



# A Uniform Programming Language for Implementing XML Standards

Pavel Labath, Joachim Niehren

► **To cite this version:**

Pavel Labath, Joachim Niehren. A Uniform Programming Language for Implementing XML Standards. 41st SOFSEM: International Conference on Current Trends in Theory and Practice of Computer Science, Jan 2015, Pec pod Sněžkou, Czech Republic. hal-00954692

**HAL Id: hal-00954692**

**<https://hal.inria.fr/hal-00954692>**

Submitted on 6 Oct 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.

# A Uniform Programming Language for Implementing XML Standards<sup>\*</sup>

Pavel Labath<sup>1</sup> and Joachim Niehren<sup>2</sup>

<sup>1</sup> Comenius University, Bratislava      <sup>2</sup> INRIA, Lille

**Abstract.** We propose X-Fun, a core language for implementing various XML standards in a uniform manner. X-Fun is a higher-order functional programming language for transforming data trees based on node selection queries. It can support the XML data model and XPATH queries as a special case. We present a lean operational semantics of X-Fun based on a typed lambda calculus that enables its in-memory implementation on top of any chosen path query evaluator. We also discuss compilers from XSLT, XQUERY and XPROC into X-Fun which cover the many details of these standardized languages. As a result, we obtain in-memory implementations of all these XML standards with large coverage and high efficiency in a uniform manner from SAXON's XPATH implementation.

**Keywords:** XML transformations, database queries, functional programming languages, compilers.

## 1 Introduction

A major drawback of query-based functional languages with data trees so far is that they either have low coverage in theory and practice or no lean operational semantics. Theory driven languages are often based on some kind of macro tree transducers [12,5,3], which have low coverage, in that they are not closed under function composition [4] and thus not Turing complete (for instance type checking is decidable [11]). The W3C standardised languages XQUERY [13] and XSLT [7], in contrast, have large coverage in practice (string operations, data joins, arithmetics, aggregation, etc.) and in theory, since they are closed by function composition and indeed Turing complete [8]. The definitions of these standards, however, consist of hundreds of pages of informal descriptions. They neither explain how to build a compiler in a principled manner nor can they be used as a basis for formal analysis.

A second drawback is the tower of languages approach, adopted for standardised XML processing languages. What happened in the case of XML was the development of a separate language for each class of use cases, which all host the XPATH language for querying data trees based on node navigation. XSLT serves for use cases with recursive document transformations such as HTML publishing, while XQUERY was developed for use cases in which XML databases are

---

<sup>\*</sup> This research was supported in part by the grant VEGA 1/0979/12

queried. Since the combination of both is needed in most larger applications, the XML pipeline language XPROC [17,16,18] was developed and standardised again by the W3C. This resulted in yet another functional programming language for processing data trees based on XPATH.

For resolving the above two drawbacks, the question is whether there exists a uniform core language for processing data trees that can cover the different XML standards in a principled manner. It should have a lean and formal operational semantics, support node selection queries as with XPATH and it should be sufficiently expressive in order to serve as a core language for implementing XQUERY, XSLT, and XPROC in a uniform manner.

*Related work.* An indicator for the existence of a uniform core language for XML processing is that the omnipresent Saxon system [14] implements XSLT and XQUERY on a common platform. However, there is no formal description of this platform as a programming language, and it does not support the XML pipeline language XPROC so far. Instead, the existing implementations of XPROC, CALABASH [16] and QUIXPROC [18], are based on Saxon’s XPATH engine directly.

The recent work from Castagna et. al. [2] gives further hope that our question will find a positive answer. They present an XPATH-based functional programming language with a lean formal model based on the lambda calculus, which thus satisfies our first two conditions above and can serve as a core language for implementing a subset of XQUERY 3.0. We believe that relevant parts of XSLT and XPROC can also be compiled into this language, even though this is not shown there. The coverage, however, will remain limited, in particular on the XPATH core (priority is given to strengthening type systems). Therefore, our last requirement is not satisfied.

*Contributions.* In this paper, we present the first positive answer to the above question based on X-Fun. This is a new purely functional programming language. X-Fun is a higher-order language and it supports the evaluation of path-based queries that select nodes in data trees. The path queries are mapped to X-Fun expressions, whose values can be computed dynamically. In contrast to most previous interfaces between databases and programming languages, we overload variables of path queries with variables of X-Fun. In this manner, the variables in path queries are always bound to tree nodes, before the path query is evaluated itself. We note in particular, that path queries are not simply mapped to X-Fun expressions of type string.

The formal model of the operational semantics of X-Fun is a lambda calculus with a parallel call-by-value reduction strategy. Parallel evaluation is possible due to the absence of imperative data structures. The main novelty in X-Fun admission of tree nodes as values of type **node**. Which precise nodes are admitted depends on a tree store. New nodes can be created dynamically by adding new trees to the tree store. The same tree can be added twice to the store but with different nodes. How nodes are represented internally can be freely chosen by the X-Fun implementation and is hidden from the programmer.

X-Fun can serve as a uniform core language for implementing XQUERY, XSLT and XPROC. In order to do so, we have developed compilers of all three languages into X-Fun. We also discuss how to implement X-Fun in an in-memory fashion on top of any in-memory XPATH evaluator. Based on our compilers, we thus obtain new in-memory implementations of XQUERY, XSLT and XPROC with large coverage. Our implementation has very good efficiency and outperforms the most widely used XPROC implementation by a wide margin.

*Outline.* In Section 2 we introduce our general model of data trees, alongside its application to XML documents. The syntax and type system of the X-Fun language is introduced in Section 3. The applications of X-Fun to XML document transformation is studied in Section 4, where we discuss compilers from other XML processing languages into X-Fun. Section 5 contains our notes on the implementation of X-Fun and the results of our experiments.

## 2 Preliminaries

We introduce a general concept of data trees which will be used in the X-Fun language. We also show how to instantiate the trees to the XML data model.

### 2.1 Data values and data trees

We fix a finite set  $Char$  whose elements will be called characters. A data value " $c_1 \dots c_n$ " is a word of characters for  $c_1, \dots, c_n \in Char$ . We define  $String = Char^*$  to be the set of all data values, and  $nil = ""$  to be the empty data value.

Next, we will fix a natural number  $k \geq 1$  and introduce data trees in which each node contains exactly  $k$  data values with characters in  $Char$ .

A *node label* is a  $k$ -tuple of data values, i.e., an element of  $(String)^k$ . The set of data trees  $\mathcal{T}$  of label size  $k$  over  $Char$  is the least set that contains all pairs of node labels and sequences of data trees in  $\mathcal{T}$ . That is, it contains all unranked trees  $t$  with the abstract syntax  $t ::= l(t_1, \dots, t_n)$ , where  $n \geq 0$  and  $l \in String^k$ . It should be noticed that the set of node labels is infinite, but that each node label can be represented finitely.

### 2.2 XML data model

For XML data trees, we can fix  $k = 4$  and  $Char$  the set of Unicode characters, and restrict ourselves to node labels of the following forms, where all  $v_i$  are data values:

("element", $v_1, v_2, nil$ )	("attribute", $v_1, v_2, v_3$ )
("comment", $nil, nil, v_3$ )	("processing-instruction", $v_1, nil, v_3$ )
("document", $nil, nil, nil$ )	("text", $v_1, nil, nil$ )

An element ("element",  $v_1, v_2, nil$ ) has three non-nil data values: its type "element", a name  $v_1$  and a namespace  $v_2$ . An attribute has four data values: its

type, a name  $v_1$ , a namespace  $v_2$ , and the attribute value  $v_3$ . A text node contains its type and its text value  $v_3$ . Besides these, there are comments, processing instructions and the rooting document node.

### 3 Language X-Fun

In this section, we introduce X-Fun, a new functional programming language for transforming data trees. X-Fun can be applied to all kinds of data trees with a suitable choice of its parameters. We will instantiate the case of data trees satisfying the XML data model concomitant with XPATH as a query language.

We start with introducing the types and values of X-Fun (Section 3.1). Then we explain how to map path queries to X-Fun values, by using particular X-Fun expressions with variables (Section 3.2). The general syntax of X-Fun expressions is given in Section 3.3. Some syntactic sugar and an example of an X-Fun program are given in Sections 3.5 and 3.6. Discussion of the typing rules for X-Fun's type system and the formal semantics of X-Fun can be found in the research report [10].

#### 3.1 Types and Values

The X-Fun language supports higher-order values and expressions with the following types:

$$T ::= \mathbf{none} \mid \mathbf{node} \mid \mathbf{tree} \mid \mathbf{number} \mid \mathbf{bool} \mid \mathbf{char} \\ \mid T_1 \times \dots \times T_n \mid [T] \mid T_1 \rightarrow T_2 \mid T_1 \cup T_2$$

A value of type **char** is an element of *Char*, a value of type **tree** is an element of  $\mathcal{T}$ . A value of type **number** is a floating point number, while the values of type **bool** are the Boolean values *true* and *false*. A value of type **node** will be a node of the graph of one of the trees stored by the environment of the X-Fun evaluator. The precise node identifiers chosen by the evaluator are left internal (to the mapping from trees to graphs).

As usual, we support list types  $[T]$  which denote all lists of values of type  $T$ , product types  $T_1 \times \dots \times T_n$  whose values are all tuples of the values of types  $T_i$ , and function types  $T_1 \rightarrow T_2$  whose values are all partial functions of values of type  $T_1$  to values of type  $T_2$ . Besides these, we also support type unions in the obvious manner.

A data value " $c_1 \dots c_n$ "  $\in$  *String* is considered as a list of characters of type **string** =  $[\mathbf{char}]$ . A node label is considered a k-tuple of strings, i.e., as a value of type **label** =  $\mathbf{string}^k$ . Hedges are considered as lists of trees of type **hedge** =  $[\mathbf{tree}]$ .

Since XPATH can return sequences of items of different types, we define the type **pathresult** as  $\mathbf{node} \cup \mathbf{number} \cup \mathbf{string} \cup \mathbf{bool}$ . The result of evaluating a path expression will then be of type  $[\mathbf{pathresult}]$ . To be able to specify path expressions, we define the type **path** as  $[\mathbf{char} \cup \mathbf{pathresult} \cup [\mathbf{pathresult}]]$ , i.e., as list of characters, individual items returned by a path expression, and whole sequences of those items.

### 3.2 XPath queries as X-Fun expressions

We will consider XPath expressions as values of our programming language. This is done in such a manner that the variables in XPATH expressions can be bound to values of the programming language. For instance, if we have an XPATH expression

```
$x//book[auth=$y]
```

then one might want to evaluate this expression while variable  $x$  is bound to a node of some tree and variable  $y$  to some data value. In X-Fun, the above query will be represented by the following expression of type **path**, where  $x$  is a variable of type **node** and  $y$  a variable of type **string**:

```
 $x :: ' / ' :: ' / ' :: ' b ' :: ' o ' :: ' o ' :: ' k ' :: ' [ ' :: ' a ' :: ' u ' :: ' t ' :: ' h ' :: ' = ' :: ' y ' :: ' ] ' :: nil$ 
```

The concrete syntax of X-Fun supports syntactic sugar for values of type **path**, so that the above expression can be defined as:

```
"$x//book[auth=$y]"
```

In order to enable the evaluation of path expressions, X-Fun supports a builtin function `evalPath` of type **path**  $\rightarrow$  [**pathresult**]. In an implementation of X-Fun, this function can be mapped straightforwardly to existing XPATH evaluators.

### 3.3 Syntax of X-Fun expressions

X-Fun is a purely functional programming language whose values subsume higher-order function, trees, strings, numbers and Boolean values. The evaluation strategy of X-Fun is fully parallel, which is possible since no imperative constructs are permitted.

The syntax of X-Fun programs  $E$  is given in Figure 1. All expressions of X-Fun are standard in functional programming languages, so we only briefly describe different kinds of subexpressions of X-Fun programs.

A variable  $x$  is evaluated to the value of the corresponding type. The constant expression  $c$  returns the respective constant, which can be a Boolean value, a number or a character from *Char*. The list constructor  $E_1 :: E_2$  prepends an element to a list, while the tuple constructor  $(E_1, \dots, E_n)$  constructs tuples.

The match expression **match**  $E \{ P_1 \rightarrow E_1, \dots, P_n \rightarrow E_n \}$  selects one of the branches  $E_i$  based on the patterns  $P_i$ , which are matched against the value of  $E$ . The pattern  $x : T$  captures a matched value of type  $T$  into a variable. The pattern  $!(E)$  matches the value against the value of expression  $E$ . Here, the matching of functional values, or lists/tuples that contain functions is not permitted. Pattern  $P_1 :: P_2$  matches a list if  $P_1$  and  $P_2$  match its head and tail, while the pattern  $(P_1, \dots, P_n)$  matches tuples.

A function expression **fun**  $x : T_1 \rightarrow T_2 \{ E \}$  returns a new function, with the argument  $x : T_1$  and the return value of type  $T_2$  obtained by the evaluation of the function body  $E$ . The expression  $E_1(E_2)$  applies a function to a value. X-Fun also supports exception handling, where exceptions are values of type **string**.

Expressions	Patterns
$E ::= x$	$P ::= x : T$
$c$	$!(E)$
$E_1 :: E_2$	$P_1 :: P_2$
$(E_1, \dots, E_n), n \geq 2$	$(P_1, \dots, P_n), n \geq 2$
<b>match</b> $E \{ P_1 \rightarrow E_1, \dots, P_n \rightarrow E_n \}$	
<b>fun</b> $x:T_1 \rightarrow T_2 \{ E \}$	
$E_1(E_2)$	
<b>try</b> $E_1$ <b>catch</b> ( $x$ ) $E_2$	
<b>raise</b> ( $E$ )	

**Fig. 1.** Syntax of X-Fun's expressions

### 3.4 Builtin operators

At the beginning of the evaluation, the environment contains bindings of the global variables given in Figure 2.

Parameters	Fixed
Global variable TYPE	Global variable TYPE
<b>makeTree</b> <b>label</b> $\times$ <b>[tree]</b> $\rightarrow$ <b>tree</b>	<b>nil</b> <b>[none]</b>
<b>evalPath</b> <b>path</b> $\rightarrow$ <b>[pathresult]</b>	<b>subtree</b> <b>node</b> $\rightarrow$ <b>tree</b>
<b>less</b> <b>char</b> $\times$ <b>char</b> $\rightarrow$ <b>bool</b>	<b>label</b> <b>node</b> $\rightarrow$ <b>label</b>
	<b>addTree</b> <b>tree</b> $\rightarrow$ <b>node</b>

**Fig. 2.** Builtin operators of X-Fun

The first block contains three functions, whose semantics are parameters of the language, and depend on the query language and data model. For a label  $l$  and a sequence of trees  $h$ , the function application **makeTree**( $l, h$ ) returns the data tree  $l(h)$ , if  $l(h)$  is a well-formed data tree (e.g., in the XML data model attributes cannot have children) and raises an exception otherwise. The function **evalPath**( $p$ ) evaluates a path expression  $p$ . Whenever  $p$  is not well-formed (e.g., with respect to the XPATH 3.0 specification) an error is raised. Note that path expressions are X-Fun values, which means they can be computed dynamically by the X-Fun program using information from the input data tree. We will also define functions **evalPath** $_T$ , on top of **evalPath**, for  $T = [\mathbf{node}], [\mathbf{string}]$ , etc. These functions verify (using a **match** expression with a typecase) that the result of the path call is of type  $T$  and raise an exception otherwise.

The next four operators are generic and do not depend on the specific kind of data trees. The variable **nil** refers to the empty list. A function application **subtree**( $v$ ) returns the subtree rooted at node  $v$ , while a function application

`label(v)` returns the label of the node. The function `addTree` returns the identifier of the root node of the tree, and is used for storing the graph of the tree in the environment. This function can be used to access nodes of newly generated trees by starting path navigation from their root.

### 3.5 Syntactic sugar

In the X-Fun snippets in the rest of the paper we shall employ some syntactic shortcuts, which enable us to express more succinctly some X-Fun constructs:

**list concatenation** We shall use the binary operator `*` to concatenate two lists.

**simplified patterns** When the type of a capture variable can be deduced from the matched expression we shall omit the “:  $T$ ” in the capture pattern. This happens when the **match** expression is used to decompose lists and tuples instead of doing a typecase. For example, we shall simply write **match**  $E \{ h :: t \rightarrow E_1, e \rightarrow E_2 \}$  to get the head and tail of a list.

**let-declarations** We shall use the syntax **let**  $x_1 = E_1, \dots, x_n = E_n$  **in**  $E$  instead of **match**  $(E_1, \dots, E_n) \{ (x_1, \dots, x_n) \rightarrow E \}$  as a more familiar way to declare variables.

**tuple arguments** We shall allow tuple arguments to functions to be written without an extra pair of parentheses. I.e.,  $f(a, b)$  instead of  $f((a, b))$ . This is unambiguous since tuples always have at least two members.

### 3.6 Example

In Figure 3 we illustrate a transformation that converts an address book into HTML. The address fields are assumed to be unordered in the input data tree, while the fields of the output HTML addresses should be published in the order **name**, **street**, **city** and, **phone**.

<code>&lt;addresses&gt;</code>		<code>&lt;ol&gt;</code>
<code>&lt;address&gt;</code>		<code>&lt;li&gt;</code>
<code>&lt;name&gt;Jemal Antidze&lt;/name&gt;</code>		<code>&lt;p&gt;Jemal Antidze&lt;/p&gt;</code>
<code>&lt;phone&gt;99532 305972&lt;/phone&gt;</code>		<code>&lt;p&gt;Tblissi&lt;/p&gt;</code>
<code>&lt;city&gt;Tblissi&lt;/city&gt;</code>		<code>&lt;p&gt;Phone: 99532 305972&lt;/p&gt;</code>
<code>&lt;phone&gt;99532 231231&lt;/phone&gt;</code>		<code>&lt;p&gt;Phone: 99532 231231&lt;/p&gt;</code>
<code>&lt;/address&gt;</code>	<code>⇒</code>	<code>&lt;/li&gt;</code>
<code>&lt;address&gt;</code>		<code>&lt;li&gt;</code>
<code>&lt;name&gt;Joachim Niehren&lt;/name&gt;</code>		<code>&lt;p&gt;Joachim Niehren&lt;/p&gt;</code>
<code>&lt;city&gt;Lille&lt;/city&gt;</code>		<code>&lt;p&gt;Rue Esquermoise&lt;/p&gt;</code>
<code>&lt;street&gt;Rue Esquermoise&lt;/street&gt;</code>		<code>&lt;p&gt;Lille&lt;/p&gt;</code>
<code>&lt;/address&gt;</code>		<code>&lt;/li&gt;</code>
<code>&lt;/addresses&gt;</code>		<code>&lt;/ol&gt;</code>

**Fig. 3.** Publication of an address book in HTML except for secret entries



An X-Fun program defining this transformation is given in Figure 4. Starting at the root it first locates all address records, and applies the function `convert_address` to each of them. For each address record, the program first extracts the values of the fields `name`, `street`, and `city` located at some children of `x`. These values are then bound to variables named `alike` and later output as text nodes. The example program uses the standard map function, which can be defined in X-Fun for every  $T$  and  $T'$  as follows

```
map $T \rightarrow T'$  = fun x: (T → T') × [T] → [T'] { match x {
  (f, head :: tail) → f(head) :: map $T \rightarrow T'$ (f, tail)
  other → nil
} }
```

and the functions `element` and `text`, which are wrappers around `makeTree` which facilitate creation of nodes of the correct kind.

```
fun book : tree → tree {
  let bookroot = addTree(book) in
  let convert_address = fun x : node → tree {
    let name = evalPath[node]("$x/child::name/text()"),
        street = evalPath[node]("$x/child::street/text()"),
        city = evalPath[node]("$x/child::city/text()") in
    element("li",
      element("p", mapnode → tree(name, subtree)) ::
      element("p", mapnode → tree(street, subtree)) ::
      element("p", mapnode → tree(city, subtree)) ::
      mapstring → tree(
        fun x: string → tree {
          element("p", text("Phone: " * x)) :: nil
        }, evalPath[string]("data($x/child::phone)"))
    )
  } in
  element("ol", mapnode → tree(convert_address,
    evalPath[node]("$bookroot/descendant::address")))
}
```

**Fig. 4.** X-Fun program converting address books to HTML

## 4 Translations from other XML languages

In this section, we briefly sketch translations from the standard XML processing languages, XSLT XQUERY and XPROC. A more thorough treatment of this topic can be found in the full version of the paper.[10] By implementing these three compilers, we obtain a uniform implementation of the whole XML processing stack based on a single X-Fun evaluator.

XSLT. Each template in the XSLT stylesheet is translated to a function in X-Fun. Furthermore, for each mode, we produce an additional function which implements the selection of the correct template from the set of templates associated with that mode according to their match patterns. The `call-template` and `apply-templates` instructions are translated as calls to the template or mode functions respectively. In the `copy-of` instruction, the nodes returned by the XPATH expression are copied to the output using the `subtree` function and strings and numbers are converted to a new text node with a call to `makeTree`. The instructions constructing elements, attributes and other XML nodes translate to corresponding calls to `makeTree`. The `for-each` instruction translates to a call to `map`, where the list to map over is produced by a call to `evalPath` and the mapping function is the body of the `for-each` instruction. Other XSLT instructions like `if` and `choose` can be translated similarly.

XQUERY. The feature that most distinguishes XQUERY is the SQL-like FLWOR expression. It enables the programmer to create a stream of tuples using the `for` and `let` clauses, filter them with a `where` clause and then reorder them using the `order by` clause. There is no single expression in X-Fun which covers this functionality, but it is easy to build it piecewise. Using several `evalPath` calls we can construct the list of tuples which corresponds to the tuple stream of XQUERY. Sorting and filtering of a list are functions easily definable in a functional language, and the functionality of `where` and `order by` is translated to calls to these functions. The sort and filter conditions are given again by calls to `evalPath` with the appropriate XPATH expression. Translation of other XQUERY constructs like the `if` expressions and functions proceeds in a straight forward manner.

XPROC. By encapsulating each processing step in a function, X-Fun can easily express the multi-stage processing which is inherent in XPROC. The pipelines then become simple function compositions. XPROC steps which invoke XQUERY or XSLT processing are handled by defining a function whose body is the translation of the respective program. Simple XPROC steps like `split-sequence`, which splits a sequence of documents into two based on an XPATH criterion are defined as normal X-Fun functions and provided as a library. The pipeline then simply calls these functions to do the required processing. The rest of the constructs like choosing among alternative subpipelines (`choose`) or looping over documents in a sequence are compiled to `match` and `map` expressions in X-Fun.

## 5 Implementation and Experiments

We have implemented a proof-of-concept X-Fun language evaluator in the Java programming language. We have instantiated X-Fun with the XML data model, using standard Java libraries for manipulating XML trees. We have used XPATH as the path language, as implemented by SAXON. We have used standard techniques for implementing functional languages, using the heap to store the values

and the environment of the program and a stack for representing recursive function calls. We reduce an expression in all possible positions in an arbitrary order.

We have attempted to interface our implementation with TATOO, a highly efficient evaluator of an XPATH fragment based on [1]. Unfortunately, the penalty of crossing the language barrier (TATOO is implemented in OCAML) shadowed all performance gains from a faster implementation, so we could not perform any significant experiments. To see the difference in performance in using a faster XPATH implementation, we would need to implement X-Fun in OCAML as well.

We have also implemented the compilers of XSLT and XQUERY into X-Fun. In order to support real-world XSLT and XQUERY, they need support for additional features, like modules and various optional attributes of expressions in these languages (e.g., grouping with the `group-starting-with` attribute, etc.). However, none of these limitations are fundamental and they are not implemented because of their volume. The supported fragment is wide enough to run all queries from the XMARK [15] benchmark.

We don't have an XPROC compiler implementation, but for the purposes of testing we have run X-Fun on manually translated programs.

## 5.1 Experiments

To evaluate the performance of our implementation, we have compared it with the leading industry tool, the SAXON XSLT and XQUERY processor. To compare our performance on XPROC pipelines, we have used CALABASH, the most frequently used XPROC processor, as baseline. The tests were run on a computer with an Intel Core i7 processor running at 2.8 GHz, with 4GB of RAM and a SATA hard drive, running 64-bit Linux operating system.

First, we have compared the running time of our implementation on XQUERY programs. We used the queries from the XMARK benchmark, and the results are in Figure 5. The tests show that the running time of both tools is comparable. X-Fun is faster in case of simple queries (Q6, Q7, Q15, which contain just a simple loop), while SAXON is faster on queries involving joins (e.g., Q8, Q9, Q11). On the rest of queries our implementation of X-Fun is at most 20% slower than the competition, which we consider a good result as Saxon is a highly optimised industry tool, while we have not spent much time optimising the performance of our X-Fun implementation.

For the XSLT test, we used a transformation publishing an address book to HTML. The transformation in question is a more elaborate version of the program in Figure 4, and it includes about 40 XPATH expressions. The tests show that SAXON is about 4 times faster than our tool (for example, 15.7 vs. 63 seconds on a 200 MB document) and that the time of both tools scales linearly with the document size.

In the XPROC comparison, we have a simple pipeline consisting of 4 steps. First, it selects subtrees from the input document, splits the resulting sequence into two based on the presence of some node. The documents from the two sequences are then joined into pairs and these pairs are concatenated to form a single document again. We have compared the performance of CALABASH with

Query	X-Fun	SAXON	Query	X-Fun	SAXON	Query	X-Fun	SAXON
Q1	13.5	10.9	Q8*	962	592	Q15	12.0	14.4
Q2	13.6	12.9	Q9*	1235	705	Q16	13.6	11.8
Q3	14.0	12.5	Q10	314	222	Q17	13.9	12.4
Q4	16.7	12.8	Q11*	650	410	Q18	14.4	12.5
Q5	17.2	13.8	Q12	595	317	Q19	20.8	15.4
Q6	11.5	13.6	Q13	20.5	11.6	Q20	13.8	12.0
Q7	11.4	12.5	Q14	14.5	12.8			

**Fig. 5.** Running time in seconds of X-Fun and SAXON on queries from the XMARK benchmark on a 500 MB document. The three queries marked with “\*”, due to their complexity, were run on a 300 MB document.

our implementation of the pipeline in X-Fun. Both implementations show linear scalability with respect to size of the input and the pipeline, as can be seen in Figures 6 and 7 (for scaling the pipeline size, we simply composed the described pipeline with itself). However, our own implementation is consistently at least two times faster, and for the larger pipelines the difference is even more apparent. While the relatively low processing speed per megabyte can be explained by the need to create many small documents (the element per megabyte density is much higher compared to the previous tests), it is surprising to see an implementation specifically designed for processing XPROC be outperformed by our unoptimised implementation of the pipeline steps.

Document size	X-Fun	CALABASH
2 MB	8.7 s	16.6 s
4 MB	15.3 s	32.6 s
6 MB	23.1 s	51.8 s
8 MB	39.5 s	78.7 s

**Fig. 6.** Performance of X-Fun and CALABASH on a fixed pipeline with varying input tree size

Pipeline size	X-Fun	CALABASH
1	8.7 s	16.6 s
2	12 s	75.8 s
3	16 s	136.6 s
4	22 s	198.6 s

**Fig. 7.** Performance of X-Fun and CALABASH on a 2 MB document with varying pipeline size

## 6 Conclusion and future work.

We have presented X-Fun, a language for processing data trees and shown that can serve as a uniform programming language for XML processing and as a uniform core language for implementing XQUERY, XSLT, and XPROC on top of any existing XPATH evaluator. Our implementation based on SAXON’s in-memory XPATH evaluator yields surprisingly efficient implementations of the three W3C standards, even there is a lot of space left for optimisation. We have obtained results which are a match for the SAXON’s XQUERY and XSLT

evaluators and in the case of XPROC, first results show that we are already faster than CALABASH.

Our prime objective in future is to build streaming implementations of X-Fun, and thus of XQUERY, XSLT, and XPROC. The main ideas behind it are described in a technical report [9]. These streaming implementation will serve in the tools called QUIXQUERY, QUIXSLT, and QUIXPROC. A first version of QUIXSLT is freely available for testing on our online demo machine [6] while streaming is not yet available for our current QUIXPROC implementation.

**Acknowledgement.** We would like to thank Guisepe Castagna and Kim Nguyen for their helpful discussions about the type system of X-Fun.

## References

1. Arroyuelo, D., et al.: Fast in-memory xpath search using compressed indexes. In: ICDE. pp. 417–428. IEEE (2010)
2. Castagna, G., Im, H., Nguyen, K., Benzaken, V.: A Core Calculus for XQuery 3.0 (2013), <http://www.pps.univ-paris-diderot.fr/~gc/papers/xqueryduce.pdf>, unpublished manuscript
3. Frisch, A., Nakano, K.: Streaming XML transformation using term rewriting. In: Programming Language Technologies for XML (PLAN-X). pp. 2–13 (2007)
4. Fülöp, Z., Vogler, H.: Syntax-Directed Semantics – Formal Models based on Tree Transducers. EATCS Monographs in Theoretical CS, Springer (1998)
5. Hakuta, S., Maneth, S., Nakano, K., Iwasaki, H.: XQuery Streaming by Forest Transducers. In: ICDE. pp. 952–963. IEEE (2014)
6. Innovimax, INRIA Lille: Quix tools suite, <https://project.inria.fr/quix-tool-suite/>
7. Kay, M.: XSL Transformations (XSLT) Version 3.0. W3C Last Call Working Draft (2013), <http://www.w3.org/TR/xslt-30>
8. Kepser, S.: A simple proof for the Turing-Completeness of XSLT and XQuery. In: Extreme Markup Languages® (2004)
9. Labath, P., Niehren, J.: A Functional Language for Hyperstreaming XSLT. Research report (Mar 2013), <http://hal.inria.fr/hal-00806343>
10. Labath, P., Niehren, J.: A Uniform Programming Language for Implementing XML Standards. Research report (Jan 2015), <http://hal.inria.fr/hal-00954692>
11. Maneth, S., Berlea, A., Perst, T., Seidl, H.: XML type checking with macro tree transducers. In: PODS. pp. 283–294. ACM-Press, New York, USA (2005)
12. Neumann, A., Seidl, H.: Locating matches of tree patterns in forests. In: FSTTCS. pp. 134–145 (1998)
13. Robie, J., et al.: XQuery 3.0: An XML Query Language. W3C Proposed Recommendation (2013), <http://www.w3.org/TR/xquery-30>
14. Saxonica: SAXON 9.5: The XSLT and XQuery Processor, <http://saxonica.com>
15. Schmidt, A., et al.: XMark: A benchmark for XML data management. In: VLDB (2002), <http://www.ins.cwi.nl/projects/xmark/Assets/xmlquery.txt>
16. Walsh, N.: XML Calabash, <http://xmlcalabash.com>
17. Walsh, N., et al.: XProc: An XML Pipeline Language. W3C Recommendation (2010), <http://www.w3.org/TR/xproc>
18. Zergaoui, M., Innovimax: QuiXProc, <https://project.inria.fr/quix-tool-suite/quixproc/>