

Precision and Complexity of XQuery Type Inference

Dario Colazzo*

LRI - Université Paris Sud/INRIA

Carlo Sartiani

Dipartimento di Matematica e Informatica - Università della Basilicata

1 Introduction

XQuery [6] is a functional and Turing-complete XML data manipulation language that allows the programmer to navigate through an XML document (or a set of documents), to select relevant fragments of the document, and to combine them so to return new documents. Defined by the W3C Consortium, XQuery is also statically and strongly typed.

In XQuery, types of input data and functions are defined in terms of regular expression types, but it is quite easy to write queries that generate non-regular languages. As a consequence, any type system for XQuery has to rely on a *type inference* process that approximates the (possibly non-regular) output type of a query with a regular type. This approximation process, while mandatory and unavoidable, may significantly decrease the precision of the inferred types. This is the case of the W3C proposed type system, which relies on some over-approximating rules for expressions widely used in practice, like `for-iteration`, for instance. Another source of undesired over-approximation is given by rules to type horizontal and upward XPath axes, for which the type *any* is always inferred.

An alternative and more precise approach for typing XQuery has been proposed in [4] and used as a basis for other proposals [2, 1, 3]. This type system has more precise type inference, at the price of an increased computational complexity.

Though these approaches are well-known by the database and programming language communities, a formal, rigorous, and complete analysis showing in which cases the two proposals differ in terms of precision and complexity for type inference, is still missing. Such formal analysis, guided by existing classifications of schemas used in practice, could have a practical

relevance as well, since it would provide important information to implementation designers.

In this paper we start this missing comparative analysis. Besides providing a clean and simple formalization of main typing mechanisms, we will formally study the complexity of the two approaches, show in which cases the W3C excessively over-approximates inferred types, and show that in many of these cases the type system proposed in [4] infers quite precise types still at a reasonable time cost.

2 Type Language

Our type language, which is an abstract model of the type language of XQuery, describes forests of *unranked, node-labeled* trees. Forests and trees obey the data model shown in following grammar, where $()$ is the empty forest, b identifies atomic values, $l[f]$ represents an element labeled by l with the nodes in f as its children, and f, f denotes the ordered concatenation of forests.

$$f ::= () \mid b \mid l[f] \mid f, f$$

As usual, the concatenation operator is associative and $()$ is its neutral element: $() , f = f, () = f$.

Our type language is shown below, where $()$ is the type for the empty sequence value, B denotes the type for base values, types T, U and $T \mid U$ are, respectively, ordered product and union types, and, finally, T^* , T^+ , and $T^?$ are the usual repetition and optional types. Vertical recursion is supported through *type environments* and type variables.

$$T ::= \begin{array}{c} () \\ \mid T \mid T \\ \mid X \end{array} \mid \begin{array}{c} B \\ T^* \end{array} \mid \begin{array}{c} l[T] \\ T^+ \end{array} \mid \begin{array}{c} T, T \\ T^? \end{array}$$

$$B ::= \text{String}$$

*This work has been partially funded by *Agence Nationale de la Recherche*, decision ANR-08-DEFIS-004.

$$\begin{aligned}
Q & ::= () \mid b \mid l[Q] \mid Q, Q \\
& \quad \mid \bar{x} \text{ child} :: \text{NodeTest} \\
& \quad \mid \bar{x} \text{ dos} :: \text{NodeTest} \\
& \quad \mid \text{for } \bar{x} \text{ in } Q \text{ return } Q \\
& \quad \mid \text{let } x ::= Q \text{ return } Q \\
& \quad \mid \text{for } \bar{x} \text{ in } Q \text{ where } P \text{ return } Q \\
& \quad \mid \text{let } x ::= Q \text{ where } P \text{ return } Q \\
\text{NodeTest} & ::= 1 \mid \text{node}() \mid \text{text}() \\
\\
P & ::= \text{true} \mid \chi \delta \chi \mid \text{empty}(\chi) \mid P \text{ or } P \\
& \quad \mid \text{not } P \mid (P) \\
\chi & ::= \bar{x} \mid x \\
\delta & ::= = \mid <
\end{aligned}$$

Figure 1: The grammar of miniXQuery.

As usual, we restrict to $l[]$ -guarded type environments, which are environments where only $l[]$ -guarded vertical recursion is allowed.

The lack of horizontal recursion is counterbalanced by the presence of the Kleene star operator $*$. This restriction is canonical, and makes the type language as expressive as regular tree languages, hence expressive enough to capture the main type mechanisms of DTD and XML Schema.

As usual the semantics of type is defined as the minimal function that satisfies the following set of monotone equations:

$$\begin{aligned}
\llbracket () \rrbracket_E & \triangleq \{ () \} \\
\llbracket B \rrbracket_E & \triangleq \{ b \mid b \text{ is a base value} \} \\
\llbracket l[T] \rrbracket_E & \triangleq \{ l[f] \mid f \in \llbracket T \rrbracket_E \} \\
\llbracket T_1 \mid T_2 \rrbracket_E & \triangleq \llbracket T_1 \rrbracket_E \cup \llbracket T_2 \rrbracket_E \\
\llbracket T_1, T_2 \rrbracket_E & \triangleq \{ f_1, f_2 \mid f_i \in \llbracket T_i \rrbracket_E \} \\
\llbracket X \rrbracket_E & \triangleq \llbracket X(E) \rrbracket_E \\
\llbracket T? \rrbracket_E & \triangleq \llbracket T \mid () \rrbracket_E \\
\llbracket T+ \rrbracket_E & \triangleq \llbracket T \rrbracket_E^+ \\
\llbracket T* \rrbracket_E & \triangleq \llbracket T \rrbracket_E^*
\end{aligned}$$

3 miniXQuery

miniXQuery is a minimal language modeling the FLWR core of XQuery. miniXQuery contains **for**, **let**, **where**, and **return** clauses, and allows the user to specify both the child and the descendants-or-self axes. The predicate language comprises variable comparisons only. The syntax of miniXQuery is shown in Figure 1. The semantics of the language and the required auxiliary functions are shown in Table 3.1 and Figure 2.

There, ρ is a substitution assigning a forest to each free variable in the query; we

$$\begin{aligned}
\text{dos}(b) & \triangleq b & \text{childr}(b) & \triangleq () \\
\text{dos}(l[f]) & \triangleq l[f], \text{dos}(f) & \text{childr}(l[f]) & \triangleq f \\
\text{dos}() & \triangleq () & \text{dos}(f, f') & \triangleq \text{dos}(f), \text{dos}(f') \\
b :: l & \triangleq () & l[f] :: l & \triangleq l[f] \\
() :: l & \triangleq () & (f, f') :: l & \triangleq f :: l, f' :: l \\
m[f] :: l & \triangleq () \quad m \neq l & f :: \text{node}() & \triangleq f \\
b :: \text{text}() & \triangleq b & () :: \text{text}() & \triangleq () \\
m[f] :: \text{text}() & \triangleq () \\
(f, f') :: \text{text}() & \triangleq f :: \text{text}(), f' :: \text{text}() \\
\text{true}(\rho) & \triangleq \text{true} \\
(\chi \delta \chi)(\rho) & \triangleq \exists t \in \text{trees}(\rho(\chi)), t' \in \text{trees}(\rho(\chi')). t \delta t' \\
(P \text{ or } P)(\rho) & \triangleq P(\rho) \text{ or } P(\rho) \\
\text{empty}((\chi))(\rho) & \triangleq \text{if } \rho(\chi) = () \text{ then true else false} \\
(\text{not } P)(\rho) & \triangleq \text{not } P(\rho)
\end{aligned}$$

Figure 2: Auxiliary functions and predicate evaluation.

make the assumption that each ρ always associates a tree t to a for-variable \bar{x} it defines; also, **dos** is a shortcut for descendant-or-self. The semantics of **for** queries is defined via the operator $\prod_{t \in \text{trees}(f)} A(t)$, yielding the forest $A(t_1), \dots, A(t_n)$, when $f = t_1, \dots, t_n$, and $()$ when $f = ()$.

4 The W3C Type System

4.1 Complexity

The first result we want to discuss here is stated in the following theorem.

Theorem 4.1 *Given an input type T and a miniXQuery query Q , $E; \Gamma \vdash_m Q : U$ can be computed, according to the W3C type system, in polynomial time and space.*

For reasons of space, we cannot prove here this result. However, we can illustrate it by looking at some of the most crucial type rules of the W3C type system.

Consider the rule (TYPEFOR) shown below. (TYPEFOR)

$$\frac{
\begin{array}{c}
E; \Gamma \vdash_m Q_1 : T_1 \\
E; \Gamma, \bar{x} : \text{Prime}_E(T_1) \vdash_m Q_2 : T_2
\end{array}
}{
E; \Gamma \vdash_m \text{for } \bar{x} \text{ in } Q_1 \text{ return } Q_2 : T_2 \cdot \text{Quant}_E(T_1)
}$$

The rule describes the behaviour of the type inference system when a **for**-iteration is visited. The output type is computed as follows. The rule first computes the inferred type for Q_1 ; from this type a *prime* type is extracted by the function $\text{Prime}_E(T_1)$. $\text{Prime}_E(T_1)$ returns the union of the uppermost base or tree

Table 3.1. *miniXQuery semantics*

$\llbracket b \rrbracket_\rho$	$\triangleq b$	$\llbracket x \rrbracket_\rho$	$\triangleq \rho(x)$	$\llbracket \bar{x} \rrbracket_\rho$	$\triangleq \rho(\bar{x})$
$\llbracket () \rrbracket_\rho$	$\triangleq ()$	$\llbracket Q_1, Q_2 \rrbracket_\rho$	$\triangleq \llbracket Q_1 \rrbracket_\rho, \llbracket Q_2 \rrbracket_\rho$	$\llbracket l[Q] \rrbracket_\rho$	$\triangleq l[\llbracket Q \rrbracket_\rho]$
$\llbracket \bar{x} \text{ child} :: \text{NodeTest} \rrbracket_\rho$	$\triangleq \text{childr}(\llbracket \bar{x} \rrbracket_\rho) :: \text{NodeTest}$				
$\llbracket \bar{x} \text{ dos} :: \text{NodeTest} \rrbracket_\rho$	$\triangleq \text{dos}(\llbracket \bar{x} \rrbracket_\rho) :: \text{NodeTest}$				
$\llbracket \text{let } x ::= Q_1 \text{ return } Q_2 \rrbracket_\rho$	$\triangleq \llbracket Q_2 \rrbracket_{\rho, x \mapsto \llbracket Q_1 \rrbracket_\rho}$				
$\llbracket \text{for } \bar{x} \text{ in } Q_1 \text{ return } Q_2 \rrbracket_\rho$	$\triangleq \prod_{t \in \text{trees}(\llbracket Q_1 \rrbracket_\rho)} \llbracket Q_2 \rrbracket_{\rho, \bar{x} \mapsto t}$				
$\llbracket \text{let } x ::= Q_1 \text{ where } P \text{ return } Q_2 \rrbracket_\rho$	$\triangleq \text{if } P(\rho, x \mapsto \llbracket Q_1 \rrbracket_\rho) \text{ then } \llbracket Q_2 \rrbracket_{\rho, x \mapsto \llbracket Q_1 \rrbracket_\rho} \text{ else } ()$				
$\llbracket \text{for } \bar{x} \text{ in } Q_1 \text{ where } P \text{ return } Q_2 \rrbracket_\rho$	$\triangleq \prod_{t \in \text{trees}(\llbracket Q_1 \rrbracket_\rho)} (\text{if } P(\rho, \bar{x} \mapsto t) \text{ then } \llbracket Q_2 \rrbracket_{\rho, \bar{x} \mapsto t} \text{ else } ())$				

types inside T and it can be computed in linear space and time. $\text{Prime}_E(T_1)$ is, then, bound to \bar{x} in the variable environment and used to infer the output type T_2 , which is further refined by the application of a quantifier in $\{?, +, *\}$ ($T_2 \cdot \text{Quant}_E(T_1)$).

It should be observed that Q_2 is visited only once by the inference rule, hence the complexity is not affected by the size of the inferred type for Q_1 .

Consider now the rule (TYPEDOS) shown below.

$$\begin{array}{c}
 \text{(TYPEDOS)} \\
 \hline
 \begin{array}{l}
 \bar{x} : T \in \Gamma \wedge (T = m_1[T'_1] \mid m_2[T'_2] \mid \dots \mid m_n[T'_n] \vee T = B) \\
 E; \Gamma \vdash_m \bar{x} \text{ child} : T_1 \\
 E; \Gamma, \bar{x} : \text{Prime}_E(T_1) \vdash_m \bar{x} \text{ child} : T_2 \\
 E; \Gamma, \bar{x} : \text{Prime}_E(T_2) \vdash_m \bar{x} \text{ child} : T_3 \\
 \dots \\
 E; \Gamma, \bar{x} : \text{Prime}_E(T_n) \vdash_m \bar{x} \text{ child} : T_{n+1} \\
 E; \Gamma \vdash_m \text{Prime}_E(T_{n+1}) <: \text{Prime}_E(T_1) \mid \dots \mid \text{Prime}_E(T_n) \\
 U' = (\text{Prime}_E(T) \mid \text{Prime}_E(T_1) \mid \dots \mid \text{Prime}_E(T_n)) * \\
 E \vdash_m U' \stackrel{\text{NodeTest}}{=} U
 \end{array} \\
 \hline
 E; \Gamma \vdash_m \bar{x} \text{ dos} :: \text{NodeTest} : U
 \end{array}$$

This rule applies to **dos** selectors and infers a type for a $\bar{x} \text{ dos} :: \text{NodeTest}$ filter. The rule essentially traverses the parse tree of T and, at each step, collects the prime types it encounters. The premise $E; \Gamma \vdash_m \text{Prime}_E(T_{n+1}) <: \text{Prime}_E(T_1) \mid \dots \mid \text{Prime}_E(T_n)$ is just a formal way to express the termination of the search in the parse tree and it does not involve any subtyping operation. The rule returns a type consisting of the star-guarded union of the types collected during the exploration.

It is quite easy to see that no type is visited twice and that the repeated application of $\text{Prime}_E(\cdot)$ can be computed in polynomial time by using some form of memoization.

4.2 Precision

To illustrate the precision issues that affects the W3C type system, we focus our attention on the following type rule.

(TYPECHILDNODETEST)

$$\frac{\begin{array}{l} \bar{x} : T \in \Gamma \wedge (T = m[T'] \vee T = B) \\ E \vdash_m \text{content}(T) \stackrel{\text{NodeTest}}{=} U \end{array}}{E; \Gamma \vdash_m \bar{x} \text{ child} :: \text{NodeTest} : U}$$

This rule infers an output type for a $\bar{x} \text{ child} :: \text{NodeTest}$ filter by extracting the children types of T and filtering them according to NodeTest .

Consider now the following query:

```

for  $\bar{x}$  in  $l[b], m[b], n[b]$ 
return for  $\bar{y}$  in  $\bar{x} \text{ child} :: \text{node}()$ 
  return  $\bar{x}, \bar{y}$ 

```

When invoked on this query, the type inference system first infers a type for $l[b], m[b], n[b]$. This type ($l[B], m[B], n[B]$) is then passed to $\text{Prime}_E(\cdot)$, whose result is $(l[B] \mid m[B] \mid n[B])$.

The inference system, then, infers a type for the inner query. In this case, the rule (TYPECHILDNODETEST) is applied and the result is just $B \mid B \mid B$.

Hence, the inferred output type for the query is: $((l[B] \mid m[B] \mid n[B]), (B)) +$.

This type is a gross over-approximation of the actual type of the output, which is the following: $l[B], B, m[B], B, n[B], B$. This approximation is justified by the need to confine space and time complexity in the polynomial realm.

In the following section we will see how precision of type inference can be improved.

5 A More precise type system

The more precise type system proposed in [4] distinguishes for its rule to type **for**-iterations:

(TYPEFOR)

$$\frac{E; \Gamma \vdash Q_1 : T_1 \quad E; \Gamma \vdash \bar{x} \text{ in } T_1 \rightarrow Q_2 : T_2}{E; \Gamma \vdash \text{for } \bar{x} \text{ in } Q_1 \text{ return } Q_2 : T_2}$$

As in the W3C case, the rule first computes the inferred type for Q_1 ; this type is then iterated on by a type analysis allowing the system to prove the premise $E; \Gamma \vdash \bar{x} \text{ in } T_1 \rightarrow Q_2 : T_2$. This analysis is performed according to several rules; for reasons of space, we only show below the most significant ones.

(TYPEINSEQ)

$$\frac{E; \Gamma \vdash \bar{x} \text{ in } T_i \rightarrow Q : U_i \quad i = 1, 2}{E; \Gamma \vdash \bar{x} \text{ in } T_1, T_2 \rightarrow Q : U_1, U_2}$$

(TYPEINSTAR)

$$\frac{E; \Gamma \vdash \bar{x} \text{ in } T \rightarrow Q : U}{E; \Gamma \vdash \bar{x} \text{ in } T* \rightarrow Q : U*}$$

(TYPEINEL)

$$\frac{E; \Gamma \vdash Q : U}{E; \Gamma \vdash \bar{x} \text{ in } l[T] \rightarrow Q : l[U]}$$

The rules for auxiliary expressions $\bar{x} \text{ in } T_1 \rightarrow Q_2$ are purely structural on T_1 , and allow for a very precise type analysis, since possible over-approximations introduced by $Quant_E(T_1)$ are avoided. For example, when T_1 is a product T', T'' , we have that, according to the W3C specification, $Quant_E(T_1)$ is $*$ (or $+$ in some cases). As a consequence, a $*$ type is always inferred for the whole **for** query. Instead, the above rule TYPEINSEQ infers a more precise sequence type T'_2, T''_2 . In general a $*$ -type is inferred much less frequently, typically only when T_1 contains a top level $*$ type. As a consequence, we have a more precise type inference.

On the other hand, for nested **for**-expressions this rule may lead to $O(n^k)$ time complexity, where k is the number of nested **for** expressions, and n is the maximal size of the type T_1^i inferred for the left hand side query of **for**-expressions. In practical cases, both k and n are rather small, hence the type inference time remains reasonable. This is testified by some tests that we made on real case and quite complex queries [5].

Concerning the typing of descendant-or-self axis, the type system proposed in [4] essentially adopts the same approach as the one proposed by the W3C. However, for non-recursive types, we can do better, as shown in [5]. Given a non-recursive type T (w.l.o.g. we can assume T not containing variables), we can define type inference for the descendant-or-self axis as follows:

Definition 5.1

$$\begin{aligned} dos(()) &= () \\ dos(B) &= B \\ dos(l[T]) &= l[T], dos(T) \\ dos(T?) &= dos(T)? \\ dos(T*) &= dos(T)* \\ dos(T+) &= dos(T)+ \\ dos(T, U) &= dos(T), dos(U) \\ dos(T | U) &= dos(T) | dos(U) \end{aligned}$$

It is straightforward to prove that for any non-recursive type T :

$$dos(T) = U \Rightarrow \forall f \in \llbracket T \rrbracket. dos(f) \in \llbracket U \rrbracket$$

As it can be seen, the inference of $*$ types is reduced as much as possible, and this leads to a better precision wrt to the W3C approach.

6 Conclusions

This work is a first step towards a comparative analysis of the two current major proposals for XQuery type inference. In future work we aim at formally characterizing classes of queries and schemas for which the two systems differs in terms of precision. At the same time we will provide classes for which the two systems ensure the same degree of precision.

References

- [1] V. Benzaken, G. Castagna, D. Colazzo, and K. Nguyen. Type-based XML projection. In *VLDB*, 2006.
- [2] J. Cheney. Flux: functional updates for XML. In *ICFP*, pages 3–14, 2008.
- [3] J. Cheney. Regular expression subtyping for XML query and update languages. In *ESOP*, 2008.
- [4] D. Colazzo, G. Ghelli, P. Manghi, and C. Sartiani. Types for Path Correctness of XML Queries. In *ICFP*, 2004.
- [5] D. Colazzo and C. Sartiani. Detection of corrupted schema mappings in xml data integration systems. *ACM Trans. Internet Techn.*, 9(4), 2009.
- [6] D. Draper, P. Fankhauser, M. Fernandez, A. Malhotra, K. Rose, M. Rys, J. Siméon, and P. Wadler. XQuery 1.0 and XPath 2.0 Formal Semantics. Technical report, World Wide Web Consortium, Jan. 2007. W3C Recommendation.