

# Query Answering for Core XPath 1.0 and Beyond

## Gottlob and Koch's Algorithm Revisited

Joachim Niehren

Mostrare project, INRIA and LIFL, Lille

November 10, 2020

# 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 = " abitebool"]/  
author[not[self = " abiteboul"]]`

XPath long: `child* :: book[child :: author[self = " abiteboul"]]/  
child :: author[not[self = " abiteboul"]]`

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

# Core XPath 1.0

Syntax where  $a \in \Delta$ :

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

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

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

# Simple Filters

## Adding label tests $\ell$ for negative axis

$$\begin{aligned} f, f' &::= a \mid \ell \mid f \wedge f' \mid [r :: f] \mid \text{not}[f] \mid f \vee f' \\ \ell &::= \text{root} \mid \text{leaf} \mid \text{left} \mid \text{right} \\ r &::= \text{self} \mid \text{child} \mid \text{child}^{-1} \mid \text{child}^* \mid \text{child}^{*-1} \\ &\quad \mid \text{next\_sib} \mid \text{next\_sib}^{-1} \mid \text{next\_sib}^* \mid \text{next\_sib}^{*-1} \end{aligned}$$

# Positive Simple Filters

## Expressing label test for negative axis

The additional label tests  $\ell$  can be expressed as follows:

$$\begin{aligned} \text{root} &= \text{not}[\text{child}^{-1} :: *] \\ \text{leaf} &= \text{not}[\text{child} :: *] \\ \text{left} &= \text{not}[\text{next\_sib}^{-1} :: *] \\ \text{right} &= \text{not}[\text{next\_sib} :: *] \end{aligned}$$

## Positive simple filters

Simple filters without negation.

- Label tests are no more redundant.
- Disjunction is no more redundant.

# Eliminate Negation

## In Axis Only Filters

Function *elim'* removes negated axis-only filters *neg*( $[r :: *]$ ):

$$\text{elim}'(\text{not}[\text{child} :: *]) = \text{leaf}$$
$$\text{elim}'(\text{not}[\text{child}^{-1} :: *]) = \text{root}$$
$$\text{elim}'(\text{not}[\text{next\_sib} :: *]) = \text{left}$$
$$\text{elim}'(\text{not}[\text{next\_sib}^{-1} :: *]) = \text{right}$$
$$\text{elim}'(\text{not}[r^* :: *]) = a[b] \text{ where } a \neq b \text{ fixed arbitrarily}$$

## Eliminate Negation 2

Function *elim* maps maps simple filters to positive simple filters by pushing down negation into label tests and axis, and removing double negations:

### Push Negation Down

$elim(not[*]) = a[b]$  where  $a \neq b$  fixed arbitrarily

$elim(not[a]) = \bigvee_{b \in \Delta \setminus \{a\}} b$

$elim(not[root]) = [child^{-1} :: *]$

$elim(not[leaf]) = [child :: *]$

$elim(not[left]) = [next\_sib^{-1} :: *]$

$elim(not[right]) = [next\_sib :: *]$

$elim(not[f \wedge f']) = not[elim(f)] \vee not[elim(f')]$

$elim(not[r :: f]) = elim'(not[r :: *]) \vee [r :: *[not[elim(f)]]]$

$elim(not[not[f]]) = elim(f)$

$elim(not[f \vee f']) = elim(not[f]) \wedge elim(not[f'])$

# Eliminate Negation 3

## Continue Homomorphically

$$\text{elim}(\ast) = \ast$$

$$\text{elim}(a) = a$$

$$\text{elim}(\text{root}) = \text{root}$$

$$\text{elim}(\text{leaf}) = \text{leaf}$$

$$\text{elim}(f \wedge f') = \text{elim}(f) \wedge \text{elim}(f')$$

$$\text{elim}([r :: f]) = [r :: \text{elim}(f)]$$

$$\text{elim}([f \vee f']) = \text{elim}(f) \vee \text{elim}(f')$$



# Positive Filters in Monadic Datalog

The monadic Datalog program of a filter  $f$  consists of the set of all clauses for the predicates  $q_{f'}$  as defined below such that  $f'$  is a subexpressions  $f$ .

## Basic cases

For all filters  $f, f'$  and labels  $a \in \Delta$ :

$$q_*(x) :- .$$

$$q_a(x) :- lab_a(x).$$

$$q_\ell(x) :- \ell(x). \quad (\ell \in \{root, leaf, left, right\})$$

$$q_{f \wedge f'}(x) :- q_f(x), q_{f'}(x).$$

$$q_{f \vee f'}(x) :- q_f(x).$$

$$q_{f \vee f'}(x) :- q_{f'}(x).$$

# Positive Filters in Monadic Datalog

## Axis

For all  $r \in \{\text{child}, \text{next\_sib}\}$  we define the clauses:

$$\begin{aligned}q_{[r::f]}(x) &:- r(x, y), q_f(y). \\q_{[r^{-1}::f]}(x) &:- r(y, x), q_f(y). \\q_{[r^*::f]}(x) &:- q_f(x). \\q_{[r^*::f]}(x) &:- r(x, y), q_{[r^*::f]}(y). \\q_{[r^{-1}*::f]}(x) &:- q_f(x). \\q_{[r^{-1}*::f]}(x) &:- r(y, x), q_{[r^{-1}*::f]}(y).\end{aligned}$$

# Removing `child` in favor of `first_child`

Replace

$$q(x) :- \text{child}(x, y), q'(y).$$

by

$$q(x) :- \text{first\_child}(x, y), \tilde{q}(y).$$

$$\tilde{q}(y) :- q'(y).$$

$$\tilde{q}(y) :- \text{next\_sib}(y, z), \tilde{q}(z).$$

and similarly for the symmetric clause:

$$q(x) :- \text{child}(y, x), q'(y).$$

# Results

## Proposition

Any simple filter  $f$  of Core XPath 1.0 can be translated to an equivalent monadic Datalog program in time  $O(|f|)$ .

## Proposition

Given a monadic Datalog program  $P$  and a tree  $t$  an closed monadic Datalog program for computing the least fixed point of  $P$  on  $t$  can be computed in time  $O(|P||t|)$ .

## Theorem

For any simple filter  $f$  of Core XPath 1.0 and any tree  $t$  we can compute  $eval_t(f)$  in time  $O(|f||t|)$ .