Efficient Static Analysis of XML Paths and Types

Pierre Genevès – EPFL, Switzerland

Joint work with Nabil Layaïda and Alan Schmitt - INRIA, France

PLDI'07, San Diego, June 2007

- More and more XML data
- **Objective:** ensuring safety and efficiency of programs that manipulate XML
- Two ways for processing XML:
 - General purpose languages extended with librairies
 DSLs: e.g. XSLT, XQuery (W3C standards) that rely on XPath
- In both cases: static analysis of programs very hard (very complex to detect errors at compile-time)
- This paper: we solve important XML static analysis tasks by reduction to satisfiability of a new tree logic

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)



• Before: complexity too high, implementations out of scope..

naner: ontimal complexity + efficient implementation

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)

/descendant::b/parent::a/child::c



• Before: complexity too high, implementations out of scope..

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)

/descendant::b/parent::a/child::c



• Before: complexity too high, implementations out of scope..

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)

/descendant::b/parent::a/child::c



• Before: complexity too high, implementations out of scope..

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)

/descendant::b/parent::a/child::c



• Before: complexity too high, implementations out of scope..

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)

/descendant::b/parent::a/child::c



| for x in (q) do { |
|-------------------|
| } |
| let n = q; |

Before: complexity too high, implementations out of scope...
 This paper: optimal complexity + efficient implementation

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)

a

• queries (XPath)



$\underbrace{/\text{descendant::b/parent::a/child::c}}_{\forall T = \emptyset} \oplus T \stackrel{?}{=} \emptyset$

| for x in (q) do { |
|-------------------|
| } |
| let n = q; |

Before: complexity too high, implementations out of scope...
This paper: optimal complexity + efficient implementation

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)



$\underbrace{/\text{descendant::b/parent::a/child::c}}_{q} \oplus T \stackrel{?}{=} \emptyset$



Before: complexity too high, implementations out of scope...
 This paper: optimal complexity + efficient implementation

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)

/descendant::b/parent::a/child::c



$$\oplus T \stackrel{?}{=} ($$
child::a/**child::c**

qoptimised

| for x in (q) do { |
|-------------------|
| } |
| let n = q; |

Before: complexity too high, implementations out of scope...
 This paper: optimal complexity + efficient implementation

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)



/descendant::b/parent::a/child::c

$$\oplus T \stackrel{?}{=} \underline{/child::a/child::c}$$

qoptimised



Before: complexity too high, implementations out of scope...

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)

/descendant::b/parent::a/child::c

$$\oplus T \stackrel{?}{\equiv} ($$
child::a/**child::c**

qoptimised

F

Туре

Before: complexity too high, implementations out of scope...

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)



$$\oplus T \stackrel{?}{=} \underline{/child::a/child::c}$$

qoptimised



Before: complexity too high, implementations out of scope...

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)





Before: complexity too high, implementations out of scope...

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)





• Before: complexity too high, implementations out of scope...

- Programs that manipulate XML trees
- Analysis:
 - tree types (XML Schemas, DTDs)
 - queries (XPath)





- Before: complexity too high, implementations out of scope...
- This paper: optimal complexity + efficient implementation

Basic Tasks

- XPath typing
- 2 XPath query comparisons
 - query containment, emptiness, overlap, equivalence

Main Applications

- Static analysis of host languages: error detection, optimization (static type-checkers, optimizing compilers)
- Checking integrity constraints in XML databases

Challenges

• Query comparisons and typing are undecidable for the complete XPath language

Open Questions

- What are the largest XPath fragments with decidable static analysis?
- Which fragments can be effectively decided in a compiler?
- Is there a generic algorithm able to solve all related XPath decision problems?

Difficulties

- Considered XPath operators and their combination (e.g., multidirectional navigation, recursion)
- Checking properties on a possibly infinite set of XML documents
- Very high computational complexity

Challenges

• Query comparisons and typing are undecidable for the complete XPath language

Open Questions

- What are the largest XPath fragments with decidable static analysis?
- Which fragments can be effectively decided in a compiler?
- Is there a generic algorithm able to solve all related XPath decision problems?

Difficulties

- Considered XPath operators and their combination (e.g., multidirectional navigation, recursion)
- Checking properties on a possibly infinite set of XML documents
- Very high computational complexity

Challenges

• Query comparisons and typing are undecidable for the complete XPath language

Open Questions

- What are the largest XPath fragments with decidable static analysis?
- Which fragments can be effectively decided in a compiler?
- Is there a generic algorithm able to solve all related XPath decision problems?

Difficulties

- Considered XPath operators and their combination (e.g., multidirectional navigation, recursion)
- Checking properties on a possibly infinite set of XML documents
- Very high computational complexity

The Logical Approach: Overview

- Find an appropriate logic for reasoning on XML trees
- Formulate the problem into the logic and test satisfiability



Critical Aspects

- The logic must be expressive enough
- The algorithm must be effective in practice for XML translations

The Logical Approach: Overview

- Find an appropriate logic for reasoning on XML trees
- Formulate the problem into the logic and test satisfiability



Critical Aspects

- The logic must be expressive enough
- On the algorithm must be effective in practice for XML translations

Models for XML Documents

- Finite ordered binary trees, one label per node
- Bijective encoding of unranked trees as binary trees:





2

Formulas of the \mathcal{L}_{μ} Logic

Closed formulas

• The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

- $\mu Z.\varphi$: finite recursion
- $\{\overline{1},\overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

• The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

Translating following-sibling::a

in \mathcal{L}_{μ} :

- $\mu Z.\varphi$: finite recursion
- $\{\overline{1},\overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

• The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

Translatingfollowing-sibling::ain \mathcal{L}_{μ} :a

- $\mu Z.\varphi$: finite recursion
- $\{\overline{1},\overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

• The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

Translatingfollowing-sibling::ain \mathcal{L}_{μ} :a



- $\mu Z.\varphi$: finite recursion
 - $\{\overline{1},\overline{2}\}$ required for forward axes!
 - {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

- The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node
- Translatingfollowing-sibling::ain \mathcal{L}_{μ} :a



- $\mu Z.\varphi$: finite recursion
- $\{\overline{1},\overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

- The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node
- Translatingfollowing-sibling::ain \mathcal{L}_{μ} :a



- $\mu Z.\varphi$: finite recursion
 - $\{\overline{1},\overline{2}\}$ required for forward axes!
 - {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

• The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

Translatingfollowing-sibling::ain \mathcal{L}_{μ} : $a \land (\mu Z. \langle \overline{2} \rangle \otimes \lor \langle \overline{2} \rangle Z)$



• $\mu Z.\varphi$: finite recursion

- $\{\overline{1},\overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

 The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

Translatingfollowing-sibling::a/preceding-sibling::bin \mathcal{L}_{μ} : $a \land (\mu Z. \langle \overline{2} \rangle \circledast \lor \langle \overline{2} \rangle Z)$



• $\mu Z.\varphi$: finite recursion

- $\{\overline{1},\overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

 The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

Translatingfollowing-sibling::a/preceding-sibling::bin \mathcal{L}_{μ} : $b \land [\mu Y. \langle 2 \rangle (a \land (\mu Z. \langle \overline{2} \rangle \circledast \lor \langle \overline{2} \rangle Z)) \lor \langle 2 \rangle Y]$



• $\mu Z.\varphi$: finite recursion

- $\{\overline{1},\overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

• The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

Translatingfollowing-sibling::a/preceding-sibling::bin \mathcal{L}_{μ} : $b \land [\mu Y. \langle 2 \rangle (a \land (\mu Z. \langle \overline{2} \rangle \circledast \lor \langle \overline{2} \rangle Z)) \lor \langle 2 \rangle Y]$



- $\mu Z.\varphi$: finite recursion
- $\{\overline{1},\overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

 The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

Translatingfollowing-sibling::a/preceding-sibling::bin \mathcal{L}_{μ} : $b \land [\mu Y. \langle 2 \rangle (a \land (\mu Z. \langle \overline{2} \rangle \circledast \lor \langle \overline{2} \rangle Z)) \lor \langle 2 \rangle Y]$



- $\mu Z.\varphi$: finite recursion
- $\{\overline{1}, \overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)

Schemas can also be captured!

• The set of models of a formula φ is the set of finite binary trees for which φ is satisfied on some node

Translatingfollowing-sibling::a/preceding-sibling::bin \mathcal{L}_{μ} : $b \land [\mu Y. \langle 2 \rangle (a \land (\mu Z. \langle \overline{2} \rangle \circledast \lor \langle \overline{2} \rangle Z)) \lor \langle 2 \rangle Y]$



- $\mu Z.\varphi$: finite recursion
- $\{\overline{1}, \overline{2}\}$ required for forward axes!
- {1,2} required for reverse axes!
- Converse programs are crucial
- Almost full XPath can be translated (only variable counting constraints and data value comparisons left)
- Schemas can also be captured!

Satisfiability-Testing Algorithm: Principles

Search for a Tree that Satisfies ψ

- ψ truth status can be determined from a few of its subformulas
- A node is a ψ -type (conjunction of formulas)

Bottom-up Construction of a Tree of ψ -types

- A set T of ψ-types is repeatedly updated (least fixpoint computation)
 - Initially: ∅
 - Step 1 : all possible leaves are added
 - Step *i* : all possible parent nodes of current nodes are added

Termination

- If ψ is present in some node, then ψ is satisfiable
- Otherwise, the algorithm terminates when no more node can be added

Theorem

The satisfiability problem for a formula $\psi \in \mathcal{L}_{\mu}$ is decidable in time $2^{O(n)}$ where n is the size of ψ .

The First Implementation

- Able to handle such a large XPath fragment
- Able to handle schemas (regular tree types)

| What Can Now Be Done | | |
|----------------------|--|--|
| Time (s) | Solved Problems | |
| < 0.5 | Comparisons of XPath queries (XPathmark) without tree types | |
| < 1 | Medium tree types involved (\approx 30 symbols, \approx 20 variables) Example: W3C SMIL | |
| < 3 | Large tree types involved (\approx 100 symbols, \approx 400 variables) Example: W3C XHTML | |

A New Tree Logic

- Best balance known between expressiveness/complexity
- Translation of main XML concepts: linear
- Implementation already fairly efficient for static analysis

Future Work

- Extensions of the logic
 - Decidable data-value comparisons
 - Decidable counting constraints
- Type inference for XSLT/XQuery without output type annotations
- More applications in program analysis?
 - \mathcal{L}_{μ} is as expressive as MSO, and the solver is orders of magnitude faster than MONA...

Thank you!

pierre.geneves@epfl.ch