id, item_desc from Item item where item.prod = ?"
Exec SQL Declare ItemCursor Cursor For ItemStat
Exec SQL Prepare PayStat From "select pay_id, pay_desc from Payment pay where item.prod = ?"
Exec SQL Declare PayCursor Cursor For PayStat
// Now execute SQL statements in nested order of XML document result. Start with customer
XMLResult := ""
Exec SQL Open CustCursor
while (CustCursor has more rows) {
  Exec SQL Fetch CustCursor into custid, custName
  XMLResult := "<customer id=" + custid + " name=" + custName + "></customer>"
  // For each customer, issue sub-query to get account information and add to acctStats
  Exec SQL Open AcctCursor Using custid
  while (AcctCursor has more rows) {
    Exec SQL Fetch AcctCursor into acctid, acctName
    XMLResult := "<account id=" + acctid + " name=" + acctName + "></account>"
  }
  // For each customer, issue sub-query to get purchase order information and add to prodOrders
  Exec SQL Open ProductCursor Using custid
  while (ProductCursor has more rows) {
    Exec SQL Fetch ProductCursor into prodid, product, prodDesc
    XMLResult := "<product id=" + prodid + " name=" + product + " desc=" + prodDesc + "></product>"
    // For each purchase order, issue a sub-query to get item information and add to prodItems
    Exec SQL Open ItemCursor Using prodid
    while (ItemCursor has more rows) {
      Exec SQL Fetch ItemCursor into itemid, itemDesc
      XMLResult := "<item id=" + itemid + " name=" + itemDesc + "></item>"
    }
    // For each purchase order, issue a sub-query to get payment information and add to prodPayments
    Exec SQL Open PayCursor Using prodid
    while (PayCursor has more rows) {
      Exec SQL Fetch PayCursor into payid, payDesc
      XMLResult := "<payment id=" + payid + " name=" + payDesc + "></payment>"
    }
    XMLResult := "</product>"
  }
  // End of looping over all purchase orders associated with a customer
  XMLResult := "</customer>"
  Return XMLResult as one result row; reset XMLResult := ""
} // loop until all customers are tagged and output

```

Conclusion (Contd.)

- Yes! XPERANTO is the first such system
- Allows users to ...
 - Store and query XML documents using a relational database system
 - Publish existing relational data as XML documents
- ... using a high-level XML query language
- Also provides a dramatic improvement in performance

Industry Impact

- Sorted outer union approach is used in the DB2 XML Extender product (beta version)
- XPERANTO is now an IBM initiative

Relational Database System Vendors

- IBM, Microsoft, Oracle, Informix, ...
 - SQL extensions for XML
- XML Translation Layer
 - “Pure XML” philosophy ... provides high-level XML query interface
 - SQL extensions for XML, while better than writing applications, is still low-level
 - More powerful than XML-extended SQL
 - SQL just not designed with nifty XML features in mind