

# XML Data Management

## Part 12: Core XPath

Werner Nutt

Faculty of Computer Science  
Master of Science in Computer Science

A.Y. 2012/2013



FREIE UNIVERSITÄT BOZEN

LIBERA UNIVERSITÀ DI BOLZANO

FREE UNIVERSITY OF BOZEN · BOLZANO

# A Logic View of Databases

A database has two parts: **schema** and **instance**

The schema describes *how data is organized*:

- relations with their names and with number, names and types of attributes
- Example: `Person(name, gender, age)`, `HasChild(parent, child)`

The instance contains the *actual data*:

- for every relation, there is a set of atoms complying with the schema,
- Example: `{Person(fred, 'male', 25), Person(mary, 'female', 60), Person(clara, 'female', 3), Person(paul, 'male', 1), HasChild(mary, fred), HasChild(fred, clara), HasChild(fred, paul)}`

*Often, we ignore types and sometimes, we ignore also the attribute names*

# First-Order Queries

## Definition

A **first-order query** has the form

$$Q = \{(x_1, \dots, x_n) \mid \phi\}$$

where

- $\phi$  is a predicate logic formula
- $x_1, \dots, x_n$  are the free variables of  $\phi$

We say that

- $\phi$  is the **body** of the query,
- $x_1, \dots, x_n$  are the **output variables**, and
- $n$  is the **arity** of the query.

## Example Queries

Let's try to express the following queries as first-order queries:

- Who are the male persons?
- Who are the grandmothers?
- Who are the male persons without children?

*What do these queries look like in SQL?*



# XML Documents as Relational Instances

XML docs are rooted labeled trees, a special case of labeled directed graphs.

**Directed graphs** (digraphs) consist of

- nodes
- directed edges (described by a binary relation `child(·, ·)`)

**Rooted trees** are digraphs with

- one source (no incoming edge), called “root”
- an arbitrary number of sinks, called “leaves”
- no cycles.

## XML Documents as Relational Instances (cntd)

We assume there is a set of **labels**  $\Sigma$  (labels model the element tags)

A **labeled tree** has exactly one label on each node

The set of all trees with labels from  $\Sigma$  is denoted as  $T_\Sigma$

We express that a node carries label  $a \in \Sigma$  with the unary relation  $lab_a(\cdot)$ .

We identify all labeled rooted trees  $t$   
with instances of the schema with the relations

- $child(\cdot, \cdot)$
- $lab_a(\cdot), a \in \Sigma$

If  $t$  is such a tree, then  $nod(t)$  is the set of all nodes of  $t$

*We are graciously ignoring strings and other values.*

*We could model them with unary relations “ $text_{xyz}(\cdot)$ ”*

# XPath Queries as First-order Queries

We want to query our movies documents known from the coursework

- Select all movies  $x$ : `(//movie)`
- Select all actors  $x$  of “Spider Man”  
`(//movie[title/text()='SM']/actor)`
- Select all pairs of actors  $x$  appearing in movie  $y$

*How many variables do we need to write these queries?*



# XPath Queries over Trees: Exercises

Write queries asking for

- all movies starring Kirsten Dunst
- all movies starring Kirsten Dunst and William Dafoe
- all movies with Kirsten Dunst, but not William Dafoe





# First-order Queries over Trees

## Syntax of Tree Formulas

For  $a \in \Sigma$  and  $x, y \in \text{Vars}$ :

$$\phi ::= lab_a(x) \mid child^*(x, y) \mid next\_sibling^*(x, y) \mid \neg\phi \mid \exists x\phi \mid \phi_1 \wedge \phi_2$$

## Transitive Closure

We have to add `child*` to the instances if we want to talk about descendants in FO, and similarly `next_sibling*` in order to talk about horizontal recursion.

## Tree Queries

As before, queries over trees are expressions of the form

$$\{(x_1, \dots, x_n) \mid \phi\},$$

# Core XPath 1.0

Introduced by Gottlob & Koch in [PODS'01, JACM'03]

- “logical core” of XPath 1.0, used to study theoretical properties
- many simplifications
  - removes arithmetic
  - removes functions on data content (e.g., on strings)
  - leaves only the navigational core

# Example Queries in Core XPath 1.0

## Select all movies

XPath short: `//movie`  
 XPath long: `child* :: movie`  
 FO logic:  $child^*(root, x) \wedge lab_{movie}(x)$

## Select all actors acting together with JDepp in a movie

XPath short: `//movie[actor/text() = 'JDepp']/actor[not[text() = 'JDepp']]`  
 XPath long: `child* :: movie[child :: actor[text() = 'JDepp']] / child :: actor[not[text() = 'JDepp']]`  
 FO logic:  $\exists y_1 (child^*(root, y_1) \wedge lab_{movie}(y_1)) \wedge \exists y_2 (child(y_1, y_2) \wedge lab_{actor}(y_2) \wedge \mathbf{text}_{JDepp}(y_2)) \wedge child(y_1, x) \wedge lab_{actor}(x) \wedge \neg \mathbf{text}_{JDepp}(x)$

# Core XPath 1.0

Syntax where  $a \in \Sigma$

filter  $f ::= * \mid a \mid \text{not } [f] \mid f[p]$

axis  $r ::= \text{child} \mid \text{next\_sibling} \mid r^{-1} \mid r^* \mid \text{self} \mid \dots$

paths  $p ::= r \mid p/p' \mid p \cup p' \mid /p$

Semantics for trees  $t \in T_\Sigma$

$$\text{eval}^t(f) \subseteq \text{nod}(t)$$

$$\text{eval}^t(r) \subseteq \text{nod}(t)^2$$

$$\text{eval}^t(p) \subseteq \text{nod}(t)^2$$

# Core XPath 1.0 Semantics

$$eval^t(*) = nod(t)$$

$$eval^t(a) = lab_a^t \quad (\text{extension of } lab_a \text{ in } t)$$

$$eval^t(\text{not } f) = nod(t) \setminus eval^t(f)$$

$$eval^t(f[p]) = \{n \in eval^t(f) \mid \exists n'.(n, n') \in eval^t(p)\}$$

$$eval^t(\text{child}) = \text{child}^t \quad (\text{extension of } \text{child} \text{ in } t)$$

$$eval^t(\text{next\_sibling}) = \text{next\_sibling}^t$$

$$eval^t(r^{-1}) = eval^t(r)^{-1}$$

$$eval^t(r^*) = eval^t(r)^*$$

$$eval^t(\text{self}) = \{(n, n) \mid n \in nod(t)\}$$

$$eval^t(r :: f) = \{(n, n') \in eval^t(r) \mid n' \in eval^t(f)\}$$

$$eval^t(p/p') = eval^t(p) \circ eval^t(p') \quad (\text{composition of relations})$$

$$eval^t(p \cup p') = eval^t(p) \cup eval^t(p')$$

$$eval^t(/p) = \{(root, n) \in nod(t)^2 \mid (root, n) \in eval^t(p)\}$$

# Exercises: Expressiveness of Core XPath 1.0

Can one define the following queries in Core XPath 1.0?

- All nodes reachable from the root over a path with labels in  $a^*b$ ?
- All nodes that are reachable from the root over a path with labels in  $(aa)^*$ ?
- Can one define the same query in Datalog for trees?
- And what about Datalog where all query predicates are unary (monadic datalog)?

# Translation to FO Logic

## Proposition

Every expression of XPath 1.0 can be translated in linear time to an FO formula with 2 free variables, that define the same binary query.

$$\llbracket * \rrbracket_x = \text{true}$$

$$\llbracket a \rrbracket_x = \text{lab}_a(x)$$

$$\llbracket \text{not } [f] \rrbracket_x = \neg \llbracket [f] \rrbracket_x$$

$$\llbracket [f][p] \rrbracket_x = \llbracket [f] \rrbracket_x \wedge \exists y \llbracket [p] \rrbracket_{x,y}$$

$$\llbracket \text{child} \rrbracket_{x,y} = \text{child}(x, y)$$

$$\llbracket \text{next\_sibling} \rrbracket_{x,y} = \text{next\_sibling}(x, y)$$

$$\llbracket [r^{-1}] \rrbracket_{x,y} = r(y, x)$$

$$\llbracket [r^*] \rrbracket_{x,y} = r^*(x, y)$$

$$\llbracket [(r^*)^{-1}] \rrbracket_{x,y} = r^*(y, x)$$

$$\llbracket \text{self} \rrbracket_{x,y} = (x = y)$$

$$\llbracket [r :: f] \rrbracket_{x,y} = \llbracket [r] \rrbracket_{x,y} \llbracket [f] \rrbracket_y$$

$$\llbracket [p \cup p'] \rrbracket_{x,y} = \llbracket [p] \rrbracket_{x,y} \vee \llbracket [p'] \rrbracket_{x,y}$$

$$\llbracket [p/p'] \rrbracket_{x,y} = \exists z (\llbracket [p] \rrbracket_{x,z} \wedge \llbracket [p'] \rrbracket_{z,y})$$

$$\llbracket [/p'] \rrbracket_{x,y} = (x = \text{root} \wedge \llbracket [p] \rrbracket_{z,y})$$

# Query Answering for Core XPath 1.0

For a start set  $S \subseteq \text{nod}(t)$ , let  $\text{eval}_S^t(p) = \{n' \mid n \in S, (n, n') \in \text{eval}^t(p)\}$

## Theorem (Gottlob & Koch (TODS'05))

For all expressions  $p$  of Core XPath 1.0, trees  $t$ , and start sets  $S$ , one can compute the monadic query  $\text{eval}_S^t(p)$  in time  $O(|t| \cdot |p|)$ .

Idea: Uses an algebra where one navigates from  $S$  through an operator tree corresponding to  $p$ .

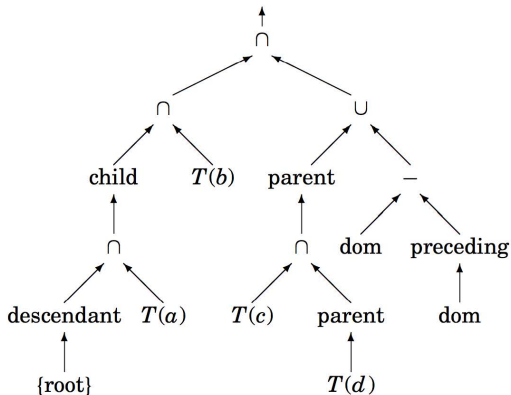


# Operator Tree of Query [Gottlob & Koch (TODS'05)]

*Example 10.3.* The Core XPath query

`/descendant::a/child::b[child::c/child::d or not(following::*)]`

is evaluated as specified by the query tree



# Complexity of Full XPath 1.0

Full XPath 1.0 is more difficult because of

- equalities (join on labels)
- data manipulation (string and arithmetic functions)
- context information (context node, position in current node set, size of node set).

## Theorem (Gottlob, Koch & Pichler, 2005)

There is an algorithm that evaluates an expression  $p$  on a document  $t$  with

- time complexity  $O(|t|^4 \cdot |p|^2)$
- space complexity  $O(|t|^2 \cdot |p|^2)$

sd



# Expressivity of Core XPath 1.0 (Marx & de Rijke)

## Theorem

For every FO tree query with 2 variables  $Q = \{x \mid \phi\}$ , there exists a Core XPath expression filter expression  $f$ , such that  $Q$  and  $f$  return the same answers over all rooted labeled trees, and conversely.

In other words, Core XPath 1.0 has the same expressivity as First-Order Logic over trees with 2 variables.

# From XPath 1.0 to XPath 2.0

XPath 2.0 extends XPath 1.0 essentially by

- stronger typing (types as in XML schema)
- sequence processing functions (e.g., remove, insert, index-of)
- explicit quantification with variables as in XQuery  
(using some  $\$x$  in  $p$  satisfies,  $e$ , every  $\$x$  in  $p$  satisfies,  $e$ )
- iteration as in XQuery  
(using for  $\$x$  in  $p$  return  $e$ )
- intersection and difference of paths

As for XPath 1.0, one has defined Core XPath 2.0, which adds quantification, iteration, iteration and path intersection and difference to XPath 1.0.

# Expressivity of Core XPath 2.0

## Theorem (Marx & ten Kaate)

- The evaluation problem for Core XPath 2.0 is PSpace-complete
- Core XPath 2.0 has the same expressivity as full first-order logice over trees.