

Foundations of XML Data Manipulation

Giorgio Ghelli

Query Languages for SSD and XML

Some names

- UnQL: BunDavHil96
- Lorel: AbiQuaMcH97
- XSLT, XPath, XQuery: google:w3c xslt, w3c xpath, w3c xquery
- XML-QL: XMLQL, xquery99
- XDuce: HosPie03
- TQL: CarGhe03
- YATL, Strudel...

Path expressions

```
<bib>
  <book year="1995">
    <author> <first>Serge</first> <last>Abiteboul</last> </author>
    <author> <first>Richard</first> <last>Hull</last> </author>
    <author> <first>Victor</first> <last>Vianu</last> </author>
    <publisher>Addison</publisher>
    <price>60</price>
  </book>
  <book year="1993">
    <title>Formal Semantics</title>
    <author> <first>Glynn</first> <last>Winskel</last> </author>
    <publisher>MIT Press</publisher>
    <price>42</price>
  </book>
</bib>
```

Path expressions

- *Document* /bib/book/author/first:
 - <first>Serge</first>, <first>Richard</first>, <first>Victor</first>, <first>Glynn</first>
- Semantics:
 - All nodes you find starting from *Document* and walking down a /bib/book/author/first path
 - More generally: $[[p]] = \{ \langle n, m \rangle \mid m \in n/p \}$
- The interesting case:
 - Data is a graph
 - Paths is a regular expression

A formal definition

- Assume a graph $G=(N,E)$ with $E \subseteq N \times A \times N$
- A word w in A^* determines a relation $|w|$ on $N \times N$:
 - $n | \epsilon | n$
 - $(n, a, n') \in E \Rightarrow n | a | n'$
 - $n | w | n'$ and $n' | w' | n'' \Rightarrow n | w.w' | n''$
- A language $L \subseteq A^*$ determines a relation:
 - $n | L | n' \Leftrightarrow n | w | n'$ for some $w \in L$
- A regexp r determines a relation:
 - $n | r | n' \Leftrightarrow n | \text{Lang}(r) | n'$
- We allow regexp on labels as well

Regular path expressions

- "book"."[T|t]title":
 - {book.title, book.Title}
- book.(Title|title):
 - {book.title, book.Title}
- ".*".title:
 - {book.title, author.title, name.title,...} i.e. {*_title}
- book.refs*.title:
 - {book.title, book.refs.title, book.refs.refs.title,...}
- (".*")*.name:
 - {name, _name, __name, ___name,...}

Syntax for path expressions

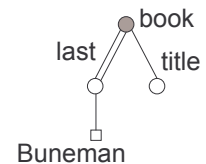
- $a ::= \varepsilon \mid \text{letter} \mid . \mid a a \mid [a]a \mid a^*$
- $r ::= \varepsilon \mid a \mid _ \mid r.r \mid r|r \mid e^*$
- Abbreviations:
 - $r? = r|\varepsilon$
 - $r+ = r.r^*$

XPath syntax

- $p, p' ::= / p \mid p/p' \mid \text{axis}::\text{nodetest}$
- $\text{nodetest} ::= * \mid \text{tag} \mid \text{text}() \mid \text{element}^* \mid \text{element}(\text{tag}) \mid \text{attribute}^* \mid \dots$
- $\text{axis} ::= \text{child} \mid \text{parent} \mid \text{d-o-s} \mid \text{a-o-s}$
 - descendant | ancestor
 - following | preceding
 - following-sibling | preceding-sibling
- $p/\text{nodetest} = p/\text{child}::\text{nodetest}$
- $p//q = p/\text{d-o-s}::*/q$
- $.. = \text{parent}::\text{node}()$

Tree Patterns

- Titles of books:
 - $\$doc // \text{book}/\text{title}$
- Titles of books by Buneman:
 - $\$doc // \text{book}[//\text{last}/\text{text}() = \text{'Buneman'}]/\text{title}$



Tree Patterns (Twigs)

- $\text{book}[\text{author}]/\text{title}$
- $\text{book}/\text{author}/../\text{title}$



- $\text{author}/../\text{book}/\text{title}$



- Up pointers are slightly more expressive than twigs

Are path expressions a query language?

The structure of XML query languages

- FROM BindingExpression(X): generate a set of bindings for X
- WHERE Condition(X): filter some bindings out
- SELECT Result(X): evaluate Result(X) once for each binding, and find a way to merge the results
- Path expressions come handy in the FROM clause

MicroXQuery

- Reference: Colazzo et al., Types for Path Correctness of XML Queries, ICALP'04.
for \$b in \$doc /bib/book,
 \$a in \$b /author
where \$b /@year > 2000
return <libro> \$b/title, \$a </libro>

MicroXQuery

- for \$b in \$doc /bib/book
 let \$a = \$b /author
 where \$b /@year > 2000
 return <libro> \$b/title, \$a </libro>
- for \$b in \$doc /bib/book
 where \$b /@year > 2000
 return <libro> \$b/title, \$b /author </libro>

Tree Patterns

- Titles of books by Buneman:
for \$x in \$doc //book
where \$x/authors/last/text() = 'Buneman'
return \$x/title
- Same as:
\$doc //book[last/text() = 'Buneman']/title

Trees to relations to trees

```
let $authors = $doc /bib/book/author
for $a in distinct($authors)
return <booksByAuth> $a,
    for $bb in $doc /bib/book
    where $a isin $bb /author
    return $bb/title
</booksByAuth>
```

Result

- From
 - <bib>
 - (<book> (<author></author>)*
 - </book>)*
 - </bib>
- To
 - (<booksByAuth>
 - <author>...</author>
 - (<title>...</title>)*
 - </booksByAuth>)*

Yet another syntax

```
BBAS[
let $authors = $doc /bib/book/author
for $a in distinct($authors)
return booksByAuth[
  $a,
  for $bb in $doc /bib/book
  where $a isin $bb /author
  return $bb/title
]
]
```

Type

- Type:
BBAS[booksbyaut[author[String],
title[String]*]*]
- If:
\$doc: bib[book [title[String],
author[String]*
]*
]

Tree manipulation

- Micro XQuery is able to do “nested relations”, i.e. trees with fixed depth
- It is unable to produce arbitrarily deeply nested trees:
 - Structural recursion
- It is unable to reverse a graph:
 - Skolem functions

Structural recursion

- $f(v) \Rightarrow v$
- $f(\text{author}[x]) \Rightarrow \text{autore}[f(x)]$
- $f(l[x]) \Rightarrow$ if $l = \text{year}$ then 0 else $l[f(x)]$
- $f(0) \Rightarrow 0$
- $f(x,y) \Rightarrow f(x),f(y)$

Structural recursion

- f : collect authors, if any
- $f(\text{author}[x],y) \Rightarrow \text{authors}[\text{author}[f(x)],$
 $g(y)], h(y)$
- $g(\text{author}[x]) \Rightarrow \text{author}[f(x)]$
- (else) $g(l[x]) \Rightarrow 0$
- $h(\text{author}[x]) \Rightarrow 0$
- (else) $h(l[x]) \Rightarrow l[f(x)]$
- Implicit: $f(a[x],y) \Rightarrow f(a[x]),f(y); f(a[x]) \Rightarrow a[f(x)]$

Structural recursion vs. query languages

- Paths go down
- for – where – selects iterates horizontally
- Structural recursion does both
- Supported by:
 - XQuery (SR only)
 - XSL (SR only)
 - XQuery (FWR + recursion)

XSL

```
<xsl:stylesheet version="1.0" ...>
  <xsl:template match="ul[parent::ul]">
    <li>
      <ul>
        <xsl:apply-templates select="@*|node()"/>
      </ul>
    </li>
  </xsl:template>
  <xsl:template match="@*|node()">
    <xsl:copy>
      <xsl:apply-templates select="@*|node()"/>
    </xsl:copy>
  </xsl:template>
</xsl:stylesheet>
```

Skolem functions

- new(a,b,c,...): a (new) node
- n-lab-n: an edge
- for \$x-\$lab-\$y in Edges
- return new(node,\$y)-\$lab-new(node,\$x)
- for \$x-\$lab-\$y in Edges
- return new(node,\$x)-\$lab-new(node,\$y),
- new(node,\$y)-home-new(root)

Logical languages: TQL

– Matching through ambient-logic formulas:

```
FROM $db |= .paper[ (.author[$X] or .autore[$X])
                    and not .editor[$X]
                  ]
SELECT author[$X]
```

Returns:

author[Cardelli] | author[Gordon] | author[Ghelli]

Find All Keys

from \$Bib

|= bib[!book[.\$k[T]]

And foreach \$X.

Not (.book.\$k[\$X] | .book.\$k[\$X])

]

select key[\$k]

The Logic

- $F \vdash_{\sigma} 0$ iff $F=0$
- $F \vdash_{\sigma} A \mid B$ iff $\exists F', F''. F = F' \mid F'', F' \vdash_{\sigma} A, F'' \vdash_{\sigma} B$
- $F \vdash_{\sigma} \eta[A]$ iff $F = \sigma(\eta)[F'], F' \vdash_{\sigma} A$
- $F \vdash_{\sigma} A \wedge B$ iff $F \vdash_{\sigma} A$ and $F \vdash_{\sigma} B$
- $F \vdash_{\sigma} \neg A$ iff not $(F \vdash_{\sigma} A)$
- $F \vdash_{\sigma} \eta = \eta'$ iff $\sigma(\eta) = \sigma(\eta')$
- $F \vdash_{\sigma} \exists x.A$ iff $\exists n. F \vdash_{\sigma\{n/x\}} A$
- $F \vdash_{\sigma} \exists X.A$ iff $\exists F'. F \vdash_{\sigma\{F'/X\}} A$
- $F \vdash_{\sigma} \mu\xi.A$ iff $F \vdash_{\sigma} A\{\mu\xi.A/\xi\}$ (is circular...)
- De Morgan duals: $\mid, \eta[\Rightarrow A], F, \vee, \neq, \forall, \vee\xi.A, !\eta[A]$

Logical languages: monadic datalog

- doc //book/author X:
 - result(Y) :- desc(doc,X), name(X,book), child(X,Y), name(Y,author)

The full hybrid mu-calculus

- $A ::= i \mid \text{lab} \mid \langle s \rangle A \mid A \wedge A \mid \neg A \mid \mu \xi. A \mid \xi$
- $E, L, w \vdash_{\sigma} i$ iff $L(i) = \{w\}$
- $E, L, w \vdash_{\sigma} \text{lab}$ iff $w \in L(\text{lab})$ ($L(w) = \text{lab}$)
- $E, L, w \vdash_{\sigma} \langle s \rangle A$ iff $\exists w'. w.s.w'$ in E and $E, L, w' \vdash_{\sigma} A$
- $E, L, w \vdash_{\sigma} \langle s \rangle A$ iff $\exists w'. w'.s.w$ in E and $E, L, w' \vdash_{\sigma} A$
- $E, L, w \vdash_{\sigma} A \wedge B$ iff $E, L, w \vdash_{\sigma} A$ and $E, L, w \vdash_{\sigma} B$
- $E, L, w \vdash_{\sigma} \neg A$ iff not $(E, L, w \vdash_{\sigma} A)$
- $E, L, w \vdash_{\sigma} \mu \xi. A$ iff $E, L, w \vdash_{\sigma} A\{\mu \xi. A/\xi\}$ (is circular...)
- Over XML:
 - Just two steps $s: \downarrow/\uparrow$ (firstchild) and \rightarrow/\leftarrow (nextsibling)
 - E is a finite tree
 - $L(\text{lab})$ is a partition of the nodes
- De Morgan duals: $[s]A, \vee, \nu \xi. A$

Encoding axes

- $\langle \text{child} \rangle A = \langle \downarrow \rangle (\mu \xi (A \vee \langle \rightarrow \rangle \xi))$
- $\langle \text{parent} \rangle A = \mu \xi (\langle \uparrow \rangle A \vee \langle \leftarrow \rangle \xi)$
- $\langle \text{desc} \rangle A = \langle \text{child} \rangle (\mu \xi (A \vee \langle \text{child} \rangle \xi))$
- $\langle \text{ancestor} \rangle$
- $\langle \text{following-sibling} \rangle$
- $[\text{everywhere}] A = \nu \xi (A \wedge [\downarrow] \xi \wedge [\rightarrow] \xi)$
- $\langle \text{somewhere} \rangle A = \mu \xi (A \vee \langle \downarrow \rangle \xi \vee \langle \rightarrow \rangle \xi)$

Logical languages: modal mu calculus

- $\langle m, n \rangle$ in $[[p]]$ iff $E, L, i \rightarrow \{n\}, m \models \langle \langle p \rangle \rangle$
- $//\text{book}/\text{author}$:
 $\langle \text{desc} \rangle (\text{book} \wedge \langle \text{child} \rangle (\text{author} \wedge i))$
- $\langle m, n \rangle$ in $//\text{book}[\text{title}]/\text{author}$
 $i \rightarrow n, m \models$
 $\langle \text{desc} \rangle (\text{book} \text{ and } (\langle \text{child} \rangle \text{title})$
 $\text{and } (\langle \text{child} \rangle \text{author and } i))$
- $\langle m, n \rangle$ in $//\text{book}[\text{not title}]/\text{author}$
 $i \rightarrow n, m \models$
 $\langle \text{desc} \rangle (\text{book} \text{ and } (\text{not } \langle \text{child} \rangle \text{title})$
 $\text{and } (\langle \text{child} \rangle \text{author and } i))$

Logical languages: MSO

- MTran: transformation language based on MSO [HinabaHosoya..]
- MSO: FO plus set quantification
- Able to express all regular tree query

MTran formulas

- $//\text{img}: x$ in $\langle \text{img} \rangle$
- $//^*[\text{date}]: \text{ex1 } y: x/y$ and y in $\langle \text{date} \rangle$
- $//a\{x\}/b\{y\}: x$ in $\langle a \rangle$ and x/y and y in $\langle b \rangle$

MTran transformation

- Add an $\langle \text{li} \rangle$ around to any $\langle \text{ul} \rangle$ whose parent is an $\langle \text{ul} \rangle$:
– $\{ \text{visit } x :: \langle \text{ul} \rangle /x \ \& \ x \text{ in } \langle \text{ul} \rangle :: \text{li}[x] \}$
- *visit* copies all unmatched nodes, *gather* does not:
– $\langle \text{ul} \rangle \{ \text{gather } x :: x \text{ in } \langle a \rangle :: \text{li}[x] \} / \langle \text{ul} \rangle$

Nesting

- {gather b :: b in <book> ::
 {gather a :: b/a & a in <author> ::
 book-author[b,a] } }

Evaluations

- Compile the formula to a tree automaton
 - Non-elementary in the worst case
 - MONA often works well
- Evaluate the automaton
 - Linear algorithm (non trivial, since the query is n-ary)

Tree Automata

- Morally, the same expressive power (over trees) as:
 - Monadic datalog
 - Tree patterns
 - MSO
 - Mu-calculus
 - XDuce type system
 - ...
- Things are not so simple...

Readings

- Xquery99: comparison of XML-QL, YATL, Lorel, XQL
- BonCer00: comparison of Lorel, XML-QL, XML-GL, XSL, XQL
- ColGheAl06: MicroXQuery
- www.w3.org/TR/xquery|xpath|xsl: XQuery, XPath, XSLT
- BunDavHil96: UnQL
- AbiQuaMcH97: Lorel
- HosPie03: XDuce
- klarlundSchweintick: XPath, XQuery, XSLT
- CarGhe03: TQL
- ToCl-lics0203: Logics, Automata