



**HAL**  
open science

# XQuery and Static Typing: Tackling the Problem of Backward Axes

Pierre Genevès, Nils Gesbert, Nabil Layaïda

► **To cite this version:**

Pierre Genevès, Nils Gesbert, Nabil Layaïda. XQuery and Static Typing: Tackling the Problem of Backward Axes. 2014. hal-00872426v2

**HAL Id: hal-00872426**

**<https://inria.hal.science/hal-00872426v2>**

Preprint submitted on 14 Apr 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# XQuery and Static Typing: Tackling the Problem of Backward Axes

Pierre Genevès  
Univ. Grenoble Alpes, LIG,  
F-38000 Grenoble, France  
CNRS, LIG, F-38000  
Grenoble, France  
Inria  
pierre.geneves@cnrs.fr

Nils Gesbert  
Univ. Grenoble Alpes, LIG,  
F-38000 Grenoble, France  
CNRS, LIG, F-38000  
Grenoble, France  
Inria  
nils.gesbert@grenoble-  
inp.fr

Nabil Layaïda  
Inria  
CNRS, LIG, F-38000  
Grenoble, France  
Univ. Grenoble Alpes, LIG,  
F-38000 Grenoble, France  
nabil.layaida@inria.fr

## ABSTRACT

XQuery is a functional language dedicated to XML data querying and manipulation. As opposed to other W3C-standardized languages for XML (e.g. XSLT), it has been intended to feature strong static typing. Currently, however, some expressions of the language cannot be statically typed with any precision. We argue that this is due to a discrepancy between the semantics of the language and its type algebra: namely, the values of the language are (possibly inner) tree nodes, which may have siblings and ancestors in the data. The types on the other hand are regular tree types, as usual in the XML world: they describe sets of trees. The type associated to a node then corresponds to the subtree whose root is that node and contains no information about the rest of the data. This makes navigation expressions using ‘backward axes,’ which return e.g. the siblings of a node, impossible to type.

We discuss how to handle this discrepancy and give first steps towards a solution. We describe a logic-based language of extended types able to represent inner tree nodes and show how it can dramatically increase the precision of typing for navigation expressions, yet is not adapted for typing copy operations which selectively extract subtrees: regular tree types are better in that case. As inclusion between the extended types and the regular tree types is checkable, in both directions, with a tool which is efficient in practice, we finally suggest that a hybrid system combining both type languages would be the way to go.

## 1. INTRODUCTION

XQuery is a functional language with some unusual features. The standard which defines it [4, 9] describes, among other things, a formal semantics for a core fragment of the language, rules to compile the full language into its core fragment, and a static type system.

Although Turing-complete, this language is not general-purpose; it is designed for manipulating XML data, in various ways. Its type system is thus built around regular tree types, as usual for XML data. The values of the language, however, are not trees or forests, but *sequences of pointers* to tree nodes. These pointers can point anywhere in the tree, not only at the root, and it is always possible, given such a pointer, to get pointers to its parent and sibling nodes. Furthermore, a sequence may contain pointers into different trees.

The formal semantics from the XQuery standard uses judgements of the form  $\text{DynEnv} \vdash \text{Expr} \Rightarrow \text{Val}$ , where  $\text{DynEnv}$  is a store of trees. Navigational expressions (e.g. getting the parent of a node) are evaluated by looking up the initial pointer in the store, navigating in there, and returning a pointer to the destination. However, XQuery is designed as a pure functional language and all the trees in the store are immutable<sup>1</sup>; the only expressions which update the store are those which create a new tree, returning a pointer to its root.

Because of this purity, it is possible to describe the semantics of Core XQuery without using an external store, but only reduction rules for expressions, if we represent tree nodes as *focused trees*, a data structure describing a whole tree ‘seen’ from a given internal node. We believe that it makes it easier to reason about programs. This will be our first contribution (Sec. 2).

This formalization will allow us to highlight a discrepancy between the semantics of XQuery and its type system (Sec. 4.2): whereas the values manipulated by the language consist of a subtree *and a context*, the types describe only the subtree and say nothing of the context. Because of this, expressions navigating upwards or between siblings are simply given the most general type, which contains no information whatsoever, regardless of the type of the initial node.

In order to solve this discrepancy, we then define (Sec. 5) a logic whose formulas denote sets of *focused trees* rather than just of trees, and discuss how it could be combined with the

<sup>1</sup>Expressions that can be used to make persistent changes to instances in the XQuery data model are defined as a separate extension of the language [22].

existing type system.

## 2. SYNTAX AND SEMANTICS OF AN XQUERY NAVIGATIONAL CORE

### 2.1 Values

#### 2.1.1 Items and Sequences

XQuery programs manipulate two ‘levels’ of values: items and sequences. In full XQuery, item values can be literals of various base types (string, boolean etc.), functions (in XQuery 3.0 [23]), and tree nodes. Base values and function values behave in a fairly standard way in XQuery, so, in order to keep this paper to the point, we consider the fragment where all items are tree nodes. Furthermore, we focus on the structure of XML trees and thus consider them composed of only element nodes (with no text content or attributes). This does not imply a loss of generality since literals and text could be encoded as trees.

In XQuery, sequence values are flat lists of items. Nested sequences do not exist. The result of evaluating an expression is always a sequence.

Tree nodes are pointers to nodes which can be anywhere in a tree, not necessarily at the root. Since the tree data structures manipulated by XQuery are always immutable, we need not however actually represent these node values as pointers into a shared data structure defined in an external environment: we may represent them as *focused trees* which contain all the information we need. We detail this structure in the next subsection.

#### 2.1.2 Focused Trees

In order to represent references to nodes of immutable trees, we use *focused trees*, inspired by Huet’s Zipper data structure [16] in the manner of [12]. Focused trees not only describe a tree but also its context: its siblings and its parent, including its parent context recursively.

Formally, we assume an alphabet  $\Sigma$  of labels, ranged over by  $\sigma$ . The syntax of our data model is as follows.

$t ::= \sigma[tl]$	tree
$tl ::= \epsilon$	list of trees
$\quad   \quad t :: tl$	empty list
$c ::= \quad   \quad t :: tl$	cons cell
	context
$\quad   \quad Top$	root of the tree
$\quad   \quad (tl, c[\sigma], tl)$	context node
$f ::= (t, c)$	focused tree

A focused tree  $(t, c)$  is a pair consisting of a tree  $t$  and its context  $c$ . The context  $(tl, c[\sigma], tl)$  comprises three components: a list of trees at the left of the current tree in reverse order (the first element of the list is the tree immediately to the left of the current tree), the context above the tree, and a list of trees at the right of the current tree. The context above the tree may be *Top* if the current tree is at the root, otherwise it is of the form  $c[\sigma]$  where  $\sigma$  is the label of the enclosing element and  $c$  is the context in which the enclosing element occurs.

We now describe how to navigate focused trees, in binary style. There are four directions that can be followed: for a focused tree  $f$ ,  $f \langle 1 \rangle$  changes the focus to the first child of the current tree,  $f \langle 2 \rangle$  changes the focus to the next sibling of the current tree,  $f \langle \bar{1} \rangle$  changes the focus to the parent of the tree *if the current tree is a leftmost sibling*, and  $f \langle \bar{2} \rangle$  changes the focus to the previous sibling.

Formally, we have:

DEFINITION 1.

$$\begin{aligned} (\sigma[t :: tl], c) \langle 1 \rangle &\stackrel{def}{=} (t, (\epsilon, c[\sigma], tl)) \\ (t, (tl_l, c[\sigma], t' :: tl_r)) \langle 2 \rangle &\stackrel{def}{=} (t', (t :: tl_l, c[\sigma], tl_r)) \\ (t, (\epsilon, c[\sigma], tl)) \langle \bar{1} \rangle &\stackrel{def}{=} (\sigma[t :: tl], c) \\ (t', (t :: tl_l, c[\sigma], tl_r)) \langle \bar{2} \rangle &\stackrel{def}{=} (t, (tl_l, c[\sigma], t' :: tl_r)) \end{aligned}$$

When the focused tree does not have the required shape, these operations are not defined.

### 2.2 Expressions

The standard defining XQuery describes how to compile (‘normalize’) expressions of the full language into a core fragment, called the XQuery Core [9]. Although this part of the specification has not been updated after XQuery 1.0, it still is a good starting point.

The formal semantics for this core fragment is defined using an external store, with node items being pointers into that store. What we propose to do is to replace these pointers with focused trees, as described in the previous subsection, which removes the need for a store. As the XQuery Core is already quite large, we will consider a much smaller fragment comprising only constructs impacted by this proposal and useful for the discussion, which we call the navigational core. It is worth noting that several other ‘core fragments’ of XQuery have already been defined and studied in research papers. We will discuss how this one relates to them in the related work section (Sec. 6).

The navigational XQuery fragment we consider is described by the abstract syntax shown in Figure 1, where *axis*  $\in \{\text{child, desc, parent, anc, psibl, nsibl, self}\}$ . The values of the language are sequences  $s$ ; we write  $[f_1, \dots, f_n]$  for  $f_1 :: \dots :: f_n :: \epsilon$ .

We only consider a particular case of the **if-then-else** expression in order to avoid adding boolean expressions to the syntax.

### 2.3 Reduction Semantics

Figure 4 gives reduction rules defining a small-step operational semantics for the focused-tree-based navigational XQuery fragment we consider.

Note that, because  $f \langle \bar{1} \rangle$  and  $f \langle \bar{2} \rangle$  are never both defined for the same  $f$ , rules R-PARENT and R-PSPARENT are mutually exclusive, and that R-NOPARENT can only apply in a case where both  $f \langle \bar{1} \rangle$  and  $f \langle \bar{2} \rangle$  are undefined. The same is true of R-ANC, R-PSANC and R-NOANC, so that the set of rules

$e ::=$	<code>&lt;<math>\sigma</math>&gt;{<math>e</math>&gt;&lt;/<math>\sigma</math>&gt;</code>	expression XML element
	<code><math>\epsilon</math></code>	empty sequence
	<code><math>e, e</math></code>	sequence concatenation
	<code>for <math>\\$v</math> in <math>e</math> return <math>e</math></code>	for loop
	<code>if empty(<math>e</math>) then <math>e</math> else <math>e</math></code>	existence test
	<code><math>\\$v</math>/axis::<math>n</math></code>	tree navigation
	<code><math>\\$v</math></code>	item variable
$n ::=$	<code><math>\sigma</math></code>	name test
	<code>*</code>	wildcard
$s ::=$	<code><math>\epsilon</math></code>	value sequence
	<code><math>f :: s</math></code>	empty sequence
$\mathcal{E} ::=$	<code><math>f :: s</math></code>	cons cell
	<code><math>\mathcal{E}</math></code>	evaluation context
	<code><math>\mathcal{E}</math></code>	context hole
	<code><math>\mathcal{E}</math></code>	
	<code>for <math>\\$v</math> in <math>\mathcal{E}</math> return <math>e</math></code>	
	<code>if empty(<math>\mathcal{E}</math>) then <math>e</math> else <math>e</math></code>	

Figure 1: Navigational core of XQuery.

is almost deterministic. The only ambiguity is the order of concatenation in expressions of the form  $s_1, s_2, s_3$ , but in that case note that the result is independent on that order.

### 3. PROBLEM STATEMENT

Now that we have a simple formal semantics for a core fragment of XQuery, we want to study the problem of type-checking on this fragment.

However, the reader may have noticed that the reduction rules we have defined cannot actually generate an error value, nor can they get stuck on a syntactically correct expression. Thus, in this small fragment, programs cannot go wrong, independently of typing. We will therefore state our problem not in terms of type *safety* but of type *conformance*: the program always generates a result, but we want to check statically whether this result conforms to what has been specified.

We did not include functions in our fragment, but will take as our use case a function definition with type annotations. We thus add to the syntax *sequence variables*  $\$v$  to represent the function's parameters, and the problem we want to solve is the following:

*Given a typing environment  $\Gamma = \$v_1 : \tau_1, \dots, \$v_n : \tau_n$ , an expression  $e$  referring to variables  $\$v_1 \dots \$v_n$ , and a result type  $\tau$ , is it the case that, for any sequence values  $s_1 \dots s_n$  matching types  $\tau_1 \dots \tau_n$ , the expression  $e^{[s_1 \dots s_n / \$v_1 \dots \$v_n]}$  reduces to a value matching type  $\tau$  ?*

As usual, we will require checking procedures for answering this question to be *sound* (give no false positives) but not necessarily *complete* (never give any false negatives). We want them however to be as *precise* as possible (as few false negatives as possible).

We now have to say what types are. The usual formal definition, close to what the standard uses, is given in the next section.

$u ::=$	<code>element <math>n</math> {<math>\tau</math>}</code>	unit type
	<code><math>u \mid u</math></code>	element
$n ::=$	<code><math>\sigma</math></code>	choice
	<code>*</code>	name test
$\tau ::=$	<code><math>u</math></code>	label
	<code>()</code>	wildcard
	<code><math>\tau, \tau</math></code>	sequence type
	<code><math>\tau \mid \tau</math></code>	unit type
	<code><math>\tau^*</math></code>	empty sequence
	<code><math>x</math></code>	concatenation
		choice
		repetition
		type reference

Figure 3: XQuery Types

## 4. THE XQUERY STATIC TYPE SYSTEM

### 4.1 Regular Tree Types

As is customary in the literature ([24, 10, 7, 8] for instance), we use a slight variant of XDuce's type language [14, 15], described in Fig. 3, to represent (core) XQuery types.

Unit types  $u$ , or 'prime types' in the XQuery terminology, correspond to items. Types  $\tau$  correspond to sequences. In the general case,  $u$  would include both element types and base types; since we removed base values from the language fragment we consider, it only includes element types.

A *type environment*  $E$  is a mapping from type references  $x$  to types  $\tau$ . An environment must respect some well-formedness constraints, so that badly-founded recursion such as  $x = () \mid x$  for example is not allowed. These constraints are detailed formally in [15].

The semantics of types is defined in terms of sets of forests, i.e. of sequences of trees (called *elements* in the XML context). A value  $s$ , which is a sequence of *items* (nodes, focused trees in our semantics), *matches* a type if the forest constituted of the subtrees rooted at all nodes of the sequence belongs to the semantics of the type. This is the same forest that is constructed in the tree creation rule R-TREE as the children of the new node.

To give a formal definition, we first define the denotation of a type depending on a function  $d$  mapping variables to sets of forests. We then define the variable denotation  $d_E$  corresponding to a type environment  $E$ . The denotation of a type containing variables is only defined if an environment providing bindings for all these variables is given.

$$\begin{aligned}
\llbracket x \rrbracket_d &= d(x) \\
\llbracket \text{element } \sigma \{ \tau \} \rrbracket_d &= \{ \llbracket \sigma[tl] \rrbracket_d \mid tl \in \llbracket \tau \rrbracket_d \} \\
\llbracket \text{element } * \{ \tau \} \rrbracket_d &= \{ \llbracket \sigma[tl] \rrbracket_d \mid \sigma \in \Sigma \text{ and } tl \in \llbracket \tau \rrbracket_d \} \\
\llbracket () \rrbracket_d &= \epsilon \\
\llbracket \tau, \tau' \rrbracket_d &= \{ \llbracket t_1, \dots, t_n, t'_1, \dots, t'_m \rrbracket_d \mid \\
&\quad \llbracket t_1 \dots t_n \rrbracket_d \in \llbracket \tau \rrbracket_d \text{ and } \llbracket t'_1 \dots t'_m \rrbracket_d \in \llbracket \tau' \rrbracket_d \} \\
\llbracket \tau \mid \tau' \rrbracket_d &= \llbracket \tau \rrbracket_d \cup \llbracket \tau' \rrbracket_d \\
\llbracket \tau^0 \rrbracket_d &= \epsilon \\
\llbracket \tau^{n+1} \rrbracket_d &= \llbracket \tau, \tau^n \rrbracket_d \\
\llbracket \tau^* \rrbracket_d &= \bigcup_{n \in \mathbb{N}} \llbracket \tau^n \rrbracket_d
\end{aligned}$$

Let  $E = (x_i = \tau_i)_{i \in I}$ . Given two mappings  $d_1$  and  $d_2$  from the  $x_i$  to sets of forests, we say that  $d_1$  is smaller than  $d_2$  if:  $\forall i \in I, d_1(x_i) \subseteq d_2(x_i)$ . The variable denotation  $d_E$

$$\begin{array}{c}
\text{(R-TREE)} \quad \langle \sigma \rangle \{[(t_1, c_1), (t_2, c_2), \dots, (t_n, c_n)]\} \langle / \sigma \rangle \longrightarrow [(\sigma[t_1 :: t_2 :: \dots :: t_n :: \epsilon], \text{Top})] \\
\text{(R-FOR)} \quad \text{for } \$v \text{ in } f_1 :: s \text{ return } e \longrightarrow e[f_1/\$v], \text{for } \$v \text{ in } s \text{ return } e \quad \text{(R-FOREMPTY)} \quad \text{for } \$v \text{ in } \epsilon \text{ return } e \longrightarrow \epsilon \\
\text{(R-SINGLETON)} \quad f \longrightarrow [f] \quad \text{(R-CONCAT)} \quad [f_1, \dots, f_n], [f'_1, \dots, f'_{n'}] \longrightarrow [f_1, \dots, f_n, f'_1, \dots, f'_{n'}] \\
\text{(R-IFT)} \quad \text{if empty}(\epsilon) \text{ then } e_1 \text{ else } e_2 \longrightarrow e_1 \quad \text{(R-IFF)} \quad \text{if empty}(f :: s) \text{ then } e_1 \text{ else } e_2 \longrightarrow e_2 \\
\text{(R-NO PARENT)} \quad (t, \text{Top})/\text{parent} :: n \longrightarrow \epsilon \quad \text{(R-NO CHILD)} \quad (\sigma[\epsilon], c)/\text{child} :: n \longrightarrow \epsilon \quad \text{(R-NO SIBL)} \quad (t, (tl, \sigma[c], \epsilon))/\text{nsibl} :: n \longrightarrow \epsilon \\
\text{(R-NO PSIBL)} \quad (t, (\epsilon, \sigma[c], tl))/\text{psibl} :: n \longrightarrow \epsilon \quad \text{(R-NO ANC)} \quad (t, \text{Top})/\text{anc} :: n \longrightarrow \epsilon \quad \text{(R-SELF STAR)} \quad f/\text{self} :: * \longrightarrow [f] \\
\text{(R-SELF MATCH)} \quad (\sigma[tl], c)/\text{self} :: \sigma \longrightarrow [(\sigma[tl], c)] \\
\text{(R-SELFDIFF)} \quad \frac{\sigma \neq \sigma'}{(\sigma[tl], c)/\text{self} :: \sigma' \longrightarrow \epsilon} \quad \text{(R-PARENT)} \quad \frac{f' = f \langle \bar{1} \rangle}{f/\text{parent} :: n \longrightarrow f'/\text{self} :: n} \\
\text{(R-PSPARENT)} \quad \frac{f' = f \langle \bar{2} \rangle}{f/\text{parent} :: n \longrightarrow f'/\text{parent} :: n} \quad \text{(R-CHILD)} \quad \frac{f' = f \langle 1 \rangle}{f/\text{child} :: n \longrightarrow f'/\text{self} :: n, f'/\text{nsibl} :: n} \\
\text{(R-NSIBL)} \quad \frac{f' = f \langle 2 \rangle}{f/\text{nsibl} :: n \longrightarrow f'/\text{self} :: n, f'/\text{nsibl} :: n} \quad \text{(R-PSIBL)} \quad \frac{f' = f \langle \bar{2} \rangle}{f/\text{psibl} :: n \longrightarrow f'/\text{psibl} :: n, f'/\text{self} :: n} \\
\text{(R-ANC)} \quad \frac{f' = f \langle \bar{1} \rangle}{f/\text{anc} :: n \longrightarrow f'/\text{anc} :: n, f'/\text{self} :: n} \quad \text{(R-PSANC)} \quad \frac{f' = f \langle \bar{2} \rangle}{f/\text{anc} :: n \longrightarrow f'/\text{anc} :: n} \\
\text{(R-DESC)} \quad f/\text{desc} :: n \longrightarrow \text{for } \$v \text{ in } f/\text{child} :: * \text{ return } \$v/\text{self} :: n, \$v/\text{desc} :: n \\
\text{(R-CONTEXT)} \quad \frac{e_1 \longrightarrow e_2}{\mathcal{E}[e_1] \longrightarrow \mathcal{E}[e_2]}
\end{array}$$

**Figure 2: Reduction Rules for the Navigational XQuery Fragment**

corresponding to the type environment  $E$  is defined as the smallest mapping such that:  $\forall i \in I, d_E(x_i) = \llbracket \tau_i \rrbracket_{d_E}$ . (It may be undefined in general, but as said above, we assume  $E$  respects well-formedness constraints such that this smallest mapping exists.)

In the following, we do not treat the case of  $x$  explicitly: we always assume we have an environment  $E$  which defines it and implicitly replace the variable with its binding when needed. We also consider that  $E$  always contains the type of all elements, `AnyElt`, defined as `element * {AnyElt*}`.

## 4.2 Typing Rules

In this section, we present an adaptation of the standard W3C XQuery typing rules to our fragment of the language. The rules use type environments  $E$ , which contain the possibly mutually recursive definitions of named types (see above), and typing environments  $\Gamma$  which map sequence variables to sequence types, as said in Sec. 3, and in addition map iteration variables to **unit types** (not types in general). Indeed, in the `for` loop, the variable is bound successively to all items in the input sequence, thus its value is an item, not a sequence.

This system constructs deterministically one type for the whole expression from the types of its subexpressions (which is sometimes called forward type inference). In general, the computed type will not be exactly the one required in our type-checking question, but we can still answer positively if the semantics of the computed type is included in the semantics of the required type, so there is an additional inclusion check (which is straightforward for this regular type

language). Here precision of typing thus means computing the smallest possible output type.

In the XQuery standard, typing the `for` loop involves two operations on types, `prime` and `quant`. The first one constructs the disjunction of all unit types which appear in a sequence type. The second one determines whether a sequence type can contain sequences with zero, one, or several items. Its result is a multiplier from the set  $\{0, 1, +, ?, *\}$ .

Typing the node test involves filtering a type to extract only the matching element types and is done by the `filter` operation. We also define `children` and `dos` operations to extract from a unit type all the possible types for, respectively, its children and its descendants plus itself.

The commutative  $\oplus$  and  $|$  operations on multipliers  $q$ , the `prime` and `quant` operators, the application of a multiplier to a sequence type, the `children` operation on unit types and the `filter` and `dos` operations on sequence types are formally defined as follows:

$$\begin{array}{l}
0 \oplus q = q \\
1 \oplus q = + \text{ if } q \neq 0 \\
+ \oplus q = + \\
q \oplus q' = * \text{ if } (q, q') \in \{*, ?\}^2 \\
q | q = q \\
0 | 1 = ? \\
1 | q = q \text{ if } q \neq 0 \\
+ | q = * \text{ if } q \notin \{1, +\} \\
* | q = *
\end{array}$$

$$\begin{aligned}
& \text{prime}(u) = u \\
& \text{prime}(\text{()}) = \text{()} \\
& \text{prime}(\tau_1 \mid \tau_2) = \text{prime}(\tau_1) \mid \text{prime}(\tau_2) \\
& \text{prime}(\tau_1, \tau_2) = \text{prime}(\tau_1) \mid \text{prime}(\tau_2) \\
& \text{prime}(\tau^*) = \text{prime}(\tau) \\
\\
& \text{quant}(u) = 1 \\
& \text{quant}(\text{()}) = 0 \\
& \text{quant}(\tau_1 \mid \tau_2) = \text{quant}(\tau_1) \mid \text{quant}(\tau_2) \\
& \text{quant}(\tau_1, \tau_2) = \text{quant}(\tau_1) \oplus \text{quant}(\tau_2) \\
& \text{quant}(\tau^*) = 0 \text{ if } \text{quant}(\tau) = 0, * \text{ otherwise} \\
\\
& \tau^0 = \text{()} \\
& \tau^1 = \tau \\
& \tau^+ = \tau, \tau^* \\
& \tau^? = \text{()} \mid \tau \\
& \tau^* = \tau^* \\
\\
& \text{filter}(\text{()}, n) = \text{()} \\
& \text{filter}(\text{element } * \{ \tau \}, n) = \text{element } * \{ \tau \} \\
& \text{filter}(\text{element } n \{ \tau \}, *) = \text{element } n \{ \tau \} \\
& \text{filter}(\text{element } \sigma \{ \tau \}, \sigma) = \text{element } \sigma \{ \tau \} \\
& \text{filter}(\text{element } \sigma \{ \tau \}, \sigma') = \text{()} \text{ if } \sigma \neq \sigma' \\
& \text{filter}(\tau_1 \mid \tau_2, n) = \text{filter}(\tau_1, n) \mid \text{filter}(\tau_2, n) \\
& \text{filter}(\tau_1, \tau_2, n) = \text{filter}(\tau_1, n), \text{filter}(\tau_2, n) \\
& \text{filter}(\tau^*, n) = \text{filter}(\tau, n)^* \\
\\
& \text{children}(\text{element } n \{ \tau \}) = \tau \\
& \text{children}(u_1 \mid u_2) = \text{children}(u_1) \mid \text{children}(u_2) \\
\\
& \text{dos}(\text{()}) = \text{()} \\
& \text{dos}(\text{element } n \{ \tau \}) = \text{element } n \{ \tau \}, \text{dos}(\tau) \\
& \text{dos}(\tau_1 \mid \tau_2) = \text{dos}(\tau_1) \mid \text{dos}(\tau_2) \\
& \text{dos}(\tau_1, \tau_2) = \text{dos}(\tau_1), \text{dos}(\tau_2) \\
& \text{dos}(\tau^*) = \text{dos}(\tau)^*
\end{aligned}$$

We can see that this type system lacks precision in several places. One of them is Rule T-FOR, which gives to the sequence resulting from the loop a *homogeneous type*, disregarding completely the fact that, since the input sequence may contain specific types in a specific order, it should be possible in most cases to deduce ordering information about the different item types which can occur in the output. This source of imprecision is well known and has been addressed in several research papers already, e.g. [10], [8]. Note that it can be addressed while staying within the standard type language presented thus far: the problem is that the type system does not fully exploit the information it has. We will not expand on this point further in this paper.

The other blatant source of imprecision is Rule T-OTHAXIS, and to a very slightly lesser extent Rule T-PARENT. Indeed, type AnyElt\* is the top type of the algebra; any value

whatsoever matches that type<sup>2</sup>. Thus the type given to navigation expressions with axes `psibl`, `nsibl` or `ancestor` contains no information at all, and it is not much better for `parent`. This means that typechecking such an expression will always fail (unless no requirement was made on the result) and that the programmer using it will be forced to by-pass static checking with a type cast.

This source of imprecision is also well known, it has however not been addressed yet as far as we know, except, indirectly (via a translation into another language), in [5]. The reason why it has not, we think, is because it is more fundamental than the problem of `for` loops. If we look at the information available to the type system, i.e. the expression and the typing environment, there is one small improvement we can easily make: if the expression is of the form  $\$v/\text{axis}::\sigma$ , then we know that the result must contain only nodes labelled  $\sigma$ —this is still very limited information. Can we do better? If we look at the semantics of Fig. 4, we know that at some point, when the surrounding `for` loop is unfolded,  $\$v$  will be replaced with a focused tree  $f = (t, c)$ . Then, depending on the axis and the shape of  $f$ , either  $f \langle 2 \rangle$ ,  $f \langle 1 \rangle$  or  $f \langle 2 \rangle$  will be computed to yield the result. The definitions of these operations (Def. 1) show that the first component (the tree) of the result is taken from  $c$ , in the case of  $\langle 2 \rangle$  and  $\langle 2 \rangle$ , or is constructed from both  $t$  and a part of  $c$  in the case of  $\langle 1 \rangle$ .

Now what is the type information we have? the environment  $\Gamma$  maps iteration variables to unit types, which describe nonempty sets of trees. The fact that  $\$v$  has type  $u$  means that when, during evaluation,  $\$v$  is replaced with  $(t, c)$ ,  $t$  will be an element of  $\llbracket u \rrbracket$ —it says nothing on what  $c$  will possibly be. Therefore it is impossible to know anything about the result of an expression such as  $\$v/\text{nsibl}::*$  using only this kind of type environment. The problem is highlighted by our focused-tree-based formalization: item values can be represented as pairs of which only the first component has a type—there is a clear discrepancy between the semantics and the type system.

It is worth noting that in the case of  $\langle 1 \rangle$  we do not have exactly zero information on the result, due to the fact that  $t$ , for which we have a type, appears as the first child of the new tree.  $\langle 1 \rangle$  is used in the reduction rules for the `parent` and `anc` axes. This fact would allow us to improve the T-PARENT rule slightly: we know that, if the parent exists, it has at least one child, which has type  $u$ . We can also use this information when typing the `anc` axis.

We summarise below the very limited improvements which are possible for the typing of navigation expressions within the standard type environment:

$$\begin{array}{c}
\text{T-PARENT}' \\
\frac{}{E; \Gamma, \$v : u \vdash \$v/\text{parent}::n : \text{() } \mid \text{element } n \{ \text{AnyElt}^*, u, \text{AnyElt}^* \}} \\
\\
\text{T-ANC} \\
\frac{E' = E, x = u \mid \text{element } * \{ \text{AnyElt}^*, x, \text{AnyElt}^* \}}{E'; \Gamma, \$v : u \vdash \$v/\text{anc}::n : \text{element } n \{ \text{AnyElt}^*, x, \text{AnyElt}^* \}} \\
\\
\text{T-OTHAXIS}' \\
\frac{\text{axis} \in \{ \text{psibl}, \text{nsibl} \}}{E; \Gamma, \$v : u \vdash \$v/\text{axis}::n : \text{element } n \{ \text{AnyElt}^* \}}
\end{array}$$

<sup>2</sup>Admittedly, if we had base types or function types, it would exclude them.

$$\begin{array}{c}
\text{(T-SEQVAR)} \quad E; \Gamma, \bar{\$v} : \tau \vdash \bar{\$v} : \tau \qquad \text{(T-EMPTY)} \quad E; \Gamma \vdash \epsilon : () \qquad \text{(T-TREE)} \quad \frac{E; \Gamma \vdash e : \tau}{E; \Gamma \vdash \langle \sigma \rangle \{e\} \langle /\sigma \rangle : \text{element } \sigma \{ \tau \}} \\
\text{(T-SEQ)} \quad \frac{E; \Gamma \vdash e_1 : \tau_1 \quad E; \Gamma \vdash e_2 : \tau_2}{E; \Gamma \vdash e_1, e_2 : \tau_1, \tau_2} \qquad \text{(T-IF)} \quad \frac{E; \Gamma \vdash e_1 : \tau_1 \quad E; \Gamma \vdash e_2 : \tau_2}{E; \Gamma \vdash \text{if empty}(e) \text{ then } e_1 \text{ else } e_2 : \tau_1 \mid \tau_2} \\
\text{(T-FOR)} \quad \frac{E; \Gamma \vdash e_1 : \tau_1 \quad E; \Gamma, \$v : \text{prime}(\tau_1) \vdash e_2 : \tau_2}{E; \Gamma \vdash \text{for } \$v \text{ in } e_1 \text{ return } e_2 : \tau_2^{\text{quant}(\tau_1)}} \\
\text{(T-SELF)} \quad E; \Gamma, \$v : u \vdash \$v/\text{self} :: n : \text{filter}(u, n) \qquad \text{(T-CHILD)} \quad E; \Gamma, \$v : u \vdash \$v/\text{child} :: n : \text{filter}(\text{children}(u), n) \\
\text{(T-DESC)} \quad E; \Gamma, \$v : u \vdash \$v/\text{desc} :: n : \text{filter}(\text{dos}(\text{children}(u)), n) \qquad \text{(T-PARENT)} \quad E; \Gamma, \$v : u \vdash \$v/\text{parent} :: n : () \mid \text{AnyElt} \\
\text{(T-OTHAXIS)} \quad \frac{\text{axis} \in \{\text{anc}, \text{psibl}, \text{nsibl}\}}{E; \Gamma, \$v : u \vdash \$v/\text{axis} :: n : \text{AnyElt}^*}
\end{array}$$

Figure 4: Standard Typing Rules for the Navigational XQuery Fragment

$\varphi, \psi ::=$	formula	$\llbracket \top \rrbracket_V \stackrel{\text{def}}{=} \mathcal{F}$	$\llbracket \sigma \rrbracket_V \stackrel{\text{def}}{=} \{f \mid \text{nm}(f) = \sigma\}$
$\top$	true	$\llbracket X \rrbracket_V \stackrel{\text{def}}{=} V(X)$	$\llbracket \neg \sigma \rrbracket_V \stackrel{\text{def}}{=} \{f \mid \text{nm}(f) \neq \sigma\}$
$\sigma \mid \neg \sigma$	atomic proposition (negated)	$\llbracket \varphi \vee \psi \rrbracket_V \stackrel{\text{def}}{=} \llbracket \varphi \rrbracket_V \cup \llbracket \psi \rrbracket_V$	$\llbracket \varphi \wedge \psi \rrbracket_V \stackrel{\text{def}}{=} \llbracket \varphi \rrbracket_V \cap \llbracket \psi \rrbracket_V$
$X$	variable		
$\varphi \vee \psi$	disjunction		
$\varphi \wedge \psi$	conjunction		
$\langle a \rangle \varphi \mid \neg \langle a \rangle \top$	existential (negated)	$\llbracket \langle a \rangle \varphi \rrbracket_V \stackrel{\text{def}}{=} \{f \langle \bar{a} \rangle \mid f \in \llbracket \varphi \rrbracket_V \wedge f \langle \bar{a} \rangle \text{ defined}\}$	
$\mu(X_i = \varphi_i)_{i \in I} \text{ in } \psi$	(least) n-ary fixpoint	$\llbracket \neg \langle a \rangle \top \rrbracket_V \stackrel{\text{def}}{=} \{f \mid f \langle a \rangle \text{ undefined}\}$	

Figure 5: Logic formulas

If we want to do better, we have to change either the type language or the semantics. We first try to do the former, in the next section.

## 5. TYPES FOR FOCUSED TREES

### 5.1 A Tree Logic

In order to describe sets of focused trees rather than just sets of trees, we use the logic language defined in [12]. It is a sub-logic of the alternation free modal  $\mu$ -calculus with converse; its syntax is given in Fig. 5, where  $a \in \{1, 2, \bar{1}, \bar{2}\}$  are *programs*, corresponding to the four directions in which trees can be navigated.

Our main reasons for choosing this formalism are: it is expressive enough to support all XQuery types, it is succinct (types are represented as formulas of linear size compared to their regular expression syntax), and the satisfiability problem for a logical formula of size  $n$  can be efficiently decided with an optimal  $2^{\mathcal{O}(n)}$  worst-case time complexity bound with the solver of [12].

Formulas include the truth predicate, atomic propositions (denoting the label of the node in focus), disjunction and conjunction of formulas, formulas under an existential modality (denoting the existence of a node, in the direction denoted by the program, satisfying the sub-formula), and a fixpoint operator.

The logic allows expressing recursion through the fixpoint binder. The recursive formula  $\langle 1 \rangle (\mu X.a \vee \langle 1 \rangle X \vee \langle 2 \rangle X)$  states the existence of some node labelled with “ $a$ ” at an arbitrary depth in the subtree. The meaning of the recursive formula  $\mu X.b \vee \langle \bar{2} \rangle X$  is that either the current node

$$\begin{array}{l}
\llbracket \mu(X_i = \varphi_i)_{i \in I} \text{ in } \psi \rrbracket_V \stackrel{\text{def}}{=} \\
\text{let } S = \{(T_i) \in \mathcal{P}(\mathcal{F})^I \mid \forall j \in I, \llbracket \varphi_j \rrbracket_{V[\bar{T}_i/\bar{X}_i]} \subseteq T_j\} \text{ in} \\
\text{let } (U_j) = (\bigcap_{(T_i) \in S} T_j)_{j \in I} \text{ in } \llbracket \psi \rrbracket_{V[\bar{U}_j/\bar{X}_j]}
\end{array}$$

where  $V[\bar{T}_i/\bar{X}_i](X) = V(X)$  if  $X \notin \{X_i\}$  and  $T_i$  if  $X = X_i$ .

Figure 6: Interpretation of formulas

is labeled  $b$  or some previous sibling of the current node is labeled  $b$ .

The interpretation of a logical formula is the set of focused trees such that the formula is satisfied at the current node. The semantics of this logic is intuitively explained through examples in [11] and formally defined in [12]. We give the formal definition in Fig. 6, where  $\mathcal{F}$  is the set of all focused trees and  $\text{nm}(f)$  is the label at the current node of  $f$ .

The lemma 4.2 of [12] says that the interpretation of a fixpoint formula is equal to the union of the interpretations of all its finite unfoldings (where unfolding is defined as usual). A consequence (detailed in [12]) is that the logic is closed under negation, i. e. for any closed  $\varphi$ ,  $\neg \varphi$  can be expressed in the syntax using De Morgan’s relations and this definition:

$$\begin{array}{l}
\neg \langle a \rangle \varphi \stackrel{\text{def}}{=} \neg \langle a \rangle \top \vee \langle a \rangle \neg \varphi \\
\neg \mu(X_i = \varphi_i) \text{ in } \psi \stackrel{\text{def}}{=} \mu(X_i = \neg \varphi_i \{ \overline{[X_i/\neg X_i]} \}) \text{ in } \neg \psi \{ \overline{[X_i/\neg X_i]} \}
\end{array}$$

In the following, we consider only closed formulas and write  $\llbracket \varphi \rrbracket$  for  $\llbracket \varphi \rrbracket_{\emptyset}$ .

### 5.2 Adding Formulas to the Type System

We now have a language which allows us to describe sets of focused trees. If we take the type language of Sec. 4 and replace unit types with formulas, we obtain a type language

for sequences of focused trees.

Recall that typechecking in the standard type system involved inferring a type for the expression using the rules then checking an inclusion between this type and the expected type. With formulas, the inclusion test translates into a satisfiability test (namely,  $\llbracket \varphi \rrbracket \subseteq \llbracket \psi \rrbracket$  if and only if  $\varphi \wedge \neg \psi$  is unsatisfiable), and as mentioned above, we have an efficient decision procedure for this test in this language. The remaining question is whether we can adapt the typing rules to this modified type language.

Looking at the type system, the change of type language does not affect the structural rules (T-SEQVAR, T-EMPTY, T-SEQ, T-IF, T-FOR) since none of them looks inside unit types, so this will only change the rules for tree construction and navigation expressions. Let us start with the parent rule. If we want the highest possible precision, it becomes these two rules:

$$\frac{\text{T-NONNULLABLEPARENT} \quad \varphi \wedge \neg \langle \bar{1} \rangle k(n) \text{ is unsatisfiable}}{\Gamma, \$v : \varphi \vdash \$v/\text{parent} :: n : k(n) \wedge \langle 1 \rangle \mu X.(\varphi \vee \langle 2 \rangle X)}$$

$$\frac{\text{T-NULLABLEPARENT} \quad \varphi \wedge \neg \langle \bar{1} \rangle k(n) \text{ is satisfiable}}{\Gamma, \$v : \varphi \vdash \$v/\text{parent} :: n : () \mid k(n) \wedge \langle 1 \rangle \mu X.(\varphi \vee \langle 2 \rangle X)}$$

where  $k(\sigma) \stackrel{\text{def}}{=} \sigma$  and  $k(*) \stackrel{\text{def}}{=} \top$ .

In the first case, the type of  $\$v$  implies that it must have a parent which matches  $n$  (this is what the premise tests). In that case, the result is always a sequence of one node and thus its type is a unit type, i. e. a logic formula. This formula says we can go from the node at focus to a node where  $\varphi$  is satisfied by going down once and then right a certain number of times. This is as precise as we can do.

The second rule deals with the fact that the result may be the empty sequence if the type of  $\$v$  allows it to not have a parent matching  $n$ . In that case the result's type is not a single formula anymore but a (very simple) regular expression of formulas. Note that it would be sound, if less precise, to use only this rule without the test on the premise. However, knowing by its type that a sequence cannot be empty is often a useful information (for example, Rule T-IF can straightforwardly be improved to take advantage of it in our particular case of the existence test).

Unfortunately, things do not go as perfectly well for other axes (except of course `self` which is trivial), because the other axes generate *sequences* of nodes of variable length. Thus the type of the expression should, again, be a regular expression of formulas. Ideally, this regular expression would reflect somehow the ordering of differently-typed nodes in the sequence, but it is certainly not possible to be exact unless the sequence has a statically known length limit (as was the case for parent)—indeed, every node in the sequence corresponds to a different focused tree and it is possible to set any of them apart from the others using a carefully crafted (and possibly large) formula. There is no most precise regular expression that we can look for. The problem of keep-

ing the correct amount of ordering information is not at all straightforward and is beyond the scope of this paper; we intend to investigate it in future work.

The simplest way to infer a type for the other axes is to consider that all items in the sequence must have the same type (as is usually the case in general-purpose programming languages). Note that it makes the new rules for `child` and `desc` slightly *less* precise than the standard rules (unless the resulting sequence is then iterated over, in which case the ordering information is lost anyway in the standard system).

We give below typing rules for the other axes using logic formulas:

$$\text{T-CHILD''} \quad E; \Gamma, \$v : \varphi \vdash \$v/\text{child} :: n : (k(n) \wedge \mu X.(\langle \bar{1} \rangle \varphi \vee \langle \bar{2} \rangle X))*$$

$$\text{T-DESC''} \quad E; \Gamma, \$v : \varphi \vdash \$v/\text{desc} :: n : (k(n) \wedge \mu X.(\langle \bar{1} \rangle \varphi \vee \langle \bar{1} \rangle X \vee \langle \bar{2} \rangle X))*$$

$$\text{T-ANC''} \quad E; \Gamma, \$v : \varphi \vdash \$v/\text{anc} :: n : (k(n) \wedge \langle 1 \rangle \mu X.(\varphi \vee \langle 1 \rangle X \vee \langle 2 \rangle X))*$$

$$\text{T-PSIBL} \quad E; \Gamma, \$v : \varphi \vdash \$v/\text{psibl} :: n : (k(n) \wedge \mu X.(\langle 2 \rangle \varphi \vee \langle 2 \rangle X))*$$

$$\text{T-NSIBL} \quad E; \Gamma, \$v : \varphi \vdash \$v/\text{nsibl} :: n : (k(n) \wedge \mu X.(\langle \bar{2} \rangle \varphi \vee \langle \bar{2} \rangle X))*$$

Note that it is perfectly possible to split each of these rules in two, just as we illustrated for the case of T-PARENT, depending whether the sequence can be statically known not to be empty, and this would play well with an improved T-IF. But contrary to T-PARENT, this would just be one arbitrary (even though sensible) way of making the rules more precise, among infinitely many. We leave this as an exercise for the reader.

Perhaps surprisingly, T-TREE, which was completely straightforward in the standard system, is the most difficult rule to adapt to formula types. Indeed, the semantics of the tree construction expression (rule R-TREE from Fig. 4) shows that the context part of the input focused trees is cut off and only the subtree part is kept. This looks deceptively simple because we represented focused trees as pairs, but logic formulas describe the focused trees as a whole, they do not separate the context and subtree components. Thus we do not know how to infer a formula describing the new element from a regular expression of formulas describing the input node sequence.

We would like to emphasize the fact that, in an actual XQuery implementation, R-TREE is indeed the most complicated operation in our small fragment: the other ones mostly manipulate pointers to navigate in a data structure. This operation, in contrast, involves the deep copy of some specific, possibly overlapping parts of possibly several structures, and the assembling of the results into a brand new structure. If we take that into account, the fact that this operation is difficult to type when we use formula-based types which try to closely describe the structure of values is less surprising.

However, this shows limits of the purely formula-based approach. Since the standard approach and the formula-based

approach are weak on different parts of the language, we can try combining the two.

### 5.3 Combining the two approaches

We saw that logic formulas are closer to the semantics of the XQuery language than standard XQuery types, and allow a reasonably precise typing of expressions involving non-downward navigation. However they are not well suited for use as the main types in a whole program, for at least two reasons. One is the difficulty of typing the element construction expression. Another one is type annotations: we stated our type-checking problem in terms of conformance to a specification. XQuery programs manipulate XML documents and, in general, both the input and output types will be specified as regular types, which are the standard for XML documents.

It therefore makes sense to investigate the possibility of a hybrid system mixing the two. It is possible, and quite straightforward, to translate standard unit types into logic formulas describing exactly the same set of trees, with or without the constraint that the node described by the unit types is at the root, as needed. We then can use the satisfiability solver to check whether a unit type is included in a logic type, as well as the converse. This means that the type-checking question described in Sec. 3 can be asked using standard regular types and answered using logic.

Furthermore, we can take advantage of the intrinsic compositionality of a type system: mostly, i.e. for all structural rules, this type-checking question is answered by splitting it into subquestions. We can answer some of the subquestions using logic and some of them using the standard type system.

## 6. RELATED WORK

Static typing for XQuery has been standardized by the W3C [9] and improved by [8]. [9] describes a type system which is polynomial (except for nested let clauses). Colazzo goes one step further by a thorough analysis of precision and complexity of this type-system and by introducing his own, more precise (but exponential) type system. The type-system proposed by the W3C has been inspired by the seminal work found in [14], which is itself based on finite tree automata containment [15]. A more precise typing of `for` loops than what made it into the standard had been studied in [10]. None of these type systems supports non-downward navigation in XML trees, despite it being a part of the XQuery standard.

In the programming languages community, the XML type-checking problem has also been studied for other particular domain-specific languages such as CDuce [2], XSLT [17] or with specific transformers like transducers [19, 20]. For a recent survey of related works on type-checking for XML, see [3] and references thereof.

The language fragment we decided to study formally is inspired by what can be found in the literature. Other formally studied fragments include XQ (‘core XQuery’) [18], recently extended into XQ<sub>H</sub> by Benedikt and Vu who added higher-order functions [1], and  $\mu$ XQ (‘micro XQuery’) [7], extended into  $\mu$ XQ<sup>+</sup> (‘mini XQuery’) in [8]. The papers defining XQ

and XQ<sub>H</sub> focus on the semantics of the language and the complexity of query evaluation. The papers defining  $\mu$ XQ and  $\mu$ XQ<sup>+</sup> focus on typing and correctness. There does not seem to be a significant difference between XQ and  $\mu$ XQ<sup>+</sup>, although the set of expressions is not exactly the same.  $\mu$ XQ does not include conditionals (`where` or, equivalently, `if-then-else` expressions), and XQ<sub>H</sub> adds  $\lambda$ -abstractions and function application. None of these fragments includes axes other than `child` and `desc`. This allows XQ and XQ<sub>H</sub> to have a formal semantics where items are simply trees without the need for a store, because it is not possible in these fragments to go from a node to its parent or siblings. Our ‘navigational XQuery’ fragment is basically XQ/ $\mu$ XQ<sup>+</sup> with the upward and sideways axes added.

Very recently, Castagna et al. [5] have defined a larger fragment which they call XQ<sub>H</sub><sup>+</sup>; it adds to XQ<sub>H</sub> value-switch and type-switch expressions, as well as non-downward axes as we do. They do not study this fragment directly: rather, they define a translation from it to an extension of the CDuce language, and study that extension, whose semantics is described by a fairly complicated system—not surprisingly since it has to integrate with CDuce’s pattern-matching, overloaded functions and other advanced features. In contrast, we tried to focus on just the unusual features of XQuery and keep things as simple as we can. Still, this work is probably the closest, in terms of objective, to the current paper, since it does treat backward axes in some way.

The logic language for XML trees which we described in Sec. 5 was described in more detail in [12], where it was used to check properties of XPath expressions. The typing rules we give in Sec. 5.2 are inspired by the way XPath subexpressions were translated into logic in that paper.

## 7. CONCLUSION

The work presented in this paper is not an accomplished type-checking system for XQuery, but what we hope will prove a useful basis and starting point to, finally, take *all* navigation expressions properly into account in such systems.

Our contribution is threefold. First, we defined a novel focused-tree-based operational semantics for a fragment of the XQuery language; this fragment was kept small here to concentrate on the core issues but can be easily extended. Second, we formulated the difficulty of typing XQuery expressions with backward axes in terms of a discrepancy between the language’s semantics and type algebra, and demonstrated that this difficulty cannot be overcome without changing, at least locally, one of the two. Third, we proposed a logic-based type language to represent the missing information and discussed how to combine it with the existing one.

The last step of this work illustrates the design of type-checking systems based on rich tree logics, where subtyping is typically decided using a logical satisfiability solver (such as the one proposed in [12]). One advantage of this approach is to inherit from the well-studied characteristics of a generic logical formalism (in particular: expressivity, worst-case complexity, robustness with respect to extensions), as well as from the availability of implementations.

One direction for further research would be to do backward type inference [21], i.e. construct types that subexpressions must match from the required type of the whole expression, then check whether the initial typing environment matches the requirements. The logic-based type language seems natural for this purpose since it provides a succinct notation for the inferred subtypes and their contexts.

Another direction of research is to extend XQuery's type algebra for higher order functions and parametric polymorphism in the light of recent breakthroughs on the subject [13, 6]. All these features combined will represent a much more powerful type system usable in practice and more attractive than W3C's proposed one, which gives too many false negatives. This will also revive programmers' interest in the language as the range of detectable errors in XQuery compilers will become much more accurate and wider.

## 8. REFERENCES

- [1] M. Benedikt and H. Vu. Higherorder functions and structured datatypes. In *WebDB: Proceedings of the 15th International Workshop on the Web and Databases*, pages 43–48, 2012.
- [2] V. Benzaken, G. Castagna, and A. Frisch. CDuce: an XML-centric general-purpose language. In *Proceedings of the eighth ACM SIGPLAN international conference on Functional programming*, ICFP '03, pages 51–63, New York, NY, USA, 2003. ACM.
- [3] V. Benzaken, G. Castagna, H. Hosoya, B. C. Pierce, and S. Vansummeren. XML typechecking. In L. Liu and M. T. Özsu, editors, *Encyclopedia of Database Systems*, pages 3646–3650. Springer US, 2009.
- [4] S. Boag, D. Chamberlin, M. Fernández, D. Florescu, J. Robie, and J. Siméon. XQuery 1.0: An XML Query Language (Second edition). W3C Recommendation, 2010. <http://www.w3.org/TR/xquery/>.
- [5] G. Castagna, H. Im, K. Nguyen, and V. Benzaken. A core calculus for XQuery 3.0. <http://www.pps.univ-paris-diderot.fr/~gc/papers/xqueryduce.pdf>, apr 2013.
- [6] G. Castagna and Z. Xu. Set-theoretic foundation of parametric polymorphism and subtyping. In *ICFP'11, Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming*, pages 94–106. ACM, sep 2011.
- [7] D. Colazzo, G. Ghelli, P. Manghi, and C. Sartiani. Types for path correctness of XML queries. In C. Okasaki and K. Fisher, editors, *ICFP*, pages 126–137. ACM, 2004.
- [8] D. Colazzo and C. Sartiani. Precision and complexity of XQuery type inference. In *Proceedings of the 13th International ACM SIGPLAN Conference on Principles and Practice of Declarative Programming, July 20-22, 2011, Odense, Denmark*, pages 89–100, 2011.
- [9] D. Draper, M. Dyck, P. Fankhauser, M. Fernández, A. Malhotra, K. Rose, M. Rys, J. Siméon, and P. Wadler. XQuery 1.0 and XPath 2.0 formal semantics, W3C recommendation, December 2010. <http://www.w3.org/TR/xquery-semantics/>.
- [10] M. F. Fernández, J. Siméon, and P. Wadler. A semi-monad for semi-structured data. In J. V. den Bussche and V. Vianu, editors, *ICDT*, volume 1973 of *Lecture Notes in Computer Science*, pages 263–300. Springer, 2001.
- [11] P. Genevès, N. Layaïda, and V. Quint. Impact of XML schema evolution. *ACM Trans. Internet Technol.*, 11:4:1–4:27, July 2011.
- [12] P. Genevès, N. Layaïda, and A. Schmitt. Efficient static analysis of XML paths and types. In *Proceedings of the 2007 ACM SIGPLAN conference on Programming language design and implementation, PLDI '07*, pages 342–351, New York, NY, USA, 2007. ACM.
- [13] N. Gesbert, P. Genevès, and N. Layaïda. Parametric polymorphism and semantic subtyping: the logical connection. In *ICFP'11, Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming*, pages 107–116. ACM, sep 2011.
- [14] H. Hosoya and B. C. Pierce. XDuce: A statically typed XML processing language. *ACM Trans. Internet Technol.*, 3(2):117–148, May 2003.
- [15] H. Hosoya, J. Vouillon, and B. C. Pierce. Regular expression types for XML. *ACM Trans. Program. Lang. Syst.*, 27(1):46–90, 2005.
- [16] G. P. Huet. The zipper. *J. Funct. Program.*, 7(5):549–554, 1997.
- [17] C. Kirkegaard, A. Møller, and M. I. Schwartzbach. Static analysis of XML transformations in Java. *IEEE Transactions on Software Engineering*, 30(3):181–192, March 2004.
- [18] C. Koch. On the complexity of nonrecursive XQuery and functional query languages on complex values. *ACM Trans. Database Syst.*, 31(4):1215–1256, Dec. 2006.
- [19] S. Maneth, A. Berlea, T. Perst, and H. Seidl. XML type checking with macro tree transducers. In *Proceedings of the twenty-fourth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems, PODS '05*, pages 283–294, New York, NY, USA, 2005. ACM.
- [20] S. Maneth, T. Perst, and H. Seidl. Exact XML type checking in polynomial time. In T. Schwentick and D. Suciú, editors, *ICDT*, volume 4353 of *Lecture Notes in Computer Science*, pages 254–268. Springer, 2007.
- [21] T. Milo, D. Suciú, and V. Vianu. Typechecking for XML transformers. *J. Comput. Syst. Sci.*, 66(1):66–97, 2003.
- [22] J. Robie, D. Chamberlin, M. Dyck, D. Florescu, J. Melton, and J. Siméon. XQuery update facility 1.0, W3C recommendation, March 2011. <http://www.w3.org/TR/xquery-update-10/>.
- [23] J. Robie, D. Chamberlin, M. Dyck, and J. Snelson. XQuery 3.0: An XML Query Language. W3C Last Call Working, July 2013. <http://www.w3.org/TR/xquery-30/>.
- [24] P. Wadler. XQuery: A typed functional language for querying XML. In J. Jeuring and S. L. P. Jones, editors, *Advanced Functional Programming*, volume 2638 of *Lecture Notes in Computer Science*, pages 188–212. Springer, 2002.