

A Declarative Embedding of XQuery in a Functional-Logic Language*

Jesús M. Almendros-Jiménez¹, Rafael Caballero²,
Yolanda García-Ruiz², and Fernando Sáenz-Pérez³

¹ Dpto. Lenguajes y Computación, Universidad de Almería, Spain

² Dpto. de Sistemas Informáticos y Computación, UCM, Spain

³ Dpto. de Ingeniería del Software e Inteligencia Artificial, UCM, Spain
jalmen@ual.es, {rafa,fernan}@sip.ucm.es , ygarciar@fdi.ucm.es

Abstract. This paper addresses the problem of integrating a fragment of XQuery, a language for querying XML documents, into the functional-logic language \mathcal{TCY} . The queries are evaluated by an interpreter, and the declarative nature of the proposal allows us to prove correctness and completeness with respect to the semantics of the subset of XQuery considered. The different fragments of XML that can be produced by XQuery expressions are obtained using the non-deterministic features of functional-logic languages. As an application of this proposal we show how the typical *generate and test* techniques of logic languages can be used for generating test-cases for XQuery expressions.

1 Introduction

XQuery has been defined as a query language for finding and extracting information from XML [15] documents. Originally designed to meet the challenges of large-scale electronic publishing, XML also plays an important role in the exchange of a wide variety of data on the Web and elsewhere. For this reason many modern languages include libraries or encodings of XQuery, including logic programming [1] and functional programming [6]. In this paper we consider the introduction of a simple subset of XQuery [18,20] into the functional-logic language \mathcal{TCY} [11].

One of the key aspects of declarative languages is the emphasis they pose on the logic semantics underpinning declarative computations. This is important for reasoning about computations, proving properties of the programs or applying declarative techniques such as abstract interpretation, partial evaluation or algorithmic debugging [14]. There are two different declarative alternatives that can be chosen for incorporating XML into a (declarative) language:

1. Use a domain-specific language and take advantage of the specific features of the host language. This is the approach taken in [9], where a rule-based

* Work partially supported by the Spanish projects STAMP TIN2008-06622-C03-01, DECLARAWEB TIN2008-06622-C03-03, Prometidos-CM S2009TIC-1465 and GPD UCM-BSCH-GR58/08-910502.

language for processing semi-structured data that is implemented and embedded into the functional-logic language Curry, and also in [13] for the case of logic programming.

2. Consider an existing query language such as XQuery, and embed a fragment of the language in the host language, in this case \mathcal{TOY} . This is the approach considered in this paper.

Thus, our goal is to include XQuery using the purely declarative features of the host languages. This allows us to prove that the semantics of the fragment of XQuery has been correctly included in \mathcal{TOY} . To the best of our knowledge, it is the first time a fragment of XQuery has been encoded in a functional-logic language. A first step in this direction was proposed in [5], where XPath [16] expressions were introduced in \mathcal{TOY} . XPath is a subset of XQuery that allows navigating and returning fragments of documents in a similar way as the path expressions used in the *chdir* command of many operating systems. The contributions of this paper with respect to [5] are:

- The setting has been extended to deal with a simple fragment of XQuery, including *for* statements for traversing XML sequences, *if/where* conditions, and the possibility of returning XML elements as results. Some basic XQuery constructions such as *let* statements are not considered, but we think that the proposal is powerful enough for representing many interesting queries.
- The soundness of the approach is formally proved, checking that the semantics of the fragment of XQuery is correctly represented in \mathcal{TOY} .

Next section introduces the fragment of XQuery considered and a suitable operational semantics for evaluating queries. The language \mathcal{TOY} and its semantics are presented in Section 3. Section 4 includes the interpreter that performs the evaluation of simple XQuery expressions in \mathcal{TOY} . The theoretical results establishing the soundness of the approach with respect to the operational semantics of Section 2 are presented in Section 4.1. Section 5 explains the automatic generation of test cases for simple XQuery expressions. Finally, Section 6 concludes summarizing the results and proposing future work.

An extended version of the paper including proofs of the theoretical results can be found at [2].

2 XQuery and Its Operational Semantics

XQuery allows the user to query several documents, applying join conditions, generating new XML fragments, and using many other features [18,20]. The syntax and semantics of the language are quite complex [19], and thus only a small subset of the language is usually considered. The next subsection introduces the fragment of XQuery considered in this paper.

2.1 The Subset SXQ

In [4] a declarative subset of XQuery, called XQ, is presented. This subset is a core language for XQuery expressions consisting of *for*, *let* and *where/if* statements.

```

query ::= ( ) | query query | tag
        | doc(File) | doc(File)/axis ::  $\nu$  | var | var/axis ::  $\nu$ 
        | for var in query return query
        | if cond then query
cond ::= var=var | query
tag ::= <a> var ... var </a> | <a> tag </a>

```

Fig. 1. Syntax of SXQ, a simplified version of XQ

In this paper we consider a simplified version of XQ which we call SXQ and whose syntax can be found in Figure 1. where *axis* can be one of *child*, *self*, *descendant* or *dos* (i.e. descendant or self), and ν is a node test. The differences of SXQ with respect to XQ are:

1. XQ includes the possibility of using variables as tag names using a constructor *lab(\$x)*.
2. XQ permits enclosing any query *Q* between tag labels *<a>Q*. SXQ only admits either variables or other tags inside a tag.

Our setting can be easily extended to support the *lab(\$x)* feature, but we omit this case for the sake of simplicity in this presentation. The second restriction is more severe: although *lets* are not part of XQ, they could be simulated using *for* statements inside tags. In our case, forbidding other queries different from variables inside tag structures imply that our core language cannot represent *let* expressions. This limitation is due to the non-deterministic essence of our embedding, since a *let* expression means collecting all the results of a query instead of producing them separately using non-determinism. In spite of these limitations, the language SXQ is still useful for solving many common queries as the following example shows.

Example 1. Consider an XML file “*bib.xml*” containing data about books, and another file “*reviews.xml*” containing reviews for some of these books (see [17], sample data 1.1.2 and 1.1.4 to check the structure of these documents and an example). Then we can list the reviews corresponding to books in “*bib.xml*” as follows:

```

for $b in doc("bib.xml")/bib/book,
    $r in doc("reviews.xml")/reviews/entry
where $b/title = $r/title
for $booktitle in $r/title,
    $revtext in $r/review
return <rev> $booktitle $revtext </rev>

```

The variable *\$b* takes the value of the different books, and *\$r* the different reviews. The *where* condition ensures that only reviews corresponding to the book are considered. Finally, the last two variables are only employed to obtain

the book title and the text of the review, the two values that are returned as output of the query by the *return* statement.

It can be argued that the code of this example does not follow the syntax of Figure 1. While this is true, it is very easy to define an algorithm that converts a query formed by *for*, *where* and *return* statements into a SXQ query (as long as it only includes variables inside tags, as stated above). The idea is simply to convert the *where* into *ifs*, following each *for* by a *return*, and decomposing XPath expressions including several steps into several *for* expressions by introducing a new auxiliary variable and each one consisting of a single step.

Example 2. The query of Example 1 using SXQ syntax:

```

for $x1 in doc("bib.xml")/child::bib return
for $x2 in $x1/child::book return
for $x3 in doc("reviews.xml")/child::reviews return
for $x4 in $x3/entry return
if ($x2/title = $x4/title) then
  for $x5 in $x4/title return
    for $x6 in $x4/review return <rev> $x5 $x6 </rev>

```

We end this subsection with a few definitions that are useful for the rest of the paper. The set of variables in a query Q is represented as $Var(Q)$. Given a query Q , we use the notation $Q|_p$ for representing the subquery Q' that can be found in Q at position p . Positions are defined as usual in syntax trees:

Definition 1. Given a query Q and a position p , $Q|_p$ is defined as follows:

$$\begin{aligned}
Q|_{\varepsilon} &= Q \\
(Q_1 Q_2)|_{(i.p)} &= (Q_i)|_p \quad i \in \{1, 2\} \\
(\text{for var in } Q_1 \text{ return } Q_2)|_{(i.p)} &= (Q_i)|_p \quad i \in \{1, 2\} \\
(\text{if } Q_1 \text{ then } Q_2)|_{(i.p)} &= (Q_i)|_p \quad i \in \{1, 2\} \\
(\text{if var=var then } Q_1)|_{(1.p)} &= (Q_1)|_p
\end{aligned}$$

Hence the position of a subquery is the path in the syntax tree represented as the concatenation of children positions $p_1 \cdot p_2 \dots \cdot p_n$. For every position p , $\varepsilon \cdot p = p \cdot \varepsilon = p$. In general $Q|_p$ is not a proper SXQ query, since it can contain *free variables*, which are variables defined previously in *for* statements in Q . The set of variables of Q that are *relevant* for $Q|_p$ is the subset of $Var(Q)$ that can appear free in any subquery at position p . This set, denoted as $Rel(Q, p)$ is defined recursively as follows:

Definition 2. Given a query Q , and a position p , $Rel(Q, p)$ is defined as:

1. \emptyset , if $p = \varepsilon$.
2. $Rel(Q_1, p')$, if $Q \equiv Q_1 Q_2$, $p = 1 \cdot p'$.
3. $Rel(Q_2, p')$, if $Q \equiv Q_1 Q_2$, $p = 2 \cdot p'$.
4. $Rel(Q_1, p')$, if $Q \equiv \text{for var in } Q_1 \text{ return } Q_2$, $p = 1 \cdot p'$.
5. $\{\text{var}\} \cup Rel(Q_2, p')$, if $Q \equiv \text{for var in } Q_1 \text{ return } Q_2$, $p = 2 \cdot p'$.
6. $Rel(Q_1, p')$, if $Q \equiv \text{if } Q_1 \text{ then } Q_2$, $p = 1 \cdot p'$.

7. $Rel(Q_2, p')$, $\text{if } Q \equiv \text{if } Q_1 \text{ then } Q_2, p = 2 \cdot p'$.

Observe that cases $Q \equiv ()$, $Q \equiv \text{tag}$, $Q \equiv \text{var}$, $Q \equiv \text{var}/\chi :: \nu$, and $\text{var} = \text{var}$ correspond to $p \equiv \varepsilon$.

Without loss of generality we assume that all the relevant variables for a given position are indexed starting from 1 at the outer level. We also assume that every `for` statement introduces a new variable. A query like `for X in ((for Y in ...) (for Y in ...)) ...` is then renamed to an equivalent query of the form `for X1 in ((for X2 in ...) (for X3 in ...)) ...` (notice that the two `Y` variables occurred in different scopes).

2.2 XQ Operational Semantics

Figure 2 introduces the operational semantics of XQ that can be found in [4]. The only difference with respect to the semantics of this paper is that there is no rule for the constructor `lab`, for the sake of simplicity.

As explained in [4], the previous semantics defines the denotation of an XQ expression Q with k relevant variables, under a graph-like representation of a data forest \mathcal{F} , and a list of indexes \bar{e} in \mathcal{F} , denoted by $\llbracket Q \rrbracket_k(\mathcal{F}, \bar{e})$. In particular, each relevant variable $\$x_i$ of Q has as value the tree of \mathcal{F} indexed at position e_i . $\chi^{\mathcal{F}}(e_i, v)$ is a boolean function that returns *true* whenever v is the subtree of \mathcal{F} indexed at position e_i . The operator $\text{construct}(a, (\mathcal{F}, [w_1 \dots w_n]))$, denotes the construction of a new tree, where a is a label, \mathcal{F} is a data forest, and $[w_1 \dots w_n]$ is a list of nodes in \mathcal{F} . When applied, construct returns an indexed forest $(\mathcal{F} \cup T', [root(T')])$, where T' is a tree with domain a new set of nodes, whose root is labeled with a , and with the subtree rooted at the i -th (in sibling order) child of $root(T')$ being an isomorphic copy of the subtree rooted by w_i in \mathcal{F} . The symbol \uplus used in the rules takes two indexed forests (\mathcal{F}_1, l_1) , (\mathcal{F}_2, l_2) and returns an indexed forest $(\mathcal{F}_1 \cup \mathcal{F}_2, l)$, where $l = l_1 \cdot l_2$. Finally, $\text{tree}(e_i)$ denotes the maximal tree within the input forest that contains the node e_i , hence $\langle_{doc}^{tree(e_i)}$ is the document order on the tree containing e_i .

$$\begin{aligned}
\llbracket () \rrbracket_k(\mathcal{F}, \bar{e}) &= (\mathcal{F}, []) \\
\llbracket Q_1 Q_2 \rrbracket_k(\mathcal{F}, \bar{e}) &= \llbracket Q_1 \rrbracket_k(\mathcal{F}, \bar{e}) \uplus \llbracket Q_2 \rrbracket_k(\mathcal{F}, \bar{e}) \\
\llbracket \text{for } \$x_{k+1} \text{ in } Q_1 \text{ return } Q_2 \rrbracket_k(\mathcal{F}, \bar{e}) &= \text{let } (\mathcal{F}', \bar{l}) = \llbracket Q_1 \rrbracket_k(\mathcal{F}, \bar{e}) \text{ in} \\
&\quad \uplus_{1 \leq i \leq |\bar{l}|} \llbracket Q_2 \rrbracket_{k+1}(\mathcal{F}', \bar{e} \cdot l_i) \\
\llbracket \$x_i \rrbracket_k(\mathcal{F}, [e_1, \dots, e_k]) &= (\mathcal{F}, [e_i]) \\
\llbracket \$x_i / \chi :: \nu \rrbracket_k(\mathcal{F}, [e_1, \dots, e_k]) &= (\mathcal{F}, \text{list of nodes } v \text{ such that } \chi^{\mathcal{F}}(e_i, v) \text{ and} \\
&\quad \text{label name of } v = \nu \text{ in order } \langle_{doc}^{tree(e_i)}) \\
\llbracket \text{if } C \text{ then } Q_1 \rrbracket_k(\mathcal{F}, \bar{e}) &= \text{if } \pi_2(\llbracket C \rrbracket_k(\mathcal{F}, \bar{e})) \neq [] \text{ then } \llbracket Q_1 \rrbracket_k(\mathcal{F}, \bar{e}) \\
&\quad \text{else } (\mathcal{F}, []) \\
\llbracket \$x_i = \$x_j \rrbracket_k(\mathcal{F}, [e_1, \dots, e_k]) &= \text{if } e_i = e_j \text{ then } \text{construct}(\text{yes}, (\mathcal{F}, [])) \\
&\quad \text{else } (\mathcal{F}, [])
\end{aligned}$$

Fig. 2. Semantics of Core XQuery

Without loss of generality this semantics assumes that all the variables relevant for a subquery are numbered consecutively starting by 1 as in Example 2. It also assumes that the documents appear explicitly in the query. That is, in Example 2 we must suppose that instead of $doc("bib.xml")$ we have the XML corresponding to this document. Of course this is not feasible in practice, but simplifies the theoretical setting and it is assumed in the rest of the paper.

These semantic rules constitute a term rewriting system (TRS in short, see [3]), with each rule defining a single reduction step. The symbol $:=^*$ represents the reflexive and transitive closure of $:=$ as usual. The TRS is terminating and confluent (the rules are not overlapping). Normal forms have the shape $(\mathcal{F}, e_1, \dots, e_n)$ where \mathcal{F} is a forest of XML fragments, and e_i are nodes in \mathcal{F} , meaning that the query returns the XML fragments (**indexed by**) e_1, \dots, e_n . The semantics evaluates a query starting with the expression $\llbracket Q \rrbracket_0(\emptyset, ())$. Along intermediate steps, expressions of the form $\llbracket Q' \rrbracket_k(\mathcal{F}, \bar{e}_k)$ are obtained. The idea is that Q' is a subquery of Q with k relevant variables (which can occur free in Q'), that must take the values \bar{e}_k . The next result formalizes these ideas.

Proposition 1. *Let Q be a SXQ query. Suppose that*

$$\llbracket Q \rrbracket_0(\emptyset, ()) :=^* \llbracket Q' \rrbracket_n(\mathcal{F}, \bar{e}_n)$$

Then:

- Q' is a subquery of Q , that is, $Q' = Q_p$ for some p .
- $Rel(Q, p) = \{X_1, \dots, X_n\}$.
- Let S be the set of free variables in Q' . Then $S \subset Rel(Q, p)$.
- $\llbracket Q' \rrbracket_n(\mathcal{F}, \bar{e}_n) = \llbracket Q'\theta \rrbracket_0(\emptyset, ())$, with $\theta = \{X_1 \mapsto e_1, \dots, X_n \mapsto e_n\}$

Proof. Straightforward from Definition 2, and from the XQ semantic rules of Figure 2.

A more detailed discussion about this semantics and its properties can be found in [4].

3 \mathcal{TOY} and Its Semantics

A \mathcal{TOY} [11] program is composed of data type declarations, type alias, infix operators, function type declarations and defining rules for functions symbols. The syntax of *partial expressions* in \mathcal{TOY} $e \in Exp_{\perp}$ is $e ::= \perp \mid X \mid h \mid (e \ e')$ where X is a variable and h either a function symbol or a data constructor. Expressions of the form $(e \ e')$ stand for the application of expression e (acting as a function) to expression e' (acting as an argument). Similarly, the syntax of *partial patterns* $t \in Pat_{\perp} \subset Exp_{\perp}$ can be defined as $t ::= \perp \mid X \mid c \ t_1 \dots t_m \mid f \ t_1 \dots t_m$ where X represents a variable, c a data constructor of arity greater or equal to m , and f a function symbol of arity greater than m , being t_i partial patterns for all $1 \leq i \leq m$. Each rule for a function f in \mathcal{TOY} has the form:

$$\underbrace{f \ t_1 \dots t_n}_{\text{left-hand side}} \rightarrow \underbrace{r}_{\text{right-hand side}} \Leftarrow \underbrace{C_1, \dots, C_k}_{\text{condition}} \text{ where } \underbrace{s_1 = u_1, \dots, s_m = u_m}_{\text{local definitions}}$$

where u_i and r are expressions (that can contain new extra variables), C_j are strict equalities, and t_i, s_i are patterns. In \mathcal{TOY} , variable names must start with either an uppercase letter or an underscore (for anonymous variables), whereas other identifiers start with lowercase.

Data type declarations and type alias are useful for representing XML documents in \mathcal{TOY} :

```
data node      = txt      string
               | comment string
               | tag      string [attribute] [node]
data attribute = att      string string
type xml      = node
```

The data type `node` represents nodes in a simple XML document. It distinguishes three types of nodes: texts, tags (element nodes), and comments, each one represented by a suitable data constructor and with arguments representing the information about the node. For instance, constructor `tag` includes the tag name (an argument of type `string`) followed by a list of attributes, and finally a list of child nodes. The data type `attribute` contains the name of the attribute and its value (both of type `string`). The last type alias, `xml`, renames the data type `node`. Of course, this list is not exhaustive, since it misses several types of XML nodes, but it is enough for this presentation.

\mathcal{TOY} includes two primitives for loading and saving XML documents, called `load_xml_file` and `write_xml_file` respectively. For convenience all the documents are started with a dummy node `root`. This is useful for grouping several XML fragments. If the file contains only one node `N` at the outer level, the `root` node is unnecessary, and can be removed using this simple function:

```
load_doc F = N <== load_xml_file F == xmlTag "root" [] [N]
```

where `F` is the name of the file containing the document. Observe that the strict equality `==` in the condition forces the evaluation of `load_xml_file F` and succeeds if the result has the form `xmlTag "root" [] [N]` for some `N`. If this is the case, `N` is returned.

The constructor-based ReWriting Logic (CRWL) [7] has been proposed as a suitable declarative semantics for functional-logic programming with lazy non-deterministic functions. The calculus is defined by five inference rules (see Figure 3): (BT) that indicates that any expression can be approximated by bottom, (RR) that establishes the reflexivity over variables, the decomposition rule (DC), the (JN) (join) rule that indicates how to prove strict equalities, and the function application rule (FA). In every inference rule, $r, e_i, a_j \in Exp_{\perp}$ are partial expressions and $t, t_k \in Pat_{\perp}$ are partial patterns. The notation $[P]_{\perp}$ of the inference rule *FA* represents the set $\{(l \rightarrow r \Leftarrow C)\theta \mid (l \rightarrow r \Leftarrow C) \in P, \theta \in Subst_{\perp}\}$ of partial instances of the rules in program P ($Subst_{\perp}$ represents the set of partial

BT	$e \rightarrow \perp$	
RR	$X \rightarrow X$	with $X \in Var$
DC	$\frac{e_1 \rightarrow t_1 \dots e_m \rightarrow t_m}{h \bar{e}_m \rightarrow h \bar{t}_m}$	$h \bar{t}_m \in Pat_{\perp}$
JN	$\frac{e \rightarrow t \quad e' \rightarrow t}{e == e'}$	$t \in Pat$ (total pattern)
FA	$\frac{e_1 \rightarrow t_1 \dots e_n \rightarrow t_n \quad C \quad r \bar{a}_k \rightarrow t}{f \bar{e}_n \bar{a}_k \rightarrow t}$	if $(f \bar{t}_n \rightarrow r \Leftarrow C) \in [P]_{\perp}, t \neq \perp$

Fig. 3. CRWL Semantic Calculus

substitutions that replace variables by partial terms). The most complex inference rule is *FA* (Function Application), which formalizes the steps for computing a *partial pattern* t as approximation of a function call $f \bar{e}_n$:

1. Obtain partial patterns t_i as suitable approximations of the arguments e_i .
2. Apply a program rule $(f \bar{t}_n \rightarrow r \Leftarrow C) \in [P]_{\perp}$, verify the condition C , and check that t approximates the right-hand side r .

In this semantic notation, local declarations $a = b$ introduced in \mathcal{TOY} syntax by the reserved word **where** are part of the condition C as approximation statements of the form $b \rightarrow a$.

The semantics in \mathcal{TOY} allows introducing non-deterministic functions, such as the following function **member** that returns all the elements in a list:

```

member :: [A] -> A
member [X | Xs] = X
member [X | Xs] = member Xs

```

Another example of \mathcal{TOY} function is the definition of the infix operator $...::$ for XPath expressions (the operator $::$ in XPath syntax):

```

(...::) :: (A -> B) -> (B -> C) -> (A -> C)
(F ...:: G) X = G (F X)

```

As the examples show, \mathcal{TOY} is a typed language. However the type declaration is optional and in the rest of the paper they are omitted for the sake of simplicity. *Goals* in \mathcal{TOY} are sequences of strict equalities. A strict equality $e_1 == e_2$ holds (inference *JN*) if both e_1 and e_2 can be reduced to the same total pattern t . For instance, the goal `member [1,2,3,4] == R` yields four answers, the four values for **R** that make the equality true: $\{R \mapsto 1\}, \dots, \{R \mapsto 4\}$.

4 Transforming SXQ into \mathcal{TOY}

In order to represent SXQ queries in \mathcal{TOY} we use some auxiliary datatypes:

```
type xPath = xml-> xml

data sxq = xfor xml sxq sxq | xif cond sxq | xmlExp xml |
          xp path | comp sxq sxq
data cond = xml := xml | cond sxq
data path = var xml | xml :/ xPath | doc string xPath
```

The structure of the datatype `sxq` allows representing any SXQ query (see SXQ syntax in Figure 1). It is worth noticing that a variable introduced by a `for` statement has type `xml`, indicating that the variable always contains a value of this type. \mathcal{TOY} includes a primitive `parse_xquery` that translates any SXQ expression into its corresponding representation as a term of this datatype, as the next example shows:

Example 3. The translation of the SXQ query of Example 2 into the datatype `sxq` produces the following \mathcal{TOY} data term:

```
Toy> parse_xquery "for $x1 in doc(\"bib.xml\")/child::bib return
                 for $x2 in ..... <rev> $x5 $x6 </rev>" == R
yes
{R --> xfor X1 (xp (doc "bib.xml" (child ::. (nameT "bib"))))
  (xfor X2 (xp ( X1 :/ (child ::.(nameT "book"))))
  (xfor X3 (xp (doc "reviews.xml" (child ::. (nameT "reviews"))))
  (xfor X4 (xp ( X3 :/ (child ::.(nameT "entry"))))
  (xif ((xp(X2 :/ (child ::.(nameT "title")))) :=
        (xp(X4 :/ (child ::.(nameT "title"))))
  (xfor X5 (xp ( X4 :/ (child ::.(nameT "title"))))
  (xfor X6 (xp ( X4 :/ (child ::.(nameT "review"))))
    (xmlExp (xmlTag "rev" [] [X5,X6]))))))))
}
```

The interpreter assumes the existence of the infix operator `::.` that connects axes and tests to build steps, defined as the sequence of applications in Section 3.

The rules of the \mathcal{TOY} interpreter that processes SXQ queries can be found in Figure 4. The main function is `sxq`, which distinguishes cases depending of the form of the query. If it is an XPath expression then the auxiliary function `sxqPath` is used. If the query is an XML expression, the expression is just returned (this is safe thanks to our constraint of allowing only variables inside XML expressions). If we have two queries (`comp` construct), the result of evaluating any of them is returned using non-determinism. The `for` statement (`xfor` construct) forces the evaluation of the query `Q1` and binds the variable `X` to the result. Then the result query `Q2` is evaluated. The case of the `if` statement is analogous. The XPath subset considered includes tests for attributes (`attr`), label names (`nameT`), general elements (`nodeT`) and text nodes (`textT`). It also includes the axes `self`, `child`, `descendant` and `dos`. Observe that we do not

```

sxq (xp E)           = sxqPath E
sxq (xmlExp X)       = X
sxq (comp Q1 Q2)     = sxq Q1
sxq (comp Q1 Q2)     = sxq Q2
sxq (xfor X Q1 Q2)   = sxq Q2 <== X== sxq Q1
sxq (xif (Q1:=Q2) Q3) = sxq Q3 <== sxq Q1 == sxq Q2
sxq (xif (cond Q1) Q2) = sxq Q2 <== sxq Q1 == _

sxqPath (var X) = X
sxqPath (X :/ S) = S X
sxqPath (doc F S) = S (load_xml_file F)

%% XPATH %%
attr A (xmlTag S Attr L) = xmlText T <== member Attr == xmlAtt A T
nameT S (xmlTag S Attr L) = xmlTag S Attr L
nodeT X = X
textT (xmlText S) = xmlText S
commentT S (xmlComment S) = xmlComment S

self X = X
child (xmlTag _Name _Attr L) = member L
descendant X = child X
descendant X = descendant Y <== child X == Y
dos = self
dos = descendant

```

Fig. 4. \mathcal{TOY} transformation rules for SXQ

include reverse axes like `ancestor` because they can be replaced by expressions including forward axes, as shown in [12]. Other constructions such as filters can be easily included (see [5]). The next example uses the interpreter to obtain the answers for the query of our running example.

Example 4. The goal `sxq (parse_xquery "for...") == R` applies the interpreter of Figure 4 to the code of Example 2 (assuming that the string after `parse_xquery` is the query in Example 2), and returns the \mathcal{TOY} representation of the expected results:

```

<rev>
  <title>TCP/IP Illustrated</title>
  <review> One of the best books on TCP/IP. </review>
</rev>
...

```

4.1 Soundness of the Transformation

One of the goals of this paper is to ensure that the embedding is semantically correct and complete. This section introduces the theoretical results establishing these properties. If V is a set of indexed variables of the form $\{X_1, \dots, X_n\}$ we

use the notation $\theta(V)$ to indicate the sequence $\theta(X_1), \dots, \theta(X_n)$. In these results it is implicitly assumed that there is a bijective mapping f from XML format to the datatype `xml` in \mathcal{TOY} . Also, variables in XQuery $\$x_i$ are assumed to be represented in \mathcal{TOY} as X_i and conversely. However, in order to simplify the presentation, we omit the explicit mention to f and to f^{-1} .

Lemma 1. *Let P be a \mathcal{TOY} program, Q' an SXQ query, and Q, p such that $Q \equiv Q'_{|p}$. Define $V = \text{Rel}(Q', p)$ (see Definition 2), and $k = |V|$. Let θ be a substitution such that $\mathcal{P} \vdash (\text{sxq } Q\theta == t)$ for some pattern t . Then $\llbracket Q \rrbracket_k(\mathcal{F}, [\theta(V)]) :=^* (\mathcal{F}', L)$, for some forests $\mathcal{F}, \mathcal{F}'$ and with L verifying $t \in L$.*

The theorem that establishes the correctness of the approach is an easy consequence of the Lemma.

Theorem 1. *Let P be the \mathcal{TOY} program of Figure 4, Q an SXQ query, t a \mathcal{TOY} pattern, and θ a substitution such that $\mathcal{P} \vdash (\text{sxq } Q\theta == t)$ for some θ . Then $\llbracket Q \rrbracket_0(\emptyset, []) :=^* (\mathcal{F}, L)$, for some forest \mathcal{F} , and L verifying $t \in L$.*

Proof. In Lemma 1 consider the position $p \equiv \varepsilon$. Then $Q' \equiv Q$, $V = \emptyset$ and $k = 0$. Without loss of generality we can restrict in the conclusion to $\mathcal{F} = \emptyset$, because $\theta(V) = \emptyset$ and therefore \mathcal{F} is not used during the rewriting process. Then the conclusion of the theorem is the conclusion of the lemma.

Thus, our approach is correct. The next Lemma allows us to prove that it is also complete, in the sense that the \mathcal{TOY} program can produce every answer obtained by the XQ operational semantics.

Lemma 2. *Let \mathcal{P} be the \mathcal{TOY} program of Figure 4. Let Q' be a SXQ query and Q, p such that $Q \equiv Q'_{|p}$. Define $V = \text{Rel}(Q', p)$ (see Definition 2) and $k = |V|$. Suppose that $\llbracket Q \rrbracket_k(\mathcal{F}, \bar{e}_k) :=^* (\mathcal{F}', \bar{a}_n)$ for some $\mathcal{F}, \mathcal{F}', \bar{e}_k, \bar{a}_n$. Then, for every a_j , $1 \leq j \leq n$, there is a substitution θ such that $\theta(X_i) = e_i$ for $X_i \in V$ and a CRWL-proof proving $\mathcal{P} \vdash \text{sxq } Q\theta == a_j$.*

As in the case of correctness, the completeness theorem is just a particular case of the Lemma:

Theorem 2. *Let \mathcal{P} be the \mathcal{TOY} program of Figure 4. Let Q be a SXQ query and suppose that $\llbracket Q \rrbracket_k(\emptyset, []) :=^* (\mathcal{F}, \bar{a}_n)$ for some \mathcal{F}, \bar{a}_n . Then for every a_j , $1 \leq j \leq n$, there is $\mathcal{P} \vdash (\text{sxq } Q)\theta == a_j$ for some substitution θ .*

Proof. As in Theorem 1, suppose $p \equiv \varepsilon$ and thus $Q' \equiv Q$. Then $V = \emptyset$ and $k = 0$. Then, if $\llbracket Q \rrbracket_0(\emptyset, \emptyset) :=^* (\mathcal{F}, \bar{a}_n)$ it is easy to check that $\llbracket Q \rrbracket_0(\mathcal{F}', \emptyset) :=^* (\mathcal{F}, \bar{a}_n)$ for any \mathcal{F}' . Then the conclusion of the lemma is the same as the conclusion of the Theorem.

The proofs of Lemmata 1 and 2 can be found in [2].

5 Application: Test Case Generation

In this section we show how an embedding of SXQ into \mathcal{TOY} can be used for obtaining test-cases for the queries. For instance, consider the erroneous query of the next example.

Example 5. Suppose that the user also wants to include the publisher of the book among the data obtained in Example 1. The following query tries to obtain this information:

```
Q = for $b in doc("bib.xml")/bib/book,
     $r in doc("reviews.xml")/reviews/entry,
     where $b/title = $r/title
     for $booktitle in $r/title,
         $revtext in $r/review,
         $publisher in $r/publisher
     return <rev> $booktitle $publisher $revtext </rev>
```

However, there is an error in this query, because in `$r/publisher` the variable `$r` should be `$b`, since the publisher is in the document `"bib.xml"`, not in `"reviews.xml"`. The user does not notice that there is an error, tries the query (in \mathcal{TOY} or in any XQuery interpreter) and receives an empty answer.

In order to check whether a query is erroneous, or even to help finding the error, it is sometimes useful to have test-cases, i.e., XML files which can produce some answer for the query. Then the test-cases and the original XML documents can be compared, and this can help finding the error. In our setting, such test-cases are obtained for free, thanks to the generate and test capabilities of logic programming. The general process can be described as follows:

1. Let Q' be the translation `parse_xquery` Q of query Q into \mathcal{TOY} .
2. Let F_1, \dots, F_k be the names of the XML documents occurring in Q' . That is, for each F_i , $1 \leq i \leq k$, there is an occurrence of an expression of the form `load_xml_file(F_i)` in Q' (which corresponds to expressions `doc(F_i)` in Q). Let Q'' be the result of replacing each `doc(F_i)` expression by a new variable D_i , for $i = 1 \dots k$.
3. Let `"expected.xml"` be a document containing an expected answer for the query Q .
4. Let E the expression `Q''==load_doc "expected.xml"`.
5. Try the goal
 $G \equiv E, \text{write_xml_file } D_1 \ F'_1, \dots, \text{write_xml_file } D_k \ F'_k$

The idea is that the goal G looks for values of the logic variables D_i fulfilling the strict equality. The result is that after solving this goal, the D_i variables contain XML documents that can produce the expected answer for this query. Then each document is saved into a new file with name F'_i . For instance F'_i can consist of the original name F_i preceded by some suitable prefix *tc*. The process can be automatized, and the result is the code of Figure 5.

```

prepareTC (xp E)          = (xp E',L)
                           where (E',L) = prepareTCPath E
prepareTC (xmlExp X)     = (xmlExp X, [])
prepareTC (comp Q1 Q2)   = (comp Q1' Q2', L1++L2)
                           where (Q1',L1) = prepareTC Q1
                                 (Q2',L2) = prepareTC Q2
prepareTC (xfor X Q1 Q2) = (xfor X Q1' Q2', L1++L2)
                           where (Q1',L1) = prepareTC Q1
                                 (Q2',L2) = prepareTC Q2
prepareTC (xif (Q1:=Q2) Q3) = (xif (Q1':=Q2') Q3',L1++(L2++L3))
                           where (Q1',L1) = prepareTC Q1
                                 (Q2',L2) = prepareTC Q2
                                 (Q3',L3) = prepareTC Q3
prepareTC (xif (cond Q1) Q2) = (xif (cond Q1) Q2, L1++L2)
                           where (Q1',L1) = prepareTC Q1
                                 (Q2',L2) = prepareTC Q2

prepareTCPath (var X)    = (var X, [])
prepareTCPath (X :/ S)  = (X :/ S, [])
prepareTCPath (doc F S) = (A :/ S, [write_xml_file A ("tc"++F)])

generateTC Q F = true <== sxq Qtc == load_doc F, L==_
                  where (Qtc,L) = prepareTC Q

```

Fig. 5. \mathcal{TCY} transformation rules for SXQ

The code uses the list concatenation operator ++ which is defined in \mathcal{TCY} as usual in functional languages such as Haskell. It is worth observing that if there are no test-case documents that can produce the expected result for the query, the call to `generateTC` will loop. The next example shows the generation of test-cases for the wrong query of Example 5.

Example 6. Consider the query of Example 5, and suppose the user writes the following document “`expected.xml`”:

```

<rev>
  <title>Some title</title>
  <review>The review</review>
  <publisher>Publisher</publisher>
</rev>

```

This is a possible expected answer for the query. Now we can try the goal:

```
Toy> Q == parse_xquery "for....", R == generateTC Q "expected.xml"
```

The first strict equality parses the query, and the second one generates the XML documents which constitute the test cases. In this example the test-cases obtained are:

```

% bibtc.xml
<bib>
  <book>
    <title>Some title</title>
  </book>
</bib>

% revtc.xml
<reviews>
  <entry>
    <title>Some title</title>
    <review>The review </review>
    <publisher>Publisher</publisher>
  </entry>
</reviews>

```

By comparing the test-case ‘`revtc.xml`’ with the file ‘`reviews.xml`’ we observe that the publisher is not in the reviews. Then it is easy to check that in the query the publisher is obtained from the reviews instead of from the *bib* document, and that this constitutes the error.

6 Conclusions

The paper shows the embedding of a fragment of the XQuery language for querying XML documents into the functional-logic language \mathcal{FOY} . Although only a small subset of XQuery consisting of *for*, *where/if* and *return* statements has been considered, the users of \mathcal{FOY} can now perform simple queries such as *join* operations. The formal definition of the embedding allows us to prove the soundness of the approach with respect to the operational semantics of XQuery. The proposal respects the declarative nature of \mathcal{FOY} , exploiting its non-deterministic nature for obtaining the different results produced by XQuery expressions. An advantage of this approach with respect to the use of lists usually employed in functional languages is that our embedding allows the user to generate test-cases automatically when possible, which is useful for testing the query, or even for helping to find the error in the query. An extended version of this paper, including the proofs of the theoretical results and more detailed explanations about how to install \mathcal{FOY} and run the prototype can be found in [2].

The most obvious future work would be introducing the *let* statement, which presents two novelties. The first is that they are *lazy*, that is, they are not evaluated if they are not required by the result. This part is easy to fulfill since we are in a lazy language. In particular, they could be introduced as local definitions (*where* statements in \mathcal{FOY}). The second novelty is more difficult to capture, and it is that the variables introduced by *let* represent an XML sequence. The natural representation in \mathcal{FOY} would be a list, but the non-deterministic nature of our proposal does not allow us to collect all the results provided by an expression in a declarative way. A possible idea would be to use the functional-logic language Curry [8] and its encapsulated-search [10], or even the non-declarative *collect* primitive included in \mathcal{FOY} . In any case, this will imply a different theoretical framework and new proofs for the results. A different line for future work is the use of test cases for finding the error in the query using some variation of declarative debugging [14] that could be applied to this setting.

References

1. Almendros-Jiménez, J.M.: An Encoding of XQuery in Prolog. In: Bellahsene, Z., Hunt, E., Rys, M., Unland, R. (eds.) XSym 2009. LNCS, vol. 5679, pp. 145–155. Springer, Heidelberg (2009)
2. Almendros-Jiménez, J., Caballero, R., García-Ruiz, Y., Sáenz-Pérez, F.: A Declarative Embedding of XQuery in a Functional-Logic Language. Technical Report SIC-04/11, Facultad de Informática, Universidad Complutense de Madrid (2011), <http://gpd.sip.ucm.es/rafa/xquery/>
3. Baader, F., Nipkow, T.: Term Rewriting and All That. Cambridge University Press (1999)
4. Benedikt, M., Koch, C.: From XQuery to relational logics. ACM Trans. Database Syst. 34, 25:1–25:48 (2009)
5. Caballero, R., García-Ruiz, Y., Sáenz-Pérez, F.: Integrating XPath with the Functional-Logic Language Toy. In: Rocha, R., Launchbury, J. (eds.) PADL 2011. LNCS, vol. 6539, pp. 145–159. Springer, Heidelberg (2011)
6. Fegaras, L.: Propagating updates through XML views using lineage tracing. In: IEEE 26th International Conference on Data Engineering (ICDE), pp. 309–320 (March 2010)
7. González-Moreno, J., Hortalá-González, M., López-Fraguas, F., Rodríguez-Artalejo, M.: A Rewriting Logic for Declarative Programming. In: Riis Nielson, H. (ed.) ESOP 1996. LNCS, vol. 1058, pp. 156–172. Springer, Heidelberg (1996)
8. Hanus, M.: Curry: An Integrated Functional Logic Language (version 0.8.2 March 28, 2006) (2003), <http://www.informatik.uni-kiel.de/~mh/curry/>
9. Hanus, M.: Declarative processing of semistructured web data. Technical report 1103, Christian-Albrechts-Universität Kiel (2011)
10. Hanus, M., Steiner, F.: Controlling Search in Declarative Programs. In: Palamidessi, C., Meinke, K., Glaser, H. (eds.) ALP 1998 and PLILP 1998. LNCS, vol. 1490, pp. 374–390. Springer, Heidelberg (1998)
11. Fraguas, F.J.L., Hernández, J.S.: \mathcal{TOY} : A Multiparadigm Declarative System. In: Narendran, P., Rusinowitch, M. (eds.) RTA 1999. LNCS, vol. 1631, pp. 244–247. Springer, Heidelberg (1999)
12. Olteanu, D., Meuss, H., Furche, T., Bry, F.: XPath: Looking Forward. In: Chaudhri, A.B., Unland, R., Djeraba, C., Lindner, W. (eds.) EDBT 2002. LNCS, vol. 2490, pp. 109–127. Springer, Heidelberg (2002)
13. Seipel, D., Baumeister, J., Hopfner, M.: Declaratively Querying and Visualizing Knowledge Bases in XML. In: Seipel, D., Hanus, M., Geske, U., Bartenstein, O. (eds.) INAP/WLP 2004. LNCS (LNAI), vol. 3392, pp. 16–31. Springer, Heidelberg (2005)
14. Shapiro, E.: Algorithmic Program Debugging. ACM Distinguished Dissertation. MIT Press (1982)
15. W3C. Extensible Markup Language (XML) (2007)
16. W3C. XML Path Language (XPath) 2.0 (2007)
17. W3C. XML Query Use Cases (2007), <http://www.w3.org/TR/xquery-use-cases/>
18. W3C. XQuery 1.0: An XML Query Language (2007)
19. W3C. XQuery 1.0 and XPath 2.0 Formal Semantics, 2nd edn. (2010), <http://www.w3.org/TR/xquery-semantics/>
20. Walmsley, P.: XQuery. O'Reilly Media, Inc. (2007)