

Foundations of Data and Knowledge Bases

XPath

Joachim Niehren

Links: Linking Dynamic Data
Inria Lille

September 7, 2022

Outline

- 1 Datalog Queries for Unranked Trees
- 2 First-order Queries for Unranked Trees
- 3 Core XPath 1.0
- 4 Core XPath 2.0
- 5 Beyond the Core

Repetition

Unranked Trees

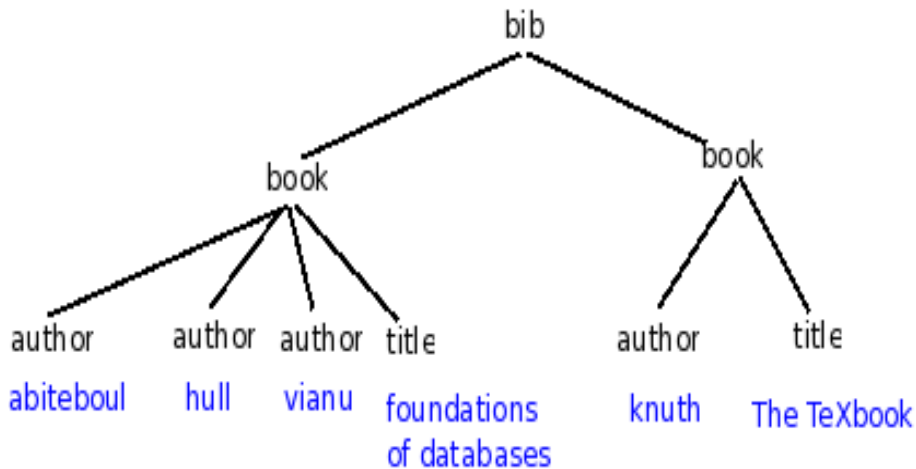
For a finite set Δ of labels (XML tags):

$$t \in T_{\Delta} ::= a(t_1, \dots, t_n) \quad n \geq 0, a \in \Delta$$

Simplifications

- textual content of nodes is not part of the tree (external annotation).
- no attributes

Example



Node Selection Queries

n-ary Queries Q

functions mapping trees to sets of n-tuples of nodes

$$\forall t \in T_{\Sigma} : Q(t) \subseteq \text{nod}(t)^n$$

Examples

monadic select all book authors

binary select all pairs of authors and titles in some book

n-ary select n-tuples of nodes in trees

Boolean where $n = 0$. Select empty tuple $()$ in documents that satisfy the schema, and nothing otherwise.

Basic Queries

monadic **leaf** selects all leaf nodes.

binary **child** selects all pairs of father-child nodes.

Monadic Datalog for Unranked Trees

Signature

Constants $\text{Consts} = \{\text{root}\}$
Relation symbols $\text{Rels} = \{\text{next_sib}(\cdot, \cdot), \text{first_child}(\cdot, \cdot),$
 $\text{last_child}(\cdot, \cdot), \text{leaf}(\cdot)\}$

Syntax (all variable $q \in \text{RelVars}$ have arity 1)

terms $t ::= x \mid a$ ($x \in \text{Vars}$ and $a \in \text{Consts}$)
literals $l ::= q(t) \mid r(t_1, \dots, t_{\text{ar}(q)})$ ($q \in \text{RelVars}$ and $r \in \text{Rels}$)
Horn clauses $q(t) :- l_1, \dots, l_n.$
Programs P finite set of Horn clauses

Example

Each second line of HTML tables: $table(tr, \dots, tr)$

$q_{even}(x) :- \text{first_child}(\text{root}, x).$

$q_{odd}(x) :- \text{next_sib}(y, x), q_{even}(y).$

$q_{even}(x) :- \text{next_sib}(y, x), q_{odd}(y).$

Normal Forms

Clause Formats

$q(\text{root}).$	test for root
$q(x) :- \text{leaf}(x)$	test for leaf
$q(x) :- \text{first_child}(x, y), q'(y).$	go down left most
$q(x) :- \text{first_child}(y, x), q'(y).$	go up left most
$q(x) :- \text{next_sib}(x, y), q'(y).$	go right
$q(x) :- \text{next_sib}(y, x), q'(y).$	go left
$q(x) :- \text{last_child}(x, y), q'(y).$	go down right most
$q(x) :- \text{last_child}(y, x), q'(y).$	go up right most
$q(x) :- q_1(x), q_2(x).$	run in parallel

Proposition (Gottlob& Koch PODS'02)

All monadic Datalog programs for unranked trees can be brought in normal form in linear time.

Lemma

Normal forms programs P can be grounded for a tree t in combined linear time $O(|t| \cdot |P|)$.

Efficient Query Answering

Define Monadic Queries

by Datalog program P and a goal predicate q

$$Q_{P,q}(t) = \text{lfp}_P^t(q)$$

Theorem (Gottlob& Koch PODS'02)

Answers of monadic queries $Q_{P,q}(t)$ defined monadic Datalog programs in unranked trees t can be computed in combined linear time $O(|t| \cdot |P|)$.

Proof.

Compute normal form of P , make it ground with respect to t and compute the least solution. □

Outline

- 1 Datalog Queries for Unranked Trees
- 2 First-order Queries for Unranked Trees**
- 3 Core XPath 1.0
- 4 Core XPath 2.0
- 5 Beyond the Core

FO Logic for Unranked Trees

Syntax (where $a \in \Delta$ and $x, y \in \text{Vars}$)

$$\phi ::= \text{lab}_a(x) \mid \text{child}^*(x, y) \mid \text{next_sib}^*(x, y) \mid \neg\phi \mid \exists x\phi \mid \phi_1 \wedge \phi_2$$

Transitive closure is not definable in FO

We have to add child^* to the structures if we want to talk about descendants in FO, and similarly next_sib^* in order to talk about horizontal recursion.

FO Queries

a FO formula ϕ and a sequence of variables $\bar{x} = (x_1, \dots, x_n)$ define an n -ary query $Q_{\phi(\bar{x})}$ such that for all trees $t \in T_{\Delta}$:

$$Q_{\phi(\bar{x})}(t) = \{(\alpha(x_1), \dots, \alpha(x_n)) \mid t, \alpha \models \phi, \alpha : \text{Vars} \rightarrow \text{nod}(t)\}$$

Outline

- 1 Datalog Queries for Unranked Trees
- 2 First-order Queries for Unranked Trees
- 3 Core XPath 1.0**
- 4 Core XPath 2.0
- 5 Beyond the Core

XPath 1.0

W3C Standard (2005)

- language for defining **monadic queries** in unranked trees.
- **no variables!**

Navigational language

- inspired by **modal logics**, i.e., variable free logics definable in FO.
- define **navigation paths** from start nodes to end nodes
- **monadic queries** select end-nodes when starting at the root

Core XPath 1.0 (Gottlob & Koch PODS'01, JACM'03)

- many simplifications
 - ▶ removes arithmetics
 - ▶ removes functions on data content
 - ▶ leaves only the navigational core
- formal semantics

Examples

Select all books in bibliographie

XPath short: `//book`

XPath long: `child*::book`

FO logic: $\text{child}^*(\text{root}, x) \wedge \text{lab}_{\text{book}}(x)$

Select all co-authors of Abiteboul in books of bibliographie

XPath short: `//book[author/text() = 'abitebool']/
author[not(text() = 'abitebool')]`

XPath long: `child*::book[child::author[text() = 'abitebool']]/
child::author[not(text() = 'abitebool')]`

FO logic: $\exists y_1 (\text{child}^*(\text{root}, y_1) \wedge \text{lab}_{\text{book}}(y_1) \\ \wedge \exists y_2 (\text{child}(y_1, y_2) \wedge \text{lab}_{\text{author}}(y_2) \wedge \text{text}_{\text{'abitebool'}}(y_2)) \\ \wedge \text{child}(y_1, x) \wedge \text{lab}_{\text{author}}(x) \wedge \neg \text{text}_{\text{'abitebool'}}(x))$

Core XPath 1.0

Syntax where $a \in \Delta$:

filter $f ::= * \mid a \mid \text{not}[f] \mid f[p]$
axis $r ::= \text{child} \mid \text{next_sib} \mid r^{-1} \mid r^* \mid \text{self} \mid \dots$
paths $p ::= r::f \mid p/p' \mid p \cup p' \mid /p$

Semantics for unranked tree $t \in T_\Delta$:

$$\begin{aligned} \text{eval}^t(f) &\subseteq \text{nod}(t) \\ \text{eval}^t(r) &\subseteq \text{nod}(t)^2 \\ \text{eval}^t(p) &\subseteq \text{nod}(t)^2 \end{aligned}$$

Semantics

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

$$eval^t(a) = lab_a^t$$

$$eval^t(not[f]) = nod(t) \setminus eval^t(f)$$

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

$$eval^t(child) = child^t$$

$$eval^t(next_sib) = next_sib^t$$

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

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

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

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

$$eval^t(p/p') = eval^t(p) \circ eval^t(p')$$

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

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

Examples

first authors of books in bibliography

/child::book/child::author[not[next_sib⁻¹::*]]*

titles of book with at least two authors

/child::book[child::author[next_sib::author]]/child::title*

all author or editor of books

/child::book/((child*::author) ∪ (child*::editor))*

test whether no book has title and booktitle

not[/child::book[[child::title] [child::booktitle]]*

test whether all books of bibliography have an author:

not[child::book[not[child::author]]]*

Exercises: Quantifiers in Core XPath 1.0

- ① Define the following FO monadic query in Core XPath 1.0

$$\text{exists_desc}_a(x) =_{\text{df}} \exists y(\text{child}^+(x, y) \wedge \text{lab}_a(y))$$

- ② Define the following FO monadic query in Core XPath 1.0

$$\text{forall_desc}_a(x) =_{\text{df}} \forall y(\text{child}^+(x, y) \rightarrow \text{lab}_a(y))$$

- ③ Suppose that the set of tree labels is $\Delta = \{a, b\}$. Define a Core XPath 1.0 filter satisfied by the root of all trees such that:

$$\begin{aligned} \text{lab}_a(\text{root}) \wedge \forall x. \forall y. \quad & (\text{lab}_a(x) \wedge \text{next_sib}(x, y)) \rightarrow \text{lab}_a(y) \\ & \wedge (\text{lab}_a(x) \wedge \text{first_child}(x, y)) \rightarrow \text{lab}_b(y) \\ & \wedge (\text{lab}_b(x) \wedge \text{next_sib}(x, y)) \rightarrow \text{lab}_b(y) \\ & \wedge (\text{lab}_b(x) \wedge \text{first_child}(x, y)) \rightarrow \text{lab}_a(y) \end{aligned}$$

Exercises: Expressiveness of Core XPath 1.0

- ④ Define the following query in Core XPath 1.0:
All nodes reachable from the root over a path with labels in a^*b
Hint: In Core XPath 1.0 you cannot do it forwards.
- ⑤ Can you define the following query in Core XPath?
All nodes that are reachable from the root over a path with labels in $(aa)^*$?
Please argue why! Can you define the same query in monadic Datalog for trees?

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.

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

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

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

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

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

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

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

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

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

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

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

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

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

$$\llbracket /p \rrbracket_{x,y} = x=\text{root} \wedge \llbracket p \rrbracket_{x,y}$$

Query Answering for Core XPath 1.0

For start set $\Pi \subseteq \text{nod}(t)$ let $\text{eval}_{\Pi}^t(p) = \{\pi' \mid \pi \in \Pi, (\pi, \pi') \in \text{eval}^t(p)\}$

Theorem (Gottlob & Koch (TODS'05))

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

- Query answering for FO formulas is PSpace complete in contrast!
- Translated path expressions $\llbracket p \rrbracket_{x,y}$ are simple FO-formulas:
 - ▶ has 3 visible variables in all position only (Immerman and Kozen's result from 89 applies)!
 - ▶ Without negation, we obtain acyclic conjunctive queries (Yanakaki's result from 81 applies)!
- But why is it in combined linear time $O(|t| \cdot |p|)$?
(Combined linear time is quadratic in the strict sence.)

Compute Filter Semantics

Proposition

For every filter expression f of Core XPath 1.0 and every tree t , one can compute $eval^t(f)$ in time $O(|t| \cdot |f|)$.

Simple Filters

$$f ::= * \mid a \mid f[f'] \mid [r::f] \mid not[f]$$

Simple Filters are Enough

replace path composition by filter composition, for instance:

$$[r_1::f_1/r_2::f_2/r_3::f_3] \quad \text{becomes} \quad [r_1::f_1[r_2::f_2[r_3::f_3]]]$$

Proof of Proposition

by induction on the structure of simple filters. Computing $eval^t(*)$, $eval^t(a)$, $eval^t(f[f'])$ and $eval^t(not[f])$ is easy. It remains to compute $eval^t([r::f])$ for all possible axis r .

Compute $eval^t([r::f])$

Generic program can be in $O(|t|^2 \cdot |f|)$

```
for  $\pi$  in  $eval^t(f)$  do
  for  $\pi'$  with  $(\pi', \pi) \in r^t$  do
    collect( $\pi'$ )
```

Interesting case: all descendants of set of nodes $r = (\text{child}^*)^{-1}$

compute $\Pi = eval^t(f)$, run down all paths of t in parallel, memoize whether node in Π has been seen, and collect all later nodes in this case. This is in $O(|t|)$ plus the time for computing $eval^t(f)$ which is in $O(|t| \cdot |f|)$ by induction hypothesis. Altogether this is $O(|t| \cdot |[r::f]|)$.

Be careful

Linear time combined complexity depends on the nature of the axis!

Compute Path Semantics

Set based computation

$$eval_{\Pi}^t(p/p') = eval_{eval_{\Pi}^t(p)}^t(p')$$

$$eval_{\Pi}^t(p \cup p') = eval_{\Pi}^t(p') \cup eval_{\Pi}^t(p)$$

$$eval_{\Pi}^t(r::f) = \{\pi' \mid \exists \pi \in \Pi : (\pi, \pi') \in r^t\} \cap eval^t(f)$$

Linear time complexity

Consider $eval_{\Pi}^t(r::f)$ above. The first set above on the right can be computed in time $O(|t|)$ with the same trick as for filters (and thus depends on the choice of axis), the second set has been computed at before hand in time $O(|t| \cdot |f|)$, see above Proposition.

Exercises [Query Answering Algorithm for Core XPath 1.0]

A node π follows another π' in a trees t , if during a preorder traversal, the first visit of π' is before the first visit of π .

- 6 Present an algorithm that inputs a tree t and a set $\Pi \subseteq \text{nod}(t)$ and computes in time $O(|t|)$ the node set $\text{eval}_{\Pi}^t(\text{following})$.
- 7 The query answering algorithm from the lecture applies to Core XPath 1.0 with the restricted set of axis above. Is it possible to extend this algorithm, such that it supports that axis `following` in addition, while preserving combined linear time complexity? Please explain your answer.

Core XPath 1.0

- 13 Define in CoreXPath 1.0 the relations:
- firstchild* relating a node to its first child if there is some
 - lastchild* relating a node to its last child if there is some
 - firstleaf* relating a node to its left most leaf descendant if there is some
 - lastleaf* relating a node to its right most leaf descendant if there is some
- 14 Define in the first order-logic of unranked trees with relations *child** and *next_sib** the relations:
- firstchild** the reflexive transitive closure of firstchild
 - lastchild** the reflexive transitive closure of lastchild
- Can you define the same relations in Core XPath 1.0 ? Or is some useful operator missing there?

Outline

- 1 Datalog Queries for Unranked Trees
- 2 First-order Queries for Unranked Trees
- 3 Core XPath 1.0
- 4 Core XPath 2.0**
- 5 Beyond the Core

XPath 2.0 (W3C, 2007)

“Few” Variables are Back

- variables to capture nodes as in [hybrid logic](#) (Blackburn 95).
- define [n-ary queries](#) by navigations paths with n-variables
- not available in XPath 1.0

Recover Missing FO-Operators

- complementation on path level
- universal quantifiers, but not added below since not needed for expressiveness.

Core XPath 2.0 (Ten Cate & Marx, ICDT 2007)

Syntax where $a \in \Delta$ and $x \in \text{Vars}$:

filter $f ::= * \mid a \mid \text{not}[f] \mid f[p] \mid [x]$
axis $r ::= \text{child} \mid \text{next_sib} \mid r^{-1} \mid r^* \mid \text{self} \mid \dots$
paths $p ::= r::f \mid p/p' \mid p \cup p' \mid /p \mid p^{\text{compl}}$

Semantics for $t \in T_\Delta$ and $\alpha : \text{Vars} \rightarrow \text{nod}(t)$:

$\text{eval}^{t,\alpha}(f) \subseteq \text{nod}(t)$ $\text{eval}^{t,\alpha}([x]) = \{\alpha(x)\}$
 $\text{eval}^{t,\alpha}(r) \subseteq \text{nod}(t)^2$
 $\text{eval}^{t,\alpha}(p) \subseteq \text{nod}(t)^2$ $\text{eval}^{t,\alpha}(p^{\text{compl}}) = \text{nod}(t)^2 \setminus \text{eval}^{t,\alpha}(p)$

Example

Select all author-title-pairs in books of a bibliography:

Core XPath 2.0 `/child*::book`
 `[child::author[x]]`
 `[child::title[y]]`

FO $\exists z_1, z_2($
 $z_1 = \text{root} \wedge \text{child}^*(z_1, z_2) \wedge \text{lab}_{\text{book}}(z_2) \wedge$
 $\text{child}(z_2, x) \wedge \text{lab}_{\text{author}}(x) \wedge$
 $\text{child}(z_2, y) \wedge \text{lab}_{\text{title}}(y))$

FO-Expressiveness

Theorem [Marx PODS'05]

All FO-queries can be expressed by Core XPath 2.0 queries and vice versa.

Proof.

- quantifier elimination (Schwentick 2000):
 - ▶ Ehrenfeucht-Frässe games
 - ▶ Shelah's composition method
- these techniques are difficult.



Theorem [Filiot&Niehren&Talbot&Tison PODS'07]

There is an FO-complete fragment of Core XPath 2.0 with polynomial time query answering:

- Syntactic restrictions on variables:
 - ▶ For all subexpressions p/p' : no variable sharing between p and p'
 - ▶ No variable below complementation

Example: FO Query Not Expressible in Core XPath 1.0

All nodes having a descendant over a path labeled by a^*b

In FO logic: select all nodes x with b -labeled descendant z such that all nodes y between x and z are labeled by a .

$\text{child}^*(\text{root}, x) \wedge \exists z. (\text{lab}_b(z) \wedge \forall y. ((\text{child}^*(x, y) \wedge \text{child}^*(y, z)) \rightarrow \text{lab}_a(y)))$

In extensions of Core XPath 1.0 with reflexive transitive closure:

$/\text{child}^*::*[(\text{child}::a)^*/\text{child}::b]$

In Core XPath 2.0:

$/\text{child}^*::*[((\text{child}^*::*) \cap (\text{child}^*::*[\text{not}[a]]/\text{child}^*)^{\text{compl}})/\text{child}::b]$

Conditional XPath (Marx 05) with Variables

Syntax where $a \in \Delta$ and $x \in \text{Vars}$:

filter	f	$::=$	$*$	$ $	a	$ $	$\text{not}[f]$	$ $	$f[p]$	$ $	$[x]$
axis	r	$::=$	child	$ $	next_sib	$ $	r^{-1}	$ $	self	$ $	\dots
paths	p	$::=$	$(r::f)^*$	$ $	$r::f$	$ $	p/p'	$ $	$p \cup p'$	$ $	$/p$

Expressiveness

Conditional XPath = F0 = Core XPath 2.0

Efficiency

Without variables, the algorithm of Gottlob and Koch extends to Conditional XPath in combined linear time, but fails for Core XPath 2.0 (where cubic time depending on the data size is needed)

Expressiveness of Core XPath 2.0

- 9 How would you define the FO-logic of words in $\{a, b\}^*$. Can you define the set of words a^*b^* by a closed formula in FO-logic? Can you define the set of all trees with the following DTD in Core XPath 2.0?

$$c \rightarrow a^*b^*, a \rightarrow \epsilon, b \rightarrow \epsilon$$

- 10 Define a query in Core XPath 2.0 that selects all pairs of teacher names and lecture titles in the database from the first lecture.
- 11 Define the conditional path $(\text{child}[f])^*$ in Core XPath 2.0 for arbitrary Core XPath 2.0 filters f .

Outline

- 1 Datalog Queries for Unranked Trees
- 2 First-order Queries for Unranked Trees
- 3 Core XPath 1.0
- 4 Core XPath 2.0
- 5 Beyond the Core**

Extensions of Data Model

Data Values

$$\begin{aligned} \textit{value} &::= \textit{string} \mid \textit{number} \mid \dots \\ \textit{string} &::= \textit{"abiteboul"} \mid \dots \\ \textit{number} &::= 0 \mid 1 \mid 2 \mid 3.4 \mid \dots \\ \textit{boolean} &::= 0 \mid 1 \end{aligned}$$

Data Trees

trees annotated with data values

every node has a content of type string

every node has a string value for all of its attributes

XPath Extension for Data Trees

Functions

access to node content values of node attributes	$text : unit \rightarrow string$ $@a : string \cup number$
counting set arithmetics	$count : 2^{value} \rightarrow number$ $+ : number \times number \rightarrow number$
equality disequality number comparison	$= : value \times value \rightarrow bool$ $\neq : value \times value \rightarrow bool$ $\leq, \geq : number \times number \rightarrow bool$

Joins, Aggregations, Functions

$$f ::= [p = p'] \mid [p \neq p'] \mid [p \leq p'] \mid [p \geq p'] \mid \dots$$
$$p ::= text() \mid @a \mid count(p) \mid p + p' \mid \dots$$

Examples in Short Notation

count number of authors of a book: *count(author)*

select last names of all teachers of lectures: *//teacher/lname/text()*

count number of teachers of a lecture

count(/teacher/lname/text()))

test whether there exist two books with the same author

[/bib[book[(child:/author/@id) = (/next_sib::*author)/@id]]]*

books by Springer after 1991:

//bib/book[publisher = "Springer" and @year ≥ 1991]