# Encoding XQuery Using *System F*

by

Yun Xia

A thesis

presented to the University of Waterloo

in fulfillment of the

thesis requirement for the degree of

Master of Mathematics

in

Computer Science

Waterloo, Ontario, Canada, 2005

**Author's Declaration for Electronic Submission of a Thesis**

I hereby declare that I am the sole author of this thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners.

I understand that my thesis may be made electronically available to the public.

# Abstract

Since the World Wide Web Consortium (W3C) has recommended XQuery as the standard XML query language, the interest in using existing relational technology to query the XML data has dramatically increased. The most significant challenge of the relational approach is how to fully support XQuery semantics in XQuery-to-SQL translation. To eliminate the implicit semantics of XQuery, an XQuery fragment must be defined with simple syntax and explicit semantics. XQ [6] is proposed as an XQuery fragment to express XML queries.

In this thesis, XQ is intensively investigated. It is encoded by *System F*, a second-order lambda calculus with a proven expressive power and a strong normalization property. Since XML data is defined as inductive data types, XML tree and XML forest, in *System F*, all basic XML operators in XQ have been successfully encoded. Also, the semantics of XQ are represented in *System F* where XQ's semantics environment is encoded by an *Environment* data type with the corresponding operators. The successful encoding of XQ by *System F* ensures the termination of XQ query evaluation.

Moreover, an extension of XQ by a new tree operator `Xtree` and a vertical `Vfor` clause is proposed in this thesis to express the *undefinable* XQ queries. It is demonstrated that this extension still allows XQ to retain its XQ-to-SQL translation property that ensures the polynomial evaluation time complexity, and its *System F* encodable property that ensures the termination of query evaluation.

iii

# Acknowledgments

I would like to acknowledge the support of several people who have made this thesis possible. First and foremost, I am very grateful to Dr. David Toman, my supervisor, for his great guidance and support during my graduate studies. His patience, valuable suggestions and comments have directly facilitated the completion of this thesis.

Also, I would like to thank Dr. Richard Trefler and Dr. Grant Weddell for reviewing my thesis and providing their valuable advice and comments.

I am particularly grateful to my husband, Xin, for his love and understanding, and to my son, Steven, for the happiness he has brought to my heart, and to my parents for their wisdom and encouragement.

# Dedication

To Xin for his love

# Contents

# Chapter 1

# Introduction

Since XML (Extensible Markup Language) has rapidly grown to be the universal format standard for exchanging data among disparate applications on the internet, it is crucial to develop a way to query and interact with XML documents.

## 1.1 XQuery

"A successful query language can enhance productivity and serves as an unifying influence in the growth of an industry" [9]. As a programming language, which depends on the tree structure of XML data, there are many desirable qualities that an XML query language should exhibit:

- Support the operations on the document order and the *axis* expressions, which are used to navigate the tree structure of XML data, and to retrieve the context of a particular document fragment;

1

- Provide the ability to express various combinations of multiple XML documents and construct new document based on the query result;

- Offer full compositionality so that the XML query operators can be compose with a full generality; for example, the result of one expression can be used as the input of another expression

As XML has grown in popularity, there has been significant research in the area of XML query languages. As a result, XQuery has been defined and recommended by W3C as the standard of XML query language to support XML applications, such as data processing, transformation, and querying tasks.

As a fully compositional and typed functional language, XQuery provides the ability to flexibly select, recombine, and restructure XML documents and fragments. XQuery is characterized by the following constituents:

- XPath expression traverses and extracts node sequences from document(s);

- FLOWR (for-let-order-where-return) expression binds selected node sequences (in order) to variables and expresses joins and filter conditions;

- Expressions that construct new XML documents or values from the binding variables returned by the FLOWR expressions;

- There is a large library of functions and operators.

To illustrate XQuery, bibliography data is taken from the XML Query Use Cases in Figure 1.1 [1].

```
<bib>
    <book year="1994">
       <title>TCP/IP Illustrated</title>
       <author><last>Stevens</last><first>W.</first></author>
        <publisher>Addison-Wesley</publisher>
        <price>65.95</price>
    </book>
    <book year="1992">
        <title>Adv. Programming in the Unix environment</title>
        <author><last>Stevens</last><first>W.</first></author>
        <publisher>Addison-Wesley</publisher>
        <price>65.95</price>
    </book>
    <book year="1999">
        <title>The Economics of Tech. for Digital TV</title>
        <editor>
                <last>Gerbarg</last><first>Darcy</first>
                 <affiliation>CITI</affiliation>
        </editor>
        <publisher>Kluwer Academic Publishers</publisher>
    </book>
</bib>
```

Figure 1.1: Bibliography data

The document contains one `bib` element that contains three `book` elements. This document is well formed because each open element tag has a corresponding closing element tag and the elements are properly nested. A common querying task is to translate the structure and content of XML data, which requires the construction of XML values. Suppose our example is in a file named `bib.xml`, the following expression, utilizing element constructors and FLOWR expressions, queries the titles and years of all books published by Addison-Wesley after 1991 in alphabetical order:

```
<bib>{
    for $b in doc("bib.xml")//book
    where $b/publisher = "Addison-Wesley" and $b/@year > 1991
    order by $b/title
    return
        <book>
            { $b/@year }
            { $b/title }
        </book>
}</bib>
```

Figure 1.2: XQuery example for the XML document in Figure 1.1

In the query expression in Figure 1.2, the `for` clause generates a sequence of `book` elements and binds the variable `$b` to each element; the `where` clause filters the sequence by retaining the books that publisher is "Addison-Wesley" and the published year is later than 1991; the `order by` clause sorts the resulting book sequence; the inner `return` clause returns a sequence of new `book` elements, each containing an attribute `year`, and an element `title`, extracted from the surviving `book` elements in the ordered sequence; the results of the entire FLOWR expression are composed in an

XML element `bib`, whose elements are the sequence returned from the inner `return` clause.

## 1.2 Challenges

The XQuery language provides such a strong expressive power that XML queries are easy to write, but XQuery complex and redundant syntax prevents an efficient XQuery implementation. Moreover, XQuery has been proven to be Turing-complete, which means the termination of its query evaluation is undecidable. A well-defined subset (fragment) of the XQuery language with a simple syntax and ensuring the query evaluation termination is pivotal to efficient XQuery implementation and optimization. This XQuery *fragment* must describe XML query semantics precisely by defining a data model and a corresponding set of simple operations. Defining the ideal subset or fragment of XQuery has proven to be difficult. Many XQuery subset have been proposed, but have problems with one of these requirements. Currently, there is no universally accepted XQuery subset which holds strong expressive, and meanwhile, ensures its evaluation termination.

## 1.3 Related Work

There are two key approaches in designing an XQuery fragment: a set-based algebra or a iterator-oriented XQuery fragment. The first approach designs an algebra with operators such as `OP(a, b)`, and the latter approach follows

the syntax of XQuery with a format

```
for $i in a, $j in b return (OP $i, $b)
```

for the operator, OP. Although the set-based algebra is easier to reason about algebraically, the iterator-oriented language is generally easier to match to the existing XQuery syntax and capture its semantics.

A proper XQuery fragment is essential for efficient XML query evaluation, but what is right for XML Query processing is still an open issue. There are several good candidates but most of them display shortcomings, such as not recognizing that XML is an ordered data structure, or not ensuring the termination of query evaluation.

The rest of this section discusses about several existing XQuery fragments or XML algebras proposed in XML query research area.

**Formal Semantics.** The XQuery Core is defined by the W3C XML Query Working Group as part of XQuery Formal Semantics document [5]. As a proper subset of XQuery, the XQuery Core is proposed in the implementation of an XQuery processor called Galax [4]. Here, the evaluation process starts with normalization which translates the XQuery expressions into the XQuery Core syntax. We illustrate the core mapping with a query for the bibliography data in Figure 1.1:

```
doc("book.xml")//book[publisher = "Addison-Wesley"],
```

which asks for the books whose publisher is "Addison-Wesley". This query expression is normalized into the XQuery Core in Figure 1.3.

As an iterator-oriented language, Formal Semantics uses `for` loops (Line 2, Line 4, and Line 7) to rewrite the path expressions to eliminate the im-

```
1.  fs:node-sequence(
2.    for $fs:dot in
3.        fs:node-sequence(
4.           for $fs:dot in doc("bib.xml")
5.    return descendant-or-self::node())
6.  return
7.        for $fs:dot in fs:node-sequence(child::book)
8.        return
9.        if (some $v1 in fn:data(child::publisher) satisfies
10.             fn:boolean(op:equal($v1,"Addison-Wesley")))
11.       then $fs:dot
12. else ())
```

Figure 1.3: XQuery Core expression for the query doc("book.xml")//book[publisher = "Addison-Wesley"]

plicit semantics of XPath. As a functional programming language, Formal Semantics designs a full semantic mapping for XQuery, and provides a set of logical optimization rules based on the presence of a schema and static type inference. Without type information, its rewriting is not feasible and makes further optimization difficult.

**TAX.** As a set-based algebra, TAX [8] is a tree-based algebra with a set of operations similar to the relational algebra operations. TAX algebra treats the XML document as a labelled order rooted tree, and decomposes a query into one or more pattern trees, shown in Fig 1.4. Each edge denotes a structural inclusion relationship between the elements, represented by the respective pattern tree nodes. The inclusion relationship can be specified to be either *pc* (parent-child relationship) or of *ad* (ancestor-descendent relationship). In a pattern tree, there is also a boolean predicate, applied

to the nodes in the pattern tree that contains the constrain information for the nodes in the pattern tree. For example, the query asks for the books published before 1993 with the publisher "Addison-Wesley", whose XQuery expression and corresponding pattern tree are reflected by Figure 1.4.

```
doc(book.xml)//book[publisher = "Addison-Wesley" AND year < 1993]
```



$1.tag = book &
$2.tag = publisher
$2.content = "Addison-Wesley" &
$3.tag = year &
$3.content < 1993

Figure 1.4: Pattern Tree for Query "doc(book.xml)//book[publisher = "Addison-Wesley" AND year < 1993]

The pattern tree expresses a projection operation, $\pi_{\mathcal{P}, PL}(\mathcal{C})$, which takes the collection, $\mathcal{C}$ (XML document), as an input, and a pattern tree, $\mathcal{P}$, and projection list, $PL$, as parameters, where a projection list is a list of node labels, appearing in pattern $\mathcal{P}$, which is $1 in Figure 1.4. By using relational algebra as a guide, TAX defines a set of operators such as *selection*, *projection*, *production*, or *grouping*, that takes a collection trees as the input and produces a set of trees as the output. The query evaluation process involves a sequence of pattern tree matchings that breaks up the input tree into small pieces, then reassembles them by grouping and projection. Because of the semantic mismatching of FLWR and TAX, even TAX is efficient in a simple XPath query evaluation, the redundant pattern matching is unavoidable when dealing with complex XQuery queries with a nested structure. An extract effort has to be taken to combine query results from multiple tree pattern queries with extra pattern matching and extra joins.

**XQ.** DeHaan et al. [6] propose a dynamic interval method that combines the relational mapping implementation with presenting an XQuery fragment XQ[1]. The most significant challenge of the relational approach in XQuery evaluation is how to fully support the XQuery semantics in the XQuery-to-SQL translating process so that the document order is retained. Most relational methods fail to provide this support, and implement only simple XPath query evaluations. The dynamic interval method proposes a novel solution that is based on a XML-specific XQ query language, which captures most of the semantics of XQuery with a fairly simple syntax. As an iterator-oriented language, XQ models an XML database as an ordered forest whose elements are rooted, node-labelled ordered trees. A simple set of operations is defined to manipulate the XML forest. Moreover, XQ keeps the FLWR syntax, which is the most powerful and core part of XQuery, and explicitly expresses its semantics using an *environment* for the variables defined in let and for expressions with their resulting values. For example, the query for the bibliography in Figure 1.1, which returns the books whose publisher is "Addison-Wesley", is expressed by XQ in Figure 1.5.

The dynamic interval method describes the interval representation of XML data and its intermediate query results as relations, which seamlessly preserve the document order required by XQuery semantics. After the XQuery expressions are represented by XQ without losing its semantic meaning, a set of templates that translates the basic operations of XQ into single SQL state-

---

[1] [6] does not give the algebra a formal name. In the thesis, it is called XQ for convenience.

```
for $a = select('book', subtreedfs(doc(book.xml))) do
let $b = select('publisher', children($a)) in
     $c = select("Addison-Wesley", children($b))
where not empty $c
return $a
```

Figure 1.5: XQ for "doc(book.xml)//book[publisher = "Addison-Wesley"]

ment is proposed based on the interval information of the XML document. Therefore, the XML queries, with the XQ-to-SQL translation, can be evaluated by using an extended SQL engine to achieve the polynomial evaluation time complexity. Even the XQuery translation targets a relational implementation, but its XQ definition can also be used within the native XML database system. The present work in this thesis continues the work of [6], and focuses on the XQ language aspect part.

## 1.4 Contributions

In the thesis, the theoretical aspects of the XQ query language, including its computability and expressiveness, with respect to its XQ-to-SQL translation property, are investigated. First, XQ is encoded using *System F*. As a second order typed lambda calculus, *System F* has a considerable expressive power [7] and a *strong normalization* property. The successfully encoding of XQ in *System F* indicates that the termination of XP query evaluation, and the decidability of query equivalence. When the expressiveness of XQ are considered, a new tree operator Xtree and a vertical Vfor clause are proposed to express the *undefinable* XQ queries. In the extension of XQ, the intention

is to retain its XQ-to-SQL translation property to achieve polynomial time complexity by utilizing relational query engine, and also its *System F* encodable property to ensure the termination of XQ query evaluation. The main contributions of the thesis are summarized as follows:

- represents XML data in *System F* as inductive data types, XML tree and XML forest, with the corresponding iterators;

- encodes the basic XML operators of XQ by *System F*, and represents the semantic equations of XQ by encoding the semantics environment using an *Environment* data type with a set of operations;

- analyzes the limitations of XQ, and proposes an extension of XQ with a vertical constructor `Xtree` and a `Vfor-do` clause;

- designs the corresponding XQ-to-SQL templates for the extension in order to utilize the relational engine to achieve a polynomial evaluation time complexity;

- encodes the extended XQ in *System F* to ensure the termination of XQ query evaluation

## 1.5  Thesis Organization

This chapter has introduced the framework of XQuery implementation by both native and relational approaches. The focus is on the challenges and research work for designing a suitable and simple XQuery fragment. The rest of the thesis is organized as follows. Chapter 2 reviews the background

and techniques of the dynamic interval method, including the XQ syntax and semantics. A brief introduction to *System F*, which is a second-order lambda calculus with an expressive power and a strong formalization property, is provided in Chapter 3. Simple data types, such as integer and list, are illustrated to describe the translation schema. Chapter 4 defines the XML document as inductive types in *System F* and represents them as lambda abstractions. The basic XML operators of XQ are encoded in *System F*, and the semantic equations of XQ are represented in *System F* by encoding the semantics environment with a set of operations of an *Environment* data type. Chapter 5 states the expressive limitations of XQ, and proposes an extension with a vertical tree constructor Xtree and a new vertical Vfor clause, which is also XQ-to-SQL translatable and *System F* encodable. Chapter 6 summarizes the thesis and suggests future work.

# Chapter 2

# Dynamic Interval Method

In the thesis, XQ, an XQuery fragment proposed for the dynamic interval method [6] is examined. The dynamic interval method aims to implement XML queries in a relational approach. Its first step is to define the XQuery fragment, XQ, to eliminate redundant syntax and implicit semantics of XQuery. Not only does XQ capture the core features of XQuery with a set of basic XML operators, but also it supports successful XQ-to-SQL translation. Although XQ supports the relational-based XML queries, it can also be used in the native XML query processors. In this chapter, a brief overview of the design concept, assumptions, and techniques, on which the thesis works build, is presented.

## 2.1   XQ: An XQuery Fragment

As a query language, XQ models an XML database as an ordered forest whose elements are rooted, labelled and ordered trees; also a simple set of

operators is presented to manipulate the forest.

## 2.1.1   XML Data Model and Operations

**Definition 1 (XML Forest)** An XML document can be described as an ordered forest XForest, and is defined inductively as

$$XForest \stackrel{\text{def}}{=} [\,] \ \mid \ [\texttt{<s>} \ XForest \ \texttt{</s>}] \ \mid \ XForest \ \texttt{@} \ XForest$$

where $s$ is a string, $[\,]$ denotes an empty forest, $[\texttt{<s>} \ XForest \ \texttt{</s>}]$ signifies a forest containing a single tree element with a root node labelled $s$ and an ordered forest XForest as its children, and $XForest \ \texttt{@} \ XForest$ indicates a new forest concatenated by two forests.

The XML forest describes an XML document without distinguishing between node identity and node types (element, attribute, and text). However, such features can be easily added by additional encoding conventions that relate either to node labelling or to subtree patterns: the text leaf node with CDATA text can be encoded by labelling the node with "text:"; an element eleName is encoded by using the label "<eleName>". An attribute of an element can be treated as a subtree of this element, and is labelled "attribute:", where as the attribute value is encoded as subtree of its attribute node and labelled as "text:". For example, the XML document in Figure 1.1 can be encoded with node type information to distinguish the element, attribute, and text nodes. Figure 2.1 depicts part of the converted XML data.

Since XML data is described as an ordered forest, *XForest*, XQ provides a set of basic operations to manipulate the forest data model so that it can navigate the structure of *XForest* data, retrieving the context of particular

```
<bib>
    <book>
      <attribute:year> <text:1994> </attribute:year>
      <title> <text:TCP/IP Illustrated/></title>
      <author>
         <last> <text:Stevens/> </last>
         <first> <text:W./> </first>
      </author>
      <publisher> <text:Addison-Wesley/> </publisher>
      <price> <text:65.95/> </price>
    </book>
    ...
</bib>
```

Figure 2.1: Converted bibliography data with encoded node information

*XForest* fragments, and joining multiple XForest forests to construct a new *XForest* document.

**Constructors.** *XForest* data (XML forest) can be empty (no tree element). Also, it can be constructed from a concatenation of two XML forests, or from an element constructor that adds a labelled root to an XML forest. The constructors and corresponding types [1]are described as follows:

**empty constructor** [ ] : *XForest*, constructs an XML forest without elements;

**element constructor** *Xnode* : *String → XForest → XForest*, which takes a string as a root label and an XML forest as children to

---

[1]In the thesis, curried functions is used instead of uncurried functions presented in [6]. The curried function $f$ of type $int \to int$ means that the function application $f\ a$ with an integer argument $a$ returns a result of integer.

construct an XML forest with a single tree node[2];

**concatenation** @ [3]: *XForest → XForest → XForest*, takes two XML
forests as arguments to construct a new XML forest by appending
the second argument to the first one, keeping the inner order of
each argument

**Horizontal Operators.** XQ defines a set of horizontal operators that ma-
nipulate the *XForest* data (XML forest) at the top element (tree) level
without recursively traversal into the structure of the trees.

**head** : *XForest → XForest*, returns the first element tree from a for-
est. Strictly, the return value should be an XML tree instead of
an XML forest. As a result, a new type *XTree* (XML Tree), in
Chapter 4, is added;

**tail** : *XForest → XForest*, takes an XML forest, and returns a new
forest without the first element. The return forest still keeps the
original order of the input forest;

**reverse** : *XForest → XForest*, reverses the element order in an XML
forest and keeps the original structure of each tree element un-
changed;

**select** : *String → XForest → XForest*, takes a string, *s*, and an XML
forest as arguments, and returns sub-forests whose tree elements

---

[2]In the thesis, the element constructor is defined to be a tree constructor that takes a
labelled node and a forest as arguments and returns a tree

[3]In the thesis, another constructor, concat, which takes a tree and a forest to construct
a new forest, is adopted instead of concatenation

have the root node labelled $s$. For example,

```
select('bib', doc(bib.xml))
```

returns the element trees whose root labels "`bib`" exists; otherwise, an empty forest is returned;

**distinct** : *XForest* → *XForest*, filters out duplicate elements in an XML forest. Although `distinct` is defined as a horizontal operator, it still needs an auxiliary boolean function `equal` that recursively compares two trees' structural relationships;

**sort** : *XForest* → *XForest*, takes an XML forest and returns a sorted forest. The `sort` operator also needs an auxiliary boolean function `less` to recursively compare two trees' structural order

**Vertical Operators.** In order to navigate the structure of each tree element of the forest, XQ provides vertical operations such as `children` and `subtreedfs` to locate the desired sub-element nodes.

**roots** : *XForest* → *XForest*, returns all the tree elements of a forest at the top level;

**children** : *XForest* → *XForest*, acting as path axis "/", selects children nodes and returns them in the original (document) order;

**subtreedfs** : *XForest* → *XForest*, acting as path axis "//" (descendent-or-self), selects self and sub-elements, and returns them in an DFS (document) order as an XML forest

**Boolean Operators.** XQ provides two binary boolean operators, `equal` and `less`, to compare the structural relation of trees. [4]

## 2.1.2 Syntax of XQ

With denoting XML data as an XML forest, dynamic interval method describes the XML query in a simple syntax, with which explicit or implicit semantics of an XQuery expression can still be captured.

**Definition 2 (XQ)** An XML query expression can be expressed by the following BNF rules:

$$
\begin{aligned}
e \quad ::= \quad & x \\
| \quad & \mathsf{XF_n}(e_1, \ldots, e_n) \\
| \quad & \mathsf{let}\ x = e_1\mathsf{in}\ e_2 \\
| \quad & \mathsf{where}\ \varphi\ \mathsf{return}\ e \\
| \quad & \mathsf{for}\ x\ \in e\ \mathsf{do}\ e'
\end{aligned}
$$

where $e$, $e_1$, $\ldots$, and $e_n$ are query expressions, $\mathsf{XF_n}$ is a basic operator defined in the previous section, $x$ is a variable, and $\varphi$ is a boolean condition.

## 2.1.3 Semantics of XQ

The semantics of XML queries are expressed by denotational semantics style equations. The denotational semantics introduces the concept of an *Environment*, which maps each free variable of an expression to its value. Since the behaviour of a query evaluation is a computation from an initial state to

---

[4][6] denotes `equal` and `less` as XML forest operators. Since we intend to compare two trees' structural relationship, they are denoted as tree operators in Chapter 4.

a final state through a sequence of intermediate states (results), the evaluation denotes as a *transformation* function $[\![ - ]\!]$, where $E' = [\![e]\!]E$ means the new state $E'$, an instance of *Environment*, is transformed from state $E$ by the computation of the expression $e$. Therefore, the semantics of FLWR-like expressions is described as semantic equations as follows:

**SEM EQ: Variables**

$$[\![x]\!]E \overset{\text{def}}{=} E\ (x)$$

The role of variable $x$ is to obtain the value (XML forest) of variable $x$ from *Environment* $E$ without changing $E$;

**SEM EQ: Let-Assignment**

$$[\![\textsf{let } x = e\textsf{in } e'\ ]\!]E \overset{\text{def}}{=} [\![e']\!](E[x = ([\![e]\!]E)])$$

Expression $e'$ is evaluated in $E'$, which is a new *Environment* extended from $E$ by binding the variable $x$ to the result of $e$;

**SEM EQ: Where-Return**

$$[\![\textsf{where } \varphi \textsf{ return } e]\!]E \overset{\text{def}}{=} \textsf{ if } ([\![\varphi]\!]E) \textsf{ then } ([\![e]\!]E) \textsf{ else } [\,]$$

where $\varphi$ is a boolean condition. If condition $\varphi$ holds, expression $e$ is evaluated in *Environment* $E$;

**SEM EQ: For-Do**

$$[\![\textsf{for } x\ \in\ e \textsf{ do } e']\!]E \overset{\text{def}}{=} [\![e']\!](E[x := v_1])\ @\ldots@\ [\![e']\!]E([x := v_k])$$

$$\textsf{where } [v_1,\ \ldots,\ v_k] = [\![e]\!]E$$

The expression for $x \in e$ do $e'$ first evaluates expression $e$ in *Environment* $E$ and obtains a list of results $[v_1, \ldots, v_k]$; and then uses assignment operation "$x \in$" to transform *Environment* $E$ into a list of new *Environments* $[E_1, \ldots, E_k]$. Finally expression $e'$ is computed in these new environments, and the results are concatenated.

**SEM EQ:** XF$_\mathtt{n}$

$$[\![\mathtt{XF_n}(e_1, \ldots, e_k)]\!]E \stackrel{\text{def}}{=} \mathtt{XF_n}\ ([\![e_1]\!]E, \ldots, [\![e_k]\!]E)$$

where XF$_\mathtt{n}$ is an XML operator defined in Section 2.1.1, such as `select`, `subtreesdfs`, with $k$ arguments. The $k$ arguments are evaluated simultaneously in *Environment* $E$, and the results are fed to the function XF$_\mathtt{n}$.

## 2.1.4  XQ Translation of XQuery

As an XQuery fragment, XQ captures all the detailed intricate nuances presented in the XQuery's FLWR and XPath expressions. Its syntax allows an arbitrary composition of the basic function invocations, local variable definitions, filtering by boolean conditions, and iteration over tree elements in XML forests [6]. Some XQueries are used for the bibliography data in Figure1.1 to illustrate the XQ translation.

**Example 2.1.1** List the books published by Addison-Wesley, including their year and title.
XQuery expression:

```
<bib>
 {
```

```
  for $b in doc("bib.xml")/bib/book
  where $b/publisher = "Addison-Wesley"
  return
    <book year="{ $b/@year }">
     { $b/title }
    </book>
 }
</bib>
```

corresponding XQ expression:

```
for $b in select(book, children(select(bib,
                                       children(doc(bib.xml))))) do
where not empty (select(Addison-Wesley, children(publisher,
                                                 children($b))))

return
  let $3 =
    let $1 = children($b)
    in (let $2 = Xnode(year, (select(year, children($1)))
             @ Xnode(title, (select(title, $1)))
        in Xnode(book, $2))
  in Xnode(bib, $3)
```

**Example 2.1.2** Create a flat list of all the title-author pairs with each pair enclosed in a "result" element.
XQuery expression:

```
<results>
  {for $b in doc(bib.xml)/bib/book,
       $t in $b/title,
       $a in $b/author
    return
       <result>
           {$t}
           {$a}
       </result>}
</results>
```

corresponding XQ expression:

```
let $2 = for $b in select(book, children(select(bib,
                  children(doc(bib.xml)))))
        do let $1 = children($b) in
           Xnode(result, (select(title, $1) @ select(author, $1)))
in Xnode(results, $2)
```

## 2.2 Dynamic Interval Relational Method

### 2.2.1 Interval Encoding

In order to manipulate XML documents in relational database, an interval-encoding method is proposed to capture enough information about XML data by encoding an XML document as triples $(s, l, r)$, where $s$ is the string of the node name (node $s$), $l$ is the left endpoint of node $s$, and $r$ is the right endpoint of node $s$. The way to generate the interval encoding is to perform a depth-first traversal of the XML tree by using a counter to mark the $l$ value of the node when it is visited first time, and to mark its $r$ value when the node is visited the last time. The interval encoding relation of bibliography XML data in Figure 1.1 is shown in Figure 2.2.

### 2.2.2 Dynamic Intervals

To evaluate an FLWR expression by a fixed relational query, a *dynamic* interval encoding is presented to represent the sequence of environments generated within an query evaluation process. A detailed definition of dynamic interval can be found in [6].

**Example 2.2.1** Figure 2.3 shows the query results for the path expression,

| s | l | r | s | l | r | s | l | r |
|---|---|---|---|---|---|---|---|---|
| bib | 0 | 79 | book | 25 | 47 | book | 52 | 78 |
| book | 1 | 24 | year | 26 | 29 | year | 53 | 56 |
| year | 2 | 5 | 1992 | 27 | 28 | 1999 | 54 | 55 |
| 1994 | 3 | 4 | title | 30 | 33 | title | 57 | 60 |
| title | 6 | 9 | Adv. Prog. | 31 | 32 | The Economics of | 58 | 59 |
| TCP/IP Illu. | 7 | 8 | author | 34 | 43 | editor | 61 | 47 |
| author | 10 | 19 | last | 35 | 38 | last | 62 | 65 |
| last | 11 | 14 | Stevens | 36 | 37 | Gerbarg | 63 | 64 |
| Stevens | 12 | 13 | first | 39 | 42 | first | 66 | 69 |
| first | 15 | 18 | W. | 40 | 41 | Darcy | 67 | 68 |
| W. | 16 | 17 | publisher | 44 | 47 | affiliation | 70 | 73 |
| author | 10 | 19 | Addison-Wesley | 45 | 46 | CITI | 71 | 72 |
| price | 20 | 23 | price | 48 | 51 | publisher | 74 | 77 |
| 65.95 | 21 | 22 | 65.95 | 49 | 50 | Kluwer Academic | 75 | 76 |

Figure 2.2: encoded Bibliography data

$T_{title} = $ `doc(bib.cml)/bib/book/title`, with the index relation $I$, representing the initial environment with $i = 0$.

### 2.2.3  XQ-to-SQL Translation

Based on the semantics of XQ expressions, SQL translation templates, including the templates for basic operations and those for FLWR patterns, are adopted to translate an XQ-syntax query expression into a relational query statement by composing SQL templates from the sub-query fragments. The detailed XQ-to-SQL translation can be found in [6].

**Proposition 1** For each XQ expression, there is a corresponding SQL query statement, such that all XQ queries can be implemented in the relational

| I | | $T_{title}$ | | |
|---|---|---|---|---|
| i | | s | l | r |
| 0 | | title | 6 | 9 |
| | | TCP/IP Illustrate | 7 | 8 |
| | | title | 30 | 33 |
| | | Advanced Programming in the Unix environment | 31 | 32 |
| | | title | 57 | 60 |
| | | The Economics of Technology and Content for Digital TV | 58 | 59 |

Figure 2.3: Titles in an initial environment

query engine.

By defining a simple XQuery fragment, XQ, to capture the core features of XQuery, and a set of XQ-to-SQL translation templates, the dynamic intervals method successfully evaluates the XML queries using the relational query engine, which ensures the evaluation polynomial time complexity.

# Chapter 3

# System F

In the thesis, *System F*, a second order lambda calculus, is adopted to represent XQ presented in [6]. In this chapter, we introduce *System F*, whose generality and strong expressive power supports faithful encoding of arbitrary composition of XQ expressions. To understand *System F*, *lambda calculus*, which underlies all functional languages and the procedural facilities of many more general languages, is introduced.

## 3.1  Lambda Calculus

Originally developed and investigated in the 1930's by logicians such as Church, Curry, and Rosser, lambda calculus has influenced the design of programming languages. Many problems of language design and implementation, especially those concerning procedure mechanisms or type structures, can be posed and investigated more easily by the lambda calculus than in

more complex languages. For example, for the long and repetitive expression,

$$(5 \times 4 \times 3 \times 2 \times 1) + (7 \times 6 \times 5 \times 4 \times 3 \times 2 \times 1) - (3 \times 2 \times 1)$$

we rewrite it to be `factorial` $5 +$ `factorial` $7 -$ `factorial` $3$, where

$$\texttt{factorial } n = \textsf{if } n = 0 \textsf{ then } 1 \textsf{ else } n \times \ (\texttt{factorial } n - 1)^1$$

For each nonnegative number $n$, instantiating the function `factorial` with the argument $n$ yields the factorial of $n$ as a result. If "$\lambda\, n.\ \ldots$" is an abbreviation for "the function that, for each $n$, yields $\ldots$", the definition of `factorial` is restated as:

$$\texttt{factorial} = \lambda\, n.\textsf{if } n = 0 \textsf{ then } n \times (\texttt{factorial } n - 1)$$

where the subexpression, $\lambda\, n. \ldots$, is called an *abstraction* or *lambda expression*. The *function application* `factorial` $0$ refers to "the result when the argument variable $n$ in the function body $(\lambda\, n.\textsf{if } n = 0 \textsf{ then } \ldots)$ is replaces by $0$"; that is, $1$.

In lambda calculus, a function is defined in the *abstraction* form, which does not specify the name of the function. The capability of denoting functions without giving them names makes abstractions particular useful for describing functions that are arguments for other "higher order" functions.

### 3.1.1   Syntax

The syntax of lambda calculus is comprised of three types of expressions (*terms*). A variable $x$ by itself is a term; an abstraction of a variable $x$ from

---

[1] Here the *curried form* "`factorial` $n$" is used instead of the *uncurried form* "`factorial` $(n)$" to deal with the higher order function that can treat functions as arguments.

a term $t_1$, written $\lambda\,x.t_1$, is a term; and the application of a term $t_1$ to another term $t_2$, written $t_1\ t_2$, is a term. These ways of forming terms are summarized in the following abstract syntax:

$$
\begin{aligned}
t ::=\ &x && \textit{variable}\\
\mid\ &\lambda\,x.t_1 && \textit{abstraction}\\
\mid\ &t_1\ t_2 && \textit{application}
\end{aligned}
$$

It is assumed that the application is left associative, that is,

$$t_1\ t_2\ t_3 \ldots\ t_n \equiv (\ldots ((t_1\ t_2)\ t_3)\ldots t_n)$$

and in $\lambda\,x.t$, the subterm $t$ extends the first stop symbol or to the end of the enclosing phrase. For example,

$$\lambda\,x.(\lambda\,y.xyx)x \equiv \lambda\,x.((\lambda\,y.((xy)x))x)$$

**Definition 3 (Bound and Free Variable)** In the term $\lambda\,x.t$, the occurrence of $x$ is a binder with scope $t$. An occurrence of variable $x$ is said to be *bound*, when it occurs in body $t$ of an abstraction $\lambda\,x.t$. An occurrence of $x$ is *free* if it appears in a position where it is not bound by an enclosing abstraction of $x$.

**Example 3.1.1** The occurrence of $x$ in $xy$ and $\lambda\,y.xy$ are free, whereas the ones in $\lambda\,x.x$ and $\lambda\,z.\,\lambda\,x.\,\lambda\,y.x(y\ z)$ are bound. In $(\lambda\,x.x)\ x$, the first occurrence of $x$ is bound and the second one is free.

The set of free variables in term $t$, denoted as $FV(t)$, is defined as

$$
\begin{aligned}
FV(x) &= \{x\}\\
FV(t_1 t_2) &= FV(t_1) \cup FV(t_2)\\
FV(\lambda\,x.t) &= FV(t) - \{x\}
\end{aligned}
$$

and term $t$ is said to be *closed*, if and only if $FV(t) = \emptyset$.

**Definition 4 ($\alpha$-conversion)** $\alpha$-conversion is the operation of replacing an occurrence of a bound variable by

$$\lambda\, x.t \equiv_\alpha \lambda\, y.t[x \mapsto y]$$

where $y \notin FV(t)$

Terms $t$ and $t'$ are $\alpha$-equivalent with $t \equiv_\alpha t'$, if $t'$ results from $t$ by a series of changes of the bound variables. The $\alpha$-conversion states that the renaming should preserve the semantic meaning.

## 3.1.2 Operational Semantics (Reduction)

The computation in lambda calculus is a sequence of applications of function to its arguments (which themselves are functions). Each step in the computation is a term rewriting that substitutes the arguments for bound variables in the abstraction's body.

**Definition 5 (Substitution)** For the arbitrary terms $t_1$, $t_2$ and variable $x$, $t_1[x \mapsto t_2]$ is defined as the result of substituting $t_2$ for each free occurrence of $x$ in $t_1$.

**Definition 6 (Redex and $\beta$-reduction)** A term matching the form

$$(\lambda\, x.t_1)\ t_2$$

is called *redex* (a reducible expression), and the operation of rewriting a redex,

$$(\lambda\, x.t_1)\ t_2 \rightsquigarrow t_1[x \mapsto t_2]$$

is called a *$\beta$-reduction*.

According to the definitions, simple $\beta$-reductions are shown as follows:

$$(\lambda x.x)\ y \ \rightsquigarrow \ x[x \mapsto y] \rightsquigarrow y$$

$$(\lambda x.x\ (\lambda x.x))\ (uv) \ \rightsquigarrow \ x\ (\lambda x.x)[x \mapsto uv] \rightsquigarrow uv\ (\lambda x.x)$$
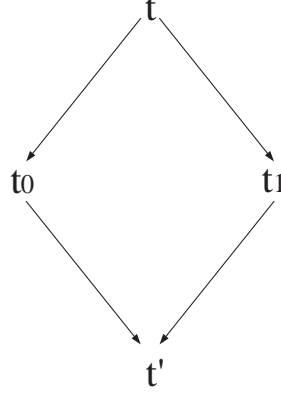
If $t'$ is obtained from $t$ by zero or more $\beta$-reductions (or just by renaming), then $t$ *reduces to* $t'$, and is written as $t \rightsquigarrow^* t'$. In lambda calculus, the reduction sequence is a type of execution, and the obvious computation termination is defined by the reduction sequence reaching a term that contains no redexes.

**Definition 7 (Normal Form)** A term $t$ is in a *normal form* if it does not have a redex as a subexpression. Moreover, if $t$ is reduced to the normal form $t'$, then $t'$ is called a normal form of term $t$.

The following are examples of reduction sequences, terminating in normal forms:

$$\underline{(\lambda x.y)(\lambda z.z)} \ \rightsquigarrow \ y$$

$$\underline{(\lambda x.x)(\lambda z.z)} \ \rightsquigarrow \ \lambda z.z$$

$$\underline{(\lambda x.xx)(\lambda z.z)} \ \rightsquigarrow \ \underline{(\lambda z.z)(\lambda z.z)} \rightsquigarrow (\lambda z.z)$$

$$\underline{(\lambda x.(\lambda y.yx)z)(t\ v)} \ \rightsquigarrow \ \underline{(\lambda x.zx)(t\ v)} \rightsquigarrow z\ (t\ v)$$

$$(\lambda x.\underline{(\lambda y.yx)z})(t\ v) \ \rightsquigarrow \ \underline{(\lambda y.y(t\ v))\ z} \rightsquigarrow z\ (t\ v)$$

Any redexes can be reduced at any time. At each step, some redexes (denoted by underlining) are chosen. The last two reduction sequences begin with the same expression and move in different directions (by reducing different redexes), but eventually return to the same normal form. The following theorem confirms that this is a general phenomenon for lambda calculus.

**Theorem 1 (Church-Rosser (Confluence) Property)** If $t \leadsto^* t_0$ and $t \leadsto^* t_1$, then there is a term $t'$ such that $t_0 \leadsto^* t'$ and $t_1 \leadsto^* t'$ [2].

This property states the uniqueness of the normal form, independent of the existence of normal form. As a special case, since normal forms can be reduced only by renaming, and if $t_0$ and $t_1$ are normal forms, then both must be $\alpha$-equivalent to $t'$, and therefore, to each other.

**Theorem 2** If a term $t$ exists a normal form $n$, then in $\alpha$-equivalence, the normal form, $n$, is unique.

Unfortunately, untyped lambda calculus terms may not have a normal form is false. In the following examples of $\beta$-reductions:

$$(\lambda\, x.xx)(\lambda\, x.xx) \quad \leadsto \quad (\lambda\, x.xx)(\lambda\, x.xx) \leadsto \ldots$$

$$(\lambda\, x.f(x\ x))(\lambda\, x.f(x\ x)) \quad \leadsto \quad f((\lambda\, x.f(x\ x))(\lambda\, x.f(x\ x))) \leadsto \ldots$$

the $\beta$-reduction sequences never terminate. For each initial term, there are no other reduction sequences, so that these initial terms have no normal form. Meanwhile, there are cases where both terminating and nonterminating

sequences begin with the same term by different redexes. For example,

$$\underline{(\lambda\,y.c)((\lambda\,x.xxx)(\lambda\,x.xxx))} \;\rightsquigarrow\; c$$

$$(\lambda\,y.c)\underline{((\lambda\,x.xxx)(\lambda\,x.xxx))} \;\rightsquigarrow\; (\lambda\,y.c)\underline{((\lambda\,x.xxx)(\lambda\,x.xxx))} \rightsquigarrow \ldots$$

For untyped lambda calculus, the existence of normal form is *undecidable* [7], which is immediate corollary of the *undecidability* of the *β-equivalence* of two lambda terms.

In summary, two important facts of untyped lambda calculus should be stated:

- any reduction sequence terminating in a normal form produces the same result. Computationally, it means that every expression describes an unique computation if it exists.

- not all reduction sequences end in a normal form even choosing the reduction strategy; the existence of normal form is undecidable.

So far, untyped lambda calculus, which is a solid core of functional programming languages, has been examined. Lambda calculus is known to be computationally equivalent in power to many other models for computation (including Turing machines); that is, any calculation that can be accomplished in any of these other models can be expressed in the lambda calculus, and vice versa [11]. But without type constraints, lambda calculus causes problems: the termination of the evaluation without type constraints is undecidable.

Due to the limitation of untyped lambda calculus, researchers have been investigating type systems for lambda calculus. In simply typed-lambda calculus, type constraints are imposed on terms to prevent self-application, and

to ensure its strong normalization property. But simple typed-lambda calculus has very little expressive power [7], and it is too weak to represent the XQ data types and operations. However, *System F*, a second-order lambda calculus, guarantees the termination of all computations, while retaining sufficient expressive power to encode XQ.

## 3.2   Introduction to System F

*System F* was first discovered by Girard [7], in 1971 as a system for constructive proofs in $PA_2$, and by Reynolds [11] in 1974 in the area of programming language. It is a extension of simple untyped lambda with treating types as variables and allowing new forms of *type abstractions* and *type applications*. Let us start with the ordinary untyped lambda calculus and add the capability to pass types as parameters. In untyped lambda calculus, the lambda *abstraction*, $(\lambda x.t)$, is used to abstract variables out of terms, and *application* is used to supply values for the abstracted parts without considering type constraints.

**Example 3.2.1** The abstraction $\lambda f. \lambda x.f(f\ x)$ [10], in untyped lambda calculus, denotes the function double, which accepts a function as an argument, and yields the composition of this function to itself. To extend this definition with type variable $T$, we rewrite it as:

$$\lambda f^{T \to T}. \lambda x^T.f(f\ x)$$

denoting the function double for type $T$. By abstracting on type variable $T$

with the symbol "Λ", a polymorphic function `double` is defined as

$$\Lambda T. \lambda f^{T \to T}. \lambda x^T. f(f\ x)$$

that can be applied to any type to obtain the corresponding "double" function for that type.

### 3.2.1 Syntax of System F

After extending lambda calculus with type variables, types become type expressions [7]:

- type variables $X$, $Y$, $Z$, . . . are types;

- if $U$ and $V$ are types, then function type $U \to V$ is a type, which denotes the set of functions that return a value of type $V$ when applied to a value of type $U$;

- if $V$ is a type and $X$ is a type variable, then the polymorphic type $\forall X.V$ is a type, which denotes the type of polymorphic function that, when applied to a type $X$, yields a result of type $V$

Therefore, there are five ways to form *terms*:

- a *variable* $x$ of type $T$ is a term;

- an *abstraction* $\lambda x^U.v$ is a term of type $U \to V$, where $x$ is a variable of type $U$ and $v$ is of type $V$;

- an *application* $t\ u$ is a term of type $V$, where $t$ is of type $U \to V$, and $u$ is type of $U$;

$$
\begin{aligned}
\textit{Type} ::= \; & X && \textit{type variable} \\
\mid \; & U \to V && \textit{function type} \\
\mid \; & \forall\, X.V && \textit{polymorphic type} \\
\textit{Term} ::= \; & x && \textit{variable} \\
\mid \; & \lambda\, x^{U}.v && \textit{abstraction} \\
\mid \; & t \; u && \textit{application} \\
\mid \; & \Lambda X.v && \textit{type abstraction} \\
\mid \; & t\,[U] && \textit{type application}
\end{aligned}
$$

Figure 3.1: Syntax definition of *System F*

- a *type abstraction* $\Lambda X.v$ is a term of type $\forall\, X.V$, where $v$ is a term of type $V$;

- a *type application* $t\,[U]$ is a term of type $V[X \mapsto U]$, where $t$ is a term of type $\forall\, X.V$, and $U$ is a type

The syntax definition of *System F* is presented in Figure 3.1. As well as the usual conversion, the $\to$ associates to the right, where

$$T_1 \to T_2 \to \cdots \to \; T_n \to \alpha \quad \text{stands for} ( T_1 \to ( T_2 \to \cdots \to \; ( T_n \to \alpha)))$$

and application associates to the left, where

$$f \; t_1 \; t_2 \dots t_n \quad \text{stands for} \quad (\dots ((f \; t_1) \; t_2) \dots) \; t_n$$

The $\beta$ reduction "$\rightsquigarrow$" is extended to denote a new rule regarding the type applications.

**Definition 8 ($\beta$-reduction)**

$$(\lambda\, x.t_1)\ t_2 \rightsquigarrow t_1[x \mapsto t_2] \tag{3.1}$$

$$(\Lambda X.t)[U] \rightsquigarrow t[X \mapsto U] \tag{3.2}$$

where $t_1$, $t_2$ are terms, $x$ is a term variable, and $X$ is a type variable.

During an evaluation, the first rule denotes the usual *$\beta$-reduction*; that is, when a function, in the abstraction form of $\lambda\, x.t_1$, is applied to term $t_2$, the pair forms a *redex*, and term $t_2$ will substitute for all the occurrences of the bound variable $x$ in $t_1$. The second rule states that, when a function, which is a type abstraction, $\Lambda X.t$, is applied to a type value $U$, the pair is also a *redex* and type $U$ will substitute for the occurrences of type variable $X$ in term $t$.

**Example 3.2.2** The double function $\Lambda T.\, \lambda f^{T \to T}.\, \lambda x^T.f\ (f\ x)$

for type `int` becomes

$$(\Lambda T.\, \lambda f^{T \to T}.\, \lambda x^T.f(f\ x))\ [\texttt{int}]$$
$$\rightsquigarrow\ (\lambda f^{T \to T}.\, \lambda x^T.f(f\ x))[T \mapsto \texttt{int}]$$
$$\rightsquigarrow\ \lambda f^{\texttt{int} \to \texttt{int}}.\, \lambda x^{\texttt{int}}.f(f\ x)$$

and for the complex function type `real` $\to$ `real`, becomes

$$(\Lambda T.\, \lambda f^{T \to T}.\, \lambda x^T.f(f\ x))\ [\texttt{real} \to \texttt{real}]$$
$$\rightsquigarrow\ (\lambda f^{T \to T}.\, \lambda x^T.f(f\ x))[T \mapsto (\texttt{real} \to \texttt{real})]$$
$$\rightsquigarrow\ \lambda f^{(\texttt{real} \to \texttt{real}) \to (\texttt{real} \to \texttt{real})}.\, \lambda x^{\texttt{real} \to \texttt{real}}.f(f\ x)$$

### 3.2.2 Why System F?

It was mentioned in Section 3.1.2 that in lambda calculus, untyped terms can not guarantee to have a normal form; that is, the termination of computation is undecidable. But when we considering types in the lambda calculus, the things have changed. It is the case that the typed-terms always possess normal forms [3]. This property also holds for *System F* [7].

**Theorem 3 (Strong Normalization Theorem)** All the well-typed terms of *System F* have normal forms. Moreover, all the terms are strongly normalizable, and there are no infinite reduction sequences, regardless of the evaluation strategy.

For languages in which each expression describes a terminating computation, such as simple typed lambda calculus, there must be computable functions that cannot be expressed [3]. Indeed we are used to taking this fact as evidence that such language are uninteresting for practical computation. Yet *System F* is such a language, in which one can express "almost everything" that one might actually want to compute. Indeed, Girard has shown that *System F* provides a surprising degree of expressiveness for the computations for a variety of data types [7].

## 3.3 Expressive Power of System F

In this section, it is demonstrated that a variety of primitive data types such as boolean, number, and list can be encoded in *System F*. These encodings are interesting for two reasons. First, they exemplify type abstraction and

application. Secondly, they demonstrate that *System F* is, computationally, a very rich language, in the sense that the pure system can express a large range of data and control structures. This means that, when we later design a full-scale language with *System F* as its core, we can add these features as primitives (for efficiency, and so that we can equip them with a more convenient concrete syntax) without disturbing the fundamental properties of the core language.

### 3.3.1   Inductive Types

The expression power of System F allows the representation of basic types such as integers, lists, and trees by defining the inductive types.

**Definition 9 (Inductive Type)** An *inductive type* $\Theta$ is a data type that can be constructed by a set of constructors $f_1$, $f_2$, ..., $f_k$, with corresponding types:

$$f_1 \; : T_1 \qquad \text{with } T_1 = T_{11} \to \cdots \to T_{1n_1} \to \Theta$$
$$f_2 \; : T_2 \qquad \text{with } T_2 = T_{21} \to \cdots \to T_{2n_2} \to \Theta$$
$$\vdots$$
$$f_k \; : T_k \qquad \text{with } T_k = T_{k1} \to \cdots \to T_{1n_k} \to \Theta$$

where $\Theta$ *occurs positively* in any $T_{ij}$, and each term of $\Theta$ can be represented by constructor $f_i$ uniquely [7].

**Definition 10 (Representation of the Inductive Type)** An inductive type, $\Theta$, which is described by the constructors $f_1$, $f_2$, ..., $f_k$ in Definition 9, is rep-

resented by

$$
\begin{aligned}
\Theta \;=\; & \forall X.\, (T_{11} \to \cdots \to T_{1n_1} \to \Theta)[\Theta \mapsto X] \\
& \to (T_{21} \to \cdots \to T_{2n_2} \to \Theta)[\Theta \mapsto X] \\
& \quad \vdots \\
& \to (T_{k1} \to \cdots \to T_{kn_k} \to \Theta)[\Theta \mapsto X] \\
& \to X
\end{aligned}
$$

where $T_i[\Theta \mapsto X]$ denotes $T_i$, where all the occurrences of $\Theta$ are replaced by a type variable $X$.

The most interesting part of *System F* is its expressive power to define the common types. Now that the general schema of the representation of type and the corresponding iterator function have been introduced, some examples are provided.

## 3.3.2 Representation of Simple Types

In this section, the general schema, defined in Section 3.3.1, is applied to describe the common types such as boolean, product type, natural number, and list.

**Example 3.3.1 (Boolean)** The boolean type, Bool, has two functions: `True` and `False`, 0-ary functions of type Bool. According the Definition 9, the boolean type is defined as

$$
\mathsf{Bool} \stackrel{\text{def}}{=} \forall X.X \to X \to X
$$

and the two 0-ary constructors are

$$\mathsf{True} \stackrel{\text{def}}{=} \Lambda X . \lambda x^X . \lambda y^X . x$$

$$\mathsf{False} \stackrel{\text{def}}{=} \Lambda X . \lambda x^X . \lambda y^X . y$$

For $e$, $e'$ of type $U$, and $b$ of type $\mathsf{Bool}$, the condition expression can be defined as an iterator function:

$$\text{if } b \text{ then } e \text{ else } e' \stackrel{\text{def}}{=} b \, [U] \, e \, e'$$

and the following reductions prescribe the right behaviour:

$$
\begin{aligned}
\text{if } \underline{\mathsf{True}} \text{ then } e \text{ else } e' \quad &\stackrel{\text{def}}{=} \quad (\Lambda X . \lambda x^X . \lambda y^X . x)[U] \, e \, e' \\
&\rightsquigarrow \quad (\lambda x^U . \lambda y^U . x) \, e \, e' \\
&\rightsquigarrow \quad (\lambda y^U . e) \, e' \\
&\rightsquigarrow \quad e \\
\text{if } \underline{\mathsf{False}} \text{ then } e \text{ else } e' \quad &\stackrel{\text{def}}{=} \quad (\Lambda X . \lambda x^X . \lambda y^X . y)[U] \, e \, e' \\
&\rightsquigarrow \quad (\lambda x^U . \lambda y^U . y) \, e \, e' \\
&\rightsquigarrow \quad (\lambda y^U . y) \, e' \\
&\rightsquigarrow \quad e'
\end{aligned}
$$

Moreover,

$$
\begin{aligned}
\mathsf{not} \quad &\stackrel{\text{def}}{=} \quad \lambda a^{\mathsf{Bool}} . \Lambda X . \lambda x^X . \lambda y^X . b \, [X] \, y \, x \\
\mathsf{and} \quad &\stackrel{\text{def}}{=} \quad \lambda a^{\mathsf{Bool}} . \lambda b^{\mathsf{Bool}} . \Lambda X . \lambda x^X . \lambda y^X . a \, [X] \, (b \, [X] \, x \, y) y
\end{aligned}
$$

where the boolean function $\mathsf{not}$ has the type $\mathsf{Bool} \to \mathsf{Bool}$, and the function $\mathsf{and}$ has type $\mathsf{Bool} \to \mathsf{Bool} \to \mathsf{Bool}$.

**Example 3.3.2 (Integer)** The integer type, Int, has two constructors: one is the zero function, zero, with type Int; the other is the successor function succ of type Int → Int. Therefore, type Int is defined as

$$\text{Int} \stackrel{\text{def}}{=} \forall\, X.X \rightarrow (X \rightarrow X) \rightarrow X$$

and we have the integer numbers as

$$0 \stackrel{\text{def}}{=} \Lambda X. \lambda\, x^X. \lambda f^{X\rightarrow X}.x$$

$$1 \stackrel{\text{def}}{=} \Lambda X. \lambda\, x^X. \lambda f^{X\rightarrow X}.f\ x$$

$$2 \stackrel{\text{def}}{=} \Lambda X. \lambda\, x^X. \lambda f^{X\rightarrow X}.f(f\ x)$$

$$\vdots$$

$$n \stackrel{\text{def}}{=} \Lambda X. \lambda\, x^X. \lambda f^{X\rightarrow X}.f^n x$$

where $f^n x$ denotes $f(f(f \ldots (f\ x)))$ with $n$ occurrences of $f$. The two basic constructors are defined as

$$\texttt{zero} \stackrel{\text{def}}{=} \Lambda X. \lambda\, x^X. \lambda f^{X\rightarrow X}.x$$

$$\texttt{succ } t \stackrel{\text{def}}{=} \Lambda X. \lambda\, x^X. \lambda f^{X\rightarrow X}.f(t\ [X]\ x\ f)$$

Clearly, $\texttt{zero} = 0$ and $\texttt{succ } n = n + 1$.

Also, there is an iterator function, $\texttt{It}_{\texttt{int}}$, with type Int → $U$,

$$\texttt{It}_{\texttt{Int}}\ t = t\ [U]\ g_1\ g_2$$

where $g_1 : U$ and $g_2 : \mathsf{Int} \to U^2$, and now

$$\mathtt{It_{Int}\ zero} \overset{\text{def}}{=} \mathtt{zero}\ [U]\ g_1\ g_2$$

$$\overset{\text{def}}{=} (\Lambda X.\lambda\, x^X.\lambda f^{X \to X}.x)[U]\ g_1\ g_2$$

$$\rightsquigarrow (\lambda\, x^U.\lambda f^{U \to U}.x)\ g_1\ g_2$$

$$\rightsquigarrow (\lambda f^{U \to U}.g_1)\ g_2$$

$$\rightsquigarrow g_1$$

$$\mathtt{It_{Int}\ succ}\ t \overset{\text{def}}{=} (\mathtt{succ}\, t)\ [U]\ g_1\ g_2$$

$$\overset{\text{def}}{=} (\Lambda X.\lambda\, x^X.\lambda f^{X \to X}.f(t\ [X]\ x\ f))[U]\ g_1\ g_2$$

$$\rightsquigarrow (\lambda\, x^U.\lambda f^{U \to U}.f\ (t\ [U]\ x\ f))\ g_1\ g_2$$

$$\rightsquigarrow (\lambda f^{U \to U}.f\ (t\ [U]\ g_1\ f))\ g_2$$

$$\rightsquigarrow g_2(t\ [U]\ g_1\ g_2)$$

$$= g_2(\mathtt{It_{Int}}\ t)$$

Based on the iterator, $\mathtt{It_{Int}}$, in order to define a function $h : \mathsf{Int} \to U$, it is necessary to decide the corresponding $g_1 : U$ and $g_2 : \mathsf{Int} \to U$, considering the results, when function $h$ is inductively applied to the constructor of type $\mathsf{Int}$:

$$h\ 0 = g_1 \tag{3.3}$$

$$h\ (\mathtt{succ}\ n) = g_2(h\ n) \tag{3.4}$$

Then $h$ is defined by using $g_1$ and $g_2$:

$$h\ n \overset{\text{def}}{=} n\ [U]g_1\ g_2 \qquad \text{that is} \qquad h \overset{\text{def}}{=} \lambda\, n^{\mathsf{Int}}.n\ [U]\ g_1\ g_2$$

---

[2] ":" means "has the type of"

The *addition* function, add : Int $\rightarrow$ Int $\rightarrow$ Int, which adds two integer numbers, with the case:

$$\text{add } m \text{ zero} = m$$

$$\text{add } m \text{ (succ } n) = \text{succ (add } m \text{ } n)$$

and matches the format of Equation 3.3 and 3.4, is defined as:

$$\text{add } m \text{ } n \overset{\text{def}}{=} n \text{ [Int] } m \text{ succ}$$

or more abstractly by

$$\text{add} \overset{\text{def}}{=} \lambda m^{\text{Int}}. \lambda n^{\text{Int}}.n \text{ [Int] } m \text{ succ}$$

**Example 3.3.3 (Product Type($\times$))** The expressive power of *System F* not only allows us to represent the primitive type, but also provides the capability to describe the various type constructors. In this thesis, the *product type*, $U \times V$, is demonstrated, allowing the construction of the pair $(a, b)$, where $a$ has the type $U$ and $b$ has the type $V$. The only constructor function for the product type is the pair function,

$$f : U \rightarrow V \rightarrow U \times V$$

which takes two arguments of type $U$ and $V$, and returns a pair of type $U \times V$. Therefore, by Definition 9, a *product type* is signified as

$$U \times V \overset{\text{def}}{=} \forall X.(U \rightarrow V \rightarrow X) \rightarrow X$$

and, given $a : U$, $b : V$, the pair $(a, b)$ is defined as

$$(a, b) \overset{\text{def}}{=} \Lambda X. \lambda f^{U \rightarrow V \rightarrow X}.f \text{ } a \text{ } b$$

Given a pair $(a, b)$, the project functions are defined as

$$(a, b).0 \overset{\text{def}}{=} (a, b) [U] (\lambda x^U . \lambda y^V . x)$$

$$(a, b).1 \overset{\text{def}}{=} (a, b) [V] (\lambda x^U . \lambda y^V . y)$$

The project function $(a, b).0$ is calculated by the $\beta$-reduction,

$$
\begin{aligned}
(a, b).0 \overset{\text{def}}{=} & \ (\Lambda X . \lambda f^{U \to V \to X} . f \ a \ b) [U] (\lambda x^U . \lambda y^V . x) \\
\rightsquigarrow & \ (\lambda f^{U \to V \to U} . f \ a \ b) (\lambda x^U . \lambda y^V . x) \\
\rightsquigarrow & \ (\lambda x^U . \lambda y^V . x) \ a \ b \\
\rightsquigarrow & \ (\lambda y^V . a) \ b \\
\rightsquigarrow & \ a
\end{aligned}
$$

which obtains the first element of the pair $(a, b)$, and $(a, b).1$ returns the second element, since

$$
\begin{aligned}
(a, b).1 \overset{\text{def}}{=} & \ (\Lambda X . \lambda f^{U \to V \to X} . f \ a \ b) [V] (\lambda x^U . \lambda y^V . y) \\
\rightsquigarrow & \ (\lambda f^{U \to V \to V} . f \ a \ b) (\lambda x^U . \lambda y^V . y) \\
\rightsquigarrow & \ (\lambda x^U . \lambda y^V . y) \ a \ b \\
\rightsquigarrow & \ (\lambda y^V . y) \ b \\
\rightsquigarrow & \ b
\end{aligned}
$$

**Example 3.3.4 (List)** A list is a finite sequence of elements. Typical lists are $[2, 5, 6]$ and ["good", "better", "best"]. The empty list, [ ], has no elements. For a given list, the order of elements is significant, and elements may appear more than once. For instance, the following lists are different:

$$[3, 5, 9] \qquad [3, 5, 9, 9] \qquad [3, 9, 5]$$

Each element in a list must have the same type. If the element $u_0$, $u_1$, ..., $u_n$ are elements of type $U$, the list $[u_0, u_1, \ldots, u_n]$ has a type $U$ list. The empty list, $[\,]$, has a polymorphic type $\alpha$ list, which is regarded as a list with any type of elements.

A list of type $U$ list is constructed by the following two functions:

- `nil` : $U$ list, a synonym of the empty list, $[\,]$;

- construct function, `cons` : $U \rightarrow U$ list $\rightarrow U$ list, returns a new list by inserting an element of type $U$ in front of an existing list of type $U$ list

Therefore, each list is either `nil`, if empty, or has the form `cons` $u$ $l$, where $u$ is the *head* of type $U$ and $l$ is the *tail* of type $U$ list. For the given elements $u_0$, $u_1$, ..., $u_n$, the list $[u_0, u_1, \ldots, u_n]$ is constructed as follows:

$$
\begin{aligned}
\texttt{nil} &= [\,] \\
\texttt{cons } u_n \texttt{ nil} &= [u_n] \\
\texttt{cons } u_{n-1} [u_n] &= [u_{n-1}, u_n] \\
&\vdots \\
\texttt{cons } u_0 [u_1, \ldots, u_n] &= [u_0, u_1, \ldots, u_{n-1}, u_n]
\end{aligned}
$$

With the schema of Definition 9, the type list $U$ is defined as:

$$U \text{ list} \stackrel{\text{def}}{=} \forall X.X \rightarrow (U \rightarrow X \rightarrow X) \rightarrow U$$

and the individual list $[u_0, u_1, \ldots, u_n]$ is

$$
\begin{aligned}
\texttt{cons } u_0 \ (\texttt{cons } u_1 \ \ldots \ (\texttt{cons } u_n \ \texttt{nil}) \ldots) &\stackrel{\text{def}}{=} \\
\Lambda X. \lambda x^X. \lambda f^{U \rightarrow X \rightarrow X}. f \ u_0 \ (f \ u_1 \ldots (f \ u_n \ x))
\end{aligned}
$$

In addition, two constructor functions are denoted as:

$$\texttt{nil} \stackrel{\text{def}}{=} \Lambda X . \lambda\, x^U . \lambda f^{U \to X \to X} . x$$

$$\texttt{cons}\ t\ ts \stackrel{\text{def}}{=} \Lambda X . \lambda\, x^U . \lambda f^{U \to X \to X} . f\ t\ (ts\ [X]\ x\ f)$$

where $t$ is an element of type $U$, and $ts$ is a list of type $\mathsf{list}\ U$.

**Definition 11 (Iterator on List)** The iterator function on lists $\texttt{It}_{\texttt{list}}$ : $U\ \mathsf{list} \to W$ takes an argument of type $U\ \mathsf{list}$ and returns the result of type $W$, and is denoted as

$$\texttt{It}_{\texttt{list}}\ ts \stackrel{\text{def}}{=} ts\ [W]\ g_1\ g_2 \tag{3.5}$$

where $g_1 : W$ and $g_2 : U \to W \to W$, and $\texttt{It}_{\texttt{list}}$ satisfies,

$$
\begin{aligned}
\texttt{It}_{\texttt{list}}\ \texttt{nil}\ &\stackrel{\text{def}}{=}\ \texttt{nil}[W]\ g_1\ g_2 \\
&\stackrel{\text{def}}{=}\ (\Lambda X . \lambda\, x^U . \lambda f^{U \to X \to X} . x)[W] g_1\ g_2 \\
&\rightsquigarrow\ (\lambda\, x^U . \lambda f^{U \to W \to W} . x)\ g_1\ g_2 \\
&\rightsquigarrow\ (\lambda f^{U \to W \to W} . g_1)\ g_2 \\
&\rightsquigarrow\ g_1
\end{aligned}
$$

and

$$
\begin{aligned}
\texttt{It}_{\texttt{list}}\ (\texttt{cons}\ t\ ts)\ &\stackrel{\text{def}}{=}\ (\texttt{cons}\ t\ ts)\ [W]\ g_1\ g_2 \\
&\stackrel{\text{def}}{=}\ (\Lambda X . \lambda\, x^U . \lambda f^{U \to X \to X} . f\ t\ (ts\ [X]\ x\ f))\ [W] g_1\ g_2 \\
&\rightsquigarrow\ (\lambda\, x^U . \lambda f^{U \to W \to W} . f\ t\ (ts\ [W]\ x\ f))\ g_1\ g_2 \\
&\rightsquigarrow\ (\lambda f^{U \to W \to W} . f\ t\ (ts\ [W]\ g_1\ f))\ g_2 \\
&\rightsquigarrow\ g_2\ t\ (ts\ [W]\ g_1\ g_2) \\
&=\ g_2\ t\ (\texttt{It}_{\texttt{list}}\ ts)
\end{aligned}
$$

The iterator definition provides a schema to define a function $f : U \text{ list} \to W$ by determining the corresponding $g_1 : W$ and $g_2 : U \text{ list} \to W \to W$, which satisfy

$$f \text{ nil} = g_1 \quad \text{and} \quad f \text{ (cons } t \text{ } ts) = g_2 \text{ } t \text{ } (f \text{ } ts) \tag{3.6}$$

such that function $f$ is defined by Definition 11 as

$$f \text{ } ts \stackrel{\text{def}}{=} ts \text{ } [W] \text{ } g_1 \text{ } g_2 \quad \text{that is,} \quad f \stackrel{\text{def}}{=} \lambda \text{ } ts^{U \text{ list}}.ts \text{ } [W] \text{ } g_1 \text{ } g_2 \tag{3.7}$$

The *length* function $\texttt{length} : U \text{ list} \to \text{Int}$ calculates the number of elements in a given list, and it is known that

$$\texttt{length nil} = 0$$

$$\texttt{length (cons } t \text{ } ts) = 1 + (\texttt{length } ts) = (\lambda \text{ } t^U . \lambda \text{ } c^{\text{Int}}.1 + c) \text{ } t \text{ } (\texttt{length } ts)$$

that is, $g_1 = 0$ with type $\text{Int}$ and $g_2 = \lambda \text{ } t^U . \lambda \text{ } c^{\text{Int}}.1 + c$ with type $U \to \text{Int} \to \text{Int}$.

Therefore, the function $\texttt{length}$ is defined using Equation 3.7 such that,

$$\texttt{length} \stackrel{\text{def}}{=} \lambda \text{ } ts^{U \text{ list}}.ts \text{ } [\text{Int}] \text{ } 0 \text{ } (\lambda \text{ } t^U . \lambda \text{ } c^{\text{Int}}.1 + c)$$

where type variable $W$ is replaced by $\text{Int}$.

Now that the expressive power of *System F* has been illustrated. In fact, all the functions, which are *System F* representable, have the common property:

**Proposition 2** The functions representable in *System F* are exactly those which are *provably total* in $\text{PA}_2$[3].

---

[3]A detained definition and proof can be found in [7]

In this chapter, the expressive power of *System F* has been demonstrated. The definition schema allows the representation of XQ and its basic operations in *System F*, as detailed in the following chapter.

# Chapter 4

# Encoding XQ in *System F*

Now that the expressive power of *System F* has been illustrated by simple data types, *System F* will presented as the core of XQ, to improve the generality and power of XQ. In this chapter, *System F* is adopted to represent the XML data model that has been defined as an ordered forest whose elements are rooted labelled trees [6]. By the representation of XML tree and forest data types in *System F*, all the basic operators of XQ are encoded by iterators defined on the XML tree and forest structures. Also a general XQ syntax is proposed in*System F* format, and the semantics are described by defining an *Environment* data type in *System F* with the corresponding operators.

## 4.1   XML data in *System F*

The dynamic interval method presents an XML document as an ordered forest whose elements are rooted labelled trees, shown in Definition 1. Without a tree definition, the XML forest constructors are not primitive, which means

that the empty node([ ]), single node forest (`<s> XF </s>`), and concatena-
tion constructor (`@`) are not independent. For example, a single node forest
can be constructed by using an empty forest and a tree; the concatenation of
two forests can also be represented recursively in a more primitive way with
a new constructor **cons**, which will be described in Section 4.2.2. In order to
precisely describe XML document, XML forest and XML tree are considered
separately.

**Definition 12 (XML Tree)**

$$XTree \stackrel{\text{def}}{=} \texttt{Xnode } String \ XForest$$

where an XML tree is constructed a *String s* as its root and an XML forest
*ts* as its children.

In this chapter, `Xnode` *s ts* is used, instead of `<s>`*ts* `</s>`, to express an XML
tree whose root is string *s*. An XML forest is described as a list of ordered
XML trees.

**Definition 13 (XML Forest)**

$$XForest \stackrel{\text{def}}{=} \texttt{nil } \mid \ \texttt{Cons } XTree \ XForest$$

where `nil` denotes an empty forest without elements, and `Cons` signifies the
construction of an XML forest by putting an XML tree *t* in front of an
existing XML forest *ts*, written as *ts'* = `Cons` *t ts*.

Simultaneous inductions occur in the *XTree* and *XForest* definitions: build-
ing an object of one type involves the object of another. For example, con-
structor `Cons` builds tree *t* by employing a forest, *ts*; `Xnode` takes tree *t* and

forest *ts* to build a new forest, *ts'*. In order to solve this problem, a new

new type, *Xtf*, which is a combination of type *XTree* and type *XForest*, con-

taining constructors of both types, to solve the simultaneous construction

inherently.

To describe type *Xtf*, a general type, called Sum type, is introduced.

**Definition 14 (Sum Type)** If $U$ and $V$ are types, the Sum of $U$ and $V$,

denoted as $U + V$, is defined as

$$U + V \stackrel{\text{def}}{=} (U \to \alpha) \to (V \to \alpha) \to \alpha^1$$

Sum type denotes an *alternative* construction, where different values are

paired with *tags*, with which the values can be processed by branch. For

example, $\iota^1 u$ is a term of $u$ of type $U$ paired with tag 1. If $t$ has type $U$ and

$f_1$, $f_2$ are functions, then $(f_1, f_2)\, \iota^1 t$ results in applying function $f_1$ to $t$.

For $u : U$, $v : V$, two constructors for the Sum type are defined as follows:

$$\iota^1 u \stackrel{\text{def}}{=} \lambda x^{U \to \alpha}.\ \lambda y^{V \to \alpha}.\ x\ u$$

$$\iota^2 v \stackrel{\text{def}}{=} \lambda x^{U \to \alpha}.\ \lambda y^{V \to \alpha}.\ y\ v$$

For functions $f_1 : U \to R$, $f_2 : V \to S$, and $t$ of type $U + V$,

$$(f_1, f_2)\ t \stackrel{\text{def}}{=} t[R + S]\ f_1\ f_2$$

and for terms $u : U$ and $v : V$, $(f_1, f_2)\, t$ satisfies

$$(f_1, f_2)\ \iota^1\ u = f_1\ u \tag{4.1}$$

$$(f_1, f_2)\ \iota^2\ v = f_2\ v \tag{4.2}$$

---

[1]The primitive types $\alpha$ and $\beta$ (denoted by Greek letter) are employed to implicitly express the type polymorphism instead of using the *universal type* expression $\forall \alpha.T$, which binds the occurrences of $\alpha$ in $T$. For example, $(U \to \alpha) \to (V \to \alpha) \to \alpha \equiv \forall \alpha.(U \to \alpha) \to (V \to \alpha) \to \alpha$.

where

$$(f_1, f_2) \; \iota^1 \; u \;\; \overset{\text{def}}{=} \;\; (\lambda \, x^{U \to \alpha}. \; \lambda \, y^{V \to \alpha}. \; x \; u)[R + S] \; f_1 \; f_2$$

$$\rightsquigarrow \;\; (\lambda \, x^{U \to R+S}. \; \lambda \, y^{V \to R+S}. \; x \; u) \; f_1 \; f_2$$

$$\rightsquigarrow \;\; (\lambda \, y^{V \to R+S}. \; f_1 \; u) \; f_2$$

$$\rightsquigarrow \;\; f_1 \; u$$

and similarly, $(f_1, f_2) \; \iota^2 \; v \rightsquigarrow^* f_2 \; v$

The problem of the simultaneous induction of types *XTree* and *XForest* can be solved by using the concept of Sum. A new type *Xtf*, which is a Sum type that combines types *XTree* and *XForest*, is built. When a function is recursively applied to an argument, its action depends on the element type it visits. For example, when the boolean operation `equal` compares the equality of two trees, at the top level, `equal` is a tree operator, comparing two root labels. Then `equal` becomes a forest operator recursively applied to the root's children. For the elements of the children forests, `equal` behaves as a tree operator, and so on.

**Definition 15 (*Xtf* Type)** *Xtf* is defined as an Sum of *XTree* and *XForest*, which is represented as follows:

$$Xtf \overset{\text{def}}{=} XTree + XForest$$

and as a polymorphic type in *System F* denoted as:

$$Xtf \overset{\text{def}}{=} \alpha \to (String \to \beta \to \alpha) \to \beta \to (\alpha \to \beta \to \beta) \to (\alpha + \beta)$$

The new type *Xtf* has three constructors, which are inherited from *XTree* and *XForest*:

- the function, `Xnode`: $String \rightarrow \beta \rightarrow \alpha$, builds a new tree with *String s* and forest *ts* in `XNode` *s ts*;

- the empty forest, `nil` : $\beta$

- the function, `Cons`: $\alpha \rightarrow \beta \rightarrow \beta$, constructs a forest by inserting tree *t* before forest *ts* with `Cons` *t ts*

With the general scheme of *System F*, the constructors of type *Xtf* are treated as certain *abstractions*.

**Definition 16 (Constructors for Type *Xtf*)**

$$\text{Xnode } s \ ts \quad \overset{\text{def}}{=} \quad \lambda\, x_1^{String \rightarrow \beta \rightarrow \alpha} \ . \ \lambda\, y_1^{\beta} \ . \ \lambda\, y_2^{\alpha \rightarrow \beta \rightarrow \beta}.$$

$$x_1 \ s \ (ts \ [\beta] \ x_1 \ y_1 \ y_2) \tag{4.3}$$

$$\text{nil} \quad \overset{\text{def}}{=} \quad \lambda\, x_1^{String \rightarrow \beta \rightarrow \alpha} \ . \ \lambda\, y_1^{\beta} \ . \ \lambda\, y_2^{\alpha \rightarrow \beta \rightarrow \beta}. \quad y_1 \tag{4.4}$$

$$\text{Cons } t \ ts \quad \overset{\text{def}}{=} \quad \lambda\, x_1^{String \rightarrow \beta \rightarrow \alpha} \ . \ \lambda\, y_1^{\beta} \ . \ \lambda\, y_2^{\alpha \rightarrow \beta \rightarrow \beta} \ .$$

$$y_2 \ (t \ [\alpha] \ x_1 \ y_1 \ y_2)(ts \ [\beta] \ x_1 \ y_1 \ y_2) \tag{4.5}$$

where (4.3) is a constructor from *XTree*, (4.4) and (4.4) are inherited from *XForest*.

With type *Xtf* and its constructors, an individual XML tree and forest can be represented by a function in *abstraction* form, shown in Example 4.1.1.

**Example 4.1.1** An XML tree in Fig 4.1 is represented as

$$t \quad \overset{\text{def}}{=} \quad \lambda\, x_1^{String \rightarrow \beta \rightarrow \alpha} \ . \ \lambda\, y_1^{\beta} \ . \ \lambda\, y_2^{\alpha \rightarrow \beta \rightarrow \beta}.$$

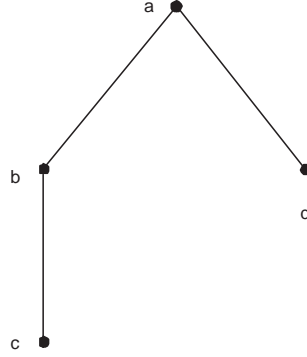$$x_1 \ a \ (y_2 \ (x_1 \ b \ (y_2 \ (x_1 \ c \ y_1) \ y_1)) \ (y_2 \ (x_1 \ d \ y_1) \ y_1))$$

Figure 4.1: An XML tree

we recognize that, with replacing $x_1$ by `Xnode`, $y_2$ by `Cons`, and $y_1$ by `nil`, $t$ becomes

$$\texttt{Xnode}\ a\ (\texttt{Cons}\ (\texttt{Xnode}\ b\ (\texttt{Cons}\ (\texttt{Xnode}\ c\ \texttt{nil})\ \texttt{nil}))$$

$$(\texttt{Cons}\ (\texttt{Xnode}\ d\ \texttt{nil})\ \texttt{nil}))$$

which is obtained by reducing $t\ [XTree]$ `Xnode nil Cons`.

As an Sum type, *Xtf* combines the properties of both *XTree* and *XForest*. Therefore, *Xtf*'s induction function, *iterator* `It`, can also manipulates both *Xtree* and *XForest* data.

**Definition 17 (Iterator for type *Xtf*)** Given object $v$ of type $V$, functions $f_T : String \rightarrow V \rightarrow U$ and $f_{XF} : U \rightarrow V \rightarrow V$, with the definition $Xtf = XTree + XFroest$, for an object $m$ of type *Xtf*, its iterator, `It`, is defined as

$$\texttt{It}\ m \overset{\text{def}}{=} m\ [U + V]\ f_T\ v\ f_{XF}$$

where the resulting type of It depends on its arguments. For object $t$ of type *XTree*,

$$\text{It } t^{\,2} \stackrel{\text{def}}{=} t\,[U]\,f_T\,v\,f_{XF}$$

 with the result type $U$. For object $ts$ of type *XForest*,

$$\text{It } ts \stackrel{\text{def}}{=} ts\,[V]\,f_T\,v\,f_{XF}$$

resulting in type $V$.

**Proposition 3** Given an object $v : V$, functions $f_T : String \rightarrow V \rightarrow U$, and $f_{XF} : U \rightarrow V \rightarrow V$, the iterator It satisfies

$$\text{It }(\text{Xnode } s\ ts)\ =\ f_T\ s\ (\text{It } ts) \tag{4.6}$$

$$\text{It nil}\ =\ v \tag{4.7}$$

$$\text{It }(\text{Cons } t\ ts)\ =\ f_{XF}\ (\text{It } t)\ (\text{It } ts) \tag{4.8}$$

**Proof:**  From the definition of iterator It, for an object $m$ of type *Xtf*, It is expressed as:

$$\text{It } m \stackrel{\text{def}}{=} m\,[U + V]\,f_T\,v\,f_{XF}$$

---

[2]Here, the sum symbol, $\iota$, is omitted for convenience.

since

$$\text{It } (\text{Xnode } s \text{ } ts) \overset{\text{def}}{=} (\text{Xnode } s \text{ } ts) \text{ } [U] \text{ } f_T \text{ } v \text{ } f_{XF}$$

$$\rightsquigarrow (\lambda \text{ } x_1^{String \to \beta \to \alpha}. \text{ } \lambda \text{ } y_1^{\beta} \text{ } . \text{ } \lambda \text{ } y_2^{\alpha \to \beta \to \beta} \text{ } . \text{ } x_1$$
$$s \text{ } (ts \text{ } [\beta] \text{ } x_1 \text{ } y_1 \text{ } y_2)) \text{ } [U] \text{ } f_T \text{ } v \text{ } f_{XF}$$

$$\rightsquigarrow (\lambda \text{ } x_1^{String \to \beta \to U}. \text{ } \lambda \text{ } y_1^{\beta} \text{ } . \text{ } \lambda \text{ } y_2^{U \to \beta \to \beta} \text{ } . \text{ } x_1$$
$$s \text{ } (ts \text{ } [\beta] \text{ } x_1 \text{ } y_1 \text{ } y_2)) \text{ } f_T \text{ } v \text{ } f_{XF}$$

$$\rightsquigarrow (\lambda \text{ } y_1^{\beta}. \text{ } \lambda \text{ } y_2^{U \to \beta \to \beta} \text{ } . \text{ } f_T \text{ } s \text{ } (ts \text{ } [\beta] \text{ } f_T \text{ } y_1 \text{ } y_2)) \text{ } v \text{ } f_{XF}$$

$$\rightsquigarrow (\lambda \text{ } y_2^{U \to V \to V}. \text{ } f_T \text{ } s \text{ } (ts \text{ } [V] \text{ } f_T \text{ } v \text{ } y_2)) \text{ } f_{XF}$$

$$\rightsquigarrow f_T \text{ } s \text{ } (ts \text{ } [V] \text{ } f_T \text{ } v \text{ } f_{XF})$$

$$\rightsquigarrow f_T \text{ } s \text{ } (\text{It } ts)$$

The following can be proved similarly:

$$\text{It nil } = \text{ } v$$

$$\text{It } (\text{Cons } t \text{ } ts) \text{ } = \text{ } f_{XF} \text{ } (\text{It } t) \text{ } (\text{It } ts)$$

$$\square$$

With iterator `It`, all the basic XML operators in XQ are represented in *System F*.

## 4.2 Encoding Basic XML Operators

### 4.2.1 General Functions

Two general functions, `select_first` $: \alpha \to \beta \to \alpha$ and `select_second` $:$ $\alpha \to \beta \to \beta$, where $\alpha$ and $\beta$ are *primitive* types that can be instantiated to

concrete types, are introduced. The first function returns the first argument, and the second function returns the second argument as follows:

$$\texttt{select\_first} = \lambda\, x^\alpha.\, \lambda\, y^\beta.\, x$$

$$\texttt{select\_second} = \lambda\, x^\alpha.\, \lambda\, y^\beta.\, y$$

## 4.2.2 Horizontal Operators

Simple horizontal operators, `head`, `tail`, `concat`, `reverse` and `select`, are defined. These operators deal only with elements at the XML forest top level (horizontal), and never recur in the tree structure (vertical) of each element. In this case, iterator `It` is simplified to a horizontal iterator $\texttt{It}_h$ as

$$\texttt{It}_\texttt{h}\ (\texttt{Xnode}\ s\ ts)\ =\ \texttt{Xnode}\ s\ ts \tag{4.9}$$

$$=\ \texttt{let}\ f_1 = \lambda\, s^{String}.\, \lambda\, b^V.\texttt{Xnode}\ s\ ts\ \texttt{in}\ f_1\ s\ (\texttt{It}_\texttt{h}\ ts)$$

$$\texttt{It}_\texttt{h}\ \texttt{nil}\ =\ c_2 \tag{4.10}$$

$$\texttt{It}_\texttt{h}\ (\texttt{Cons}\ t\ ts)\ =\ f_2\ t\ (\texttt{It}_\texttt{h}\ ts) \tag{4.11}$$

$$=\ f_2\ (\texttt{It}_\texttt{h}\ t)\ (\texttt{It}_\texttt{h}\ ts)$$

When $\texttt{It}_h$ traversals along the structure of object $m$ of type $Xtf$, for $m$ an XML tree, $\texttt{It}_h$ returns $m$ without recurring through $m$'s tree structure. For $m$ an XML forest, `Cons` $t$ $ts$, function $f_2 : Xtf \rightarrow V \rightarrow V$ takes two arguments: the first element tree $t$ and the result returned by $\texttt{It}_h$ recursively applied to the remaining forest $ts$. Equations (4.9-4.11) are rewritten (the second line of each equation) to fit the form of the standard iterator Equations (4.6-4.8). For object $m$ of type $Xtf$, function $\texttt{It}_\texttt{h}$ is defined as

$$\texttt{It} \stackrel{\text{def}}{=} \lambda\, m^{Xtf}.\, m\ [V]\ (\lambda\, s^{String}.\, \lambda\, b^V.\ \texttt{Xnode}\ s\ ts)\ c_2\ f_2 \tag{4.12}$$

With defining the corresponding functions $c_2$, $f_2$ and proper types $V$, a specified horizontal operator is expressed by $\text{It}_h$ Equation (4.12).

**head** This operator obtains the first tree element from an XML forest, and has type $\text{head}: \; Xtf \to Xtf$:

$$
\begin{aligned}
\text{head nil} \;\; &= \;\; \text{default} \\
\text{head (Cons } t \; ts) \;\; &= \;\; t \\
&= \;\; (\lambda\, t^{Xtf}.\, \lambda\, z^{Xtf}.\, t) \; t \; (\text{head } ts)
\end{aligned}
$$

In the case of empty forest, $\text{nil}$, $\text{head}$ returns a *unit* type value $\text{default}$ meaning *nothing*, where type *unit* is consistent with any other type. Using Equation (4.12) with $V$ being $Xtf$, $\text{head}$ is defined as

$$
\text{head} \quad \overset{\text{def}}{=} \quad \lambda\, m^{Xtf}.\, m \; (\lambda\, s^{String}.\, \lambda\, b^{Xtf}.\text{Xnode } s \; ts) \; \text{default}
$$
$$
(\lambda\, t^{Xtf}.\, \lambda\, z^{Xtf}.\, t)
$$

**tail** This operator returns all but the first tree element from an XML forest with type $\text{tail}: \; Xtf \to Xtf$. In the case of empty forest, $\text{nil}$, $\text{tail}$ returns only $\text{nil}$.

$$
\begin{aligned}
\text{tail nil} \;\; &= \;\; \text{nil} \\
\text{tail (Cons } t \; ts) \;\; &= \;\; ts
\end{aligned}
$$

The second equation does not match the form of Equation (4.11). This case is treated by a new approach.

Let $m : Xtf$, $h : Xtf \to Xtf \times Xtf$, and $h\ m = (m, \texttt{tail}\ m)$, then

$$h\ \texttt{nil}\ =\ (\texttt{nil}, \texttt{tail}\ \texttt{nil})$$

$$=\ (\texttt{nil},\ \texttt{nil})$$

$$h\ (\texttt{Cons}\ t\ ts)\ =\ (\texttt{Cons}\ t\ ts,\ \texttt{tail}\ (\texttt{Cons}\ t\ ts))$$

$$=\ (\texttt{Cons}\ t\ ts,\ ts)$$

$$=\ (\texttt{Cons}\ t\ ts,\ \texttt{select\_first}\ ts\ (\texttt{tail}\ ts))$$

$$=\ \texttt{let}\ f_2 = \lambda\,t^{Xtf}.\ \lambda(k,\ z)^{Xtf \times Xtf}.\ (\texttt{Cons}\ t\ k,$$

$$\texttt{select\_first}\ k\ z)$$

$$\texttt{in}\ f_2\ t\ (h\ ts)$$

which fits the forms of Equations (4.10-4.11) with $V$ as $Xtf \times Xtf$. With $\texttt{tail}\ (\texttt{Xnode}\ s\ ts) = \texttt{Xnode}\ s\ ts$, $h$ has

$$h\ \texttt{Xnode}s\ ts\ =\ (\texttt{Xnode}\ s\ ts,\ \texttt{Xnode}\ s\ ts)$$

$$=\ \lambda\,s^{String}.\ \lambda(k,z)^{Xtf \times Xtf}.(\texttt{Xnode}\ s\ k,$$

$$\texttt{Xnode}\ s\ k)\ s\ (ts, \texttt{tail}\ ts)$$

$$=\ \texttt{let}\ f_1 = \lambda\,s^{String}.\ \lambda(k,\ z)^{Xtf \times Xtf}.\ (\texttt{Xnode}\ s\ k,$$

$$\texttt{Xnode}\ s\ k)$$

$$\texttt{in}\ f_1\ s\ (h\ ts)$$

By applying Equation (4.12), the following is defined:

$$h\ \stackrel{\text{def}}{=}\ \lambda\,m^{Xtf}.\ m\ [Xtf \times Xtf]$$

$$\texttt{let}\ f_1 = \lambda\,s^{String}.\ \lambda(k,\ z)^{Xtf \times Xtf}.(\texttt{Xnode}\ s\ k,\ \texttt{Xnode}\ s\ k)$$

$$f_2 = \lambda\,t^{Xtf}.\ \lambda(k,\ z)^{Xtf \times Xtf}.(\texttt{Cons}\ t\ k, \texttt{select\_first}\ k\ z)$$

$$\texttt{in}\ f_1\ (\texttt{nil},\ \texttt{nil})\ f_2$$

and `tail` is the second element in the pair of type $[Xtf \times Xtf]$. There-fore,

$$
\text{tail} \quad \overset{\text{def}}{=} \quad \lambda\, m^{Xtf}.\, (\ m\ [Xtf \times Xtf]
$$

$$
\text{let } f_1 = \lambda\, s^{String}.\ \lambda(k,\ z)^{Xtf \times Xtf}.\ (\text{Xnode } s\ k,\ \text{Xnode } s\ k)
$$

$$
f_2 = \lambda\, t^{Xtf}.\, \lambda(k,\ z)^{Xtf \times Xtf}.(\text{Cons } t\ k, \text{select\_first } k\ z)
$$

$$
\text{in } f_1\ (\text{nil},\ \text{nil})\ f_2).1
$$

**concat (@)** For two XML forest $ts_1$ and $ts_2$, `concat` $ts_1\ ts_2$[3] concatenates $t_1$ and $t_2$ with type of `concat` : $Xtf \to Xtf \to Xtf$, represented by

$$
\text{concat nil } y \;=\; y
$$

$$
\text{concat (Cons } t\ ts)\ y \;=\; \text{Cons } t\ (\text{concat } ts\ y)
$$

These equations can be written more abstractly as

$$
\text{concat nil} \;=\; \lambda\, y^{Xtf}.\, y
$$

$$
\text{concat (Cons } t\ ts) \;=\; \text{Cons } t\ (\text{concat } ts)
$$

which matches the forms of Equations(4.10-4.11) with type $V$ as $Xtf \to Xtf$ so that the following can be defined:

$$
\text{concat} \quad \overset{\text{def}}{=} \quad \lambda\, m^{Xtf}.\, m\ [Xtf \to Xtf]\ (\lambda\, s^{String}.\ \lambda\, b^{Xtf}.\ \text{Xnode } s\ ts)
$$

$$
(\lambda\, y^{Xtf}.\, y)\ \text{Cons}
$$

**reverse** This operator reverses an XML forest on the top level with type `reverse` : $Xtf \to Xtf$. For example, `reverse` $[t_1,\ t_2,\ t_3]$ returns

---

[3]Its infix form is @, represented by $ts_1$ @ $ts_2$.

$[t_3, t_2, t_1]$.

In order to define `reverse`, a subsidiary function, `rappend` : $Xtf \rightarrow Xtf \rightarrow Xtf$, is presented, such that

$$\texttt{rappend } [a, b, c] \ [3, 2, 1] = [1, 2, 3, a, b, c]$$

If the first argument is `nil`, `rappend` *reverses* the second argument so that

$$\texttt{rappend } y \texttt{ nil} = y$$

$$\texttt{rappend } y \ (\texttt{Cons } t \ ts) = \texttt{Cons } t \ (\texttt{rappend } y \ ts)$$

With Equation (4.12) and type $V$ as $Xtf$, function `rappend` is denoted as

$$\texttt{rappend } y \quad \overset{\text{def}}{=} \quad \lambda \, m^{Xtf}. \ m \ [Xtf] \ (\lambda \, s^{String}. \ \lambda \, b^{Xtf}. \ \texttt{Xnode } s \ ts) \ y$$
$$(\lambda \, t^{Xtf}. \ \lambda \, z^{Xtf}. \ \texttt{Cons } t \ z)$$

or

$$\texttt{rappend} \quad \overset{\text{def}}{=} \quad \lambda \, y^{Xtf}. \lambda \, m^{Xtf}. m \ [Xtf] \ (\lambda \, s^{String}. \lambda \, b^{Xtf}. \texttt{Xnode } s \ ts) \ y$$
$$(\lambda \, t. \ \lambda \, z. \ \texttt{Cons } t \ z)$$

With `rappend`, `reverse` is expressed as follows:

$$\texttt{reverse nil} = \texttt{nil}$$

$$\texttt{reverse } (\texttt{Cons } t \ ts) = \texttt{rappend } (\texttt{Cons } t \ \texttt{nil}) \ (\texttt{reverse } ts)$$

$$= (\lambda \, t^{Xtf}. \ \texttt{rappend Cons } t \ \texttt{nil}) \ t \ (\texttt{reverse } ts)$$

which fit the recursive forms of Equations (4.10-4.11). By Equation

(4.12),

$$\texttt{reverse} \quad \overset{\text{def}}{=} \quad \lambda\, m^{Xtf}.\ m\ [Xtf]\ (\lambda\, s^{String}.\ \lambda\, b^{Xtf}.\ \texttt{Xnode}\ s\ ts)\ \texttt{nil}$$

$$(\lambda\, t^{Xtf}.\ \texttt{rappend Cons}\ t\ \texttt{nil})$$

**select** For XML forest $ts$, and string $s$, $\texttt{select} : String \rightarrow Xtf \rightarrow Xtf$ extracts the tree elements whose roots are labelled $s$.

First a function $\texttt{rootLabel}: Xtf \rightarrow String$ is defined to extract the roots label as follows:

$$\texttt{rootLabel}\ (\texttt{Xnode}\ s\ ts) \quad = \quad s$$

$$= \quad \texttt{select\_first}\ s\ (\texttt{rootLabel}\ ts)$$

$$\texttt{rootLabel nil} \quad = \quad \texttt{default}$$

$$\texttt{rootLabel}\ (\texttt{Cons}\ t\ ts) \quad = \quad (\lambda\, x^{String}.\ \lambda\, y^{String}.\ \texttt{default})$$

$$(\texttt{rootLabel}\ t)\ (\texttt{rootLabel}\ ts)$$

When $\texttt{rootLabel}$ applies to an XML forest, it returns $\texttt{default}$ which means *nothing*. The function $\texttt{rootLabel}$ is then signified as

$$\texttt{rootLabel} \quad \overset{\text{def}}{=} \quad \lambda\, m^{Xtf}.\ m\ [String]\ \texttt{select\_first default}$$

$$(\lambda\, x^{String}.\ \lambda\, y^{String}.\ \texttt{default})$$

With function `rootLabel`, `select` is described as

$$\texttt{select } s \texttt{ nil } = \texttt{ nil}$$

$$\texttt{select } s \texttt{ (Cons } t \texttt{ } ts) = \texttt{ if } s == (\texttt{rootLabel } t) \texttt{ then}$$

$$\texttt{Cons } t \texttt{ (select } s \texttt{ } ts) \texttt{ else select } s \texttt{ } ts$$

$$= \texttt{ if } s == (\texttt{rootLabel } t) \texttt{ then}$$

$$\texttt{let } f = \lambda z^{Xtf}. \texttt{ Cons } t \texttt{ } z$$

$$\texttt{in } f \texttt{ } t \texttt{ (select } s \texttt{ } ts)$$

$$\texttt{else select\_second } t \texttt{ (select } s \texttt{ } ts)$$

where "==" is a boolean operator comparing strings equality, and can also be represented in *System F*: when a string is represented as a list of characters, the operator "==" is defined as a polymorphic function on list. Evidently, the previous equations fit (4.10-4.11) with type $V$ as *Xtf*. Therefore, `select` is defined as

$$\texttt{select} \overset{\text{def}}{=} \lambda s^{String}. \ \lambda m^{Xtf}. \ m \ [Xtf]$$

$$\texttt{let } f_1 = \lambda s^{String}. \ \lambda b^{Xtf}.\texttt{Xnode } s \texttt{ } ts$$

$$f_2 = \lambda t^{Xtf}.(\texttt{if } s == (\texttt{rootLabel } t) \texttt{ then}$$

$$\lambda z^{Xtf}.\texttt{Cons } t \texttt{ } z \texttt{ else } \lambda z^{Xtf}.\texttt{select\_second } t \texttt{ } z)$$

$$\texttt{in } f_1 \texttt{ nil } f_2$$

### 4.2.3   Vertical Operators

*Vertical* operators are considered to be the operators that recursively deal with the tree structure of each tree in an XML forest. For example, perform-

ing a depth-first traversal for tree $t$ attains a forest, including $t$'s subtrees, in a DFS order. Now we show how vertical operators roots, children, subtreesdfs are represented by polymorphic functions in *System F*.

**roots** The operator roots : $Xtf \rightarrow Xtf$ returns all the root nodes, indicating one node for each tree, and a list of root nodes for a forest,

$$
\begin{aligned}
\text{roots (Xnode } s\ ts) \quad &= \quad \text{Xnode } s\ \text{nil} \\
&= \quad \text{let } f_1 = \lambda s^{String}.\ \lambda b^{Xtf}.\ \text{Xnode } s\ \text{nil in} \\
&\qquad f_1\ s\ (\text{roots } ts) \\
\text{roots nil} \quad &= \quad \text{nil} \\
\text{roots (Cons } t\ ts) \quad &= \quad \text{Cons (roots } t)\ (\text{roots } ts)
\end{aligned}
$$

fits Equations (4.6-4.8) with $U$ and $V$ as type $Xtf$, and roots is defined as

$$
\begin{aligned}
\text{roots} \quad &\overset{\text{def}}{=} \quad \lambda m^{Xtf}.\ m\ [Xtf]\ (\lambda s^{String}.\ \lambda b^{Xtf}.\ \text{Xnode } s\ \text{nil}) \\
&\qquad \text{nil Cons}
\end{aligned}
$$

**children** For XML forest $m$, children : $Xtf \rightarrow Xtf$ recursively obtains children $c_i$ : $Xtf$ for each tree element $t_i$, and concatenates all $c_i$, and is expressed as

$$
\begin{aligned}
\text{children (Xnode } s\ ts) \quad &= \quad ts \\
\text{children nil} \quad &= \quad \text{nil} \\
\text{children (Cons } t\ ts) \quad &= \quad \text{concat (children } t)\ (\text{children } ts)
\end{aligned}
$$

where the first equation does not match the form of Equation (4.6), but can be treated by the approach used for the tail function.

Let $g : Xtf \to Xtf \times Xtf$, such that $g\ m = (m,\ \texttt{children}\ m)$, then

$$g\ (\texttt{Xnode}\ s\ ts)\ =\ (\texttt{Xnode}\ s\ ts,\ ts)$$

$$=\ \lambda s^{String}.\ \lambda(k,\ z)^{Xtf \times Xtf}.\ (\texttt{Xnode}\ s\ k,\ k)\ s\ (g\ ts)$$

$$=\ \text{let}\ f_T = \lambda s^{String}.\ \lambda(k,\ z)^{Xtf \times Xtf}.\ (Xnode\ s\ k,\ k)$$

$$\text{in}\ f_T\ s\ (g\ ts)$$

$$g\ \texttt{nil}\ =\ (\texttt{nil}, \texttt{children}\ \texttt{nil})$$

$$=\ (\texttt{nil},\ \texttt{nil})$$

$$g\ (\texttt{Cons}\ t\ ts)\ =\ (\texttt{Cons}\ t\ ts,\ \texttt{children}\ (\texttt{Cons}\ t\ ts))$$

$$=\ (\texttt{Cons}\ t\ ts,\ \ \texttt{concat}\ (\texttt{children}\ t)\ (\texttt{children}\ ts))$$

$$=\ \lambda(k_1,\ z_1)^{Xtf\ \times Xtf}.\ \lambda(k_2,\ z_2)^{Xtf\ \times Xtf}.$$

$$(\texttt{Cons}\ k_1\ k_2,\ \texttt{concat}\ z_1\ z_2)\ (g\ t)\ (g\ ts)$$

$$=\ \text{let}\ f_{XF} = \lambda(k_1,\ z_1)^{Xtf\ \times Xtf}.\ \lambda(k_2,\ z_2)^{Xtf\ \times Xtf}.$$

$$(\texttt{Cons}\ k_1\ k_2,\ \texttt{concat}\ z_1\ z_2)$$

$$\text{in}\ f_{XF}\ (g\ t)\ (g\ ts)$$

which fit the forms of Equation (4.6-4.8), and function $g$ is defined as follows:

$$g\ \overset{\text{def}}{=}\ \lambda m^{Xtf}.\ m\ [Xtf \times Xtf]$$

$$\text{let}\ f_T = \lambda s^{String}.\ \lambda(k,\ z)^{Xtf \times Xtf}.\ (\texttt{Xnode}\ s\ k,\ k)$$

$$f_{XF} = \lambda(k_1,\ z_1)^{Xtf\ \times Xtf}.\ \lambda(k_2,\ z_2)^{Xtf\ \times Xtf}.$$

$$(\texttt{Cons}\ k_1\ k_2,\ \texttt{concat}\ z_1\ z_2)$$

$$\text{in}\ f_T\ (\texttt{nil},\ \texttt{nil})\ f_{XF}$$

and $\texttt{children}$ is the second element in the pair of type $[Xtf \times Xtf]$,

expressed as

$$\texttt{children} \quad \overset{\text{def}}{=} \quad \text{let } g = \lambda\, m^{Xtf}.\ m\ [Xtf \times Xtf]$$

$$\text{let } f_T = \lambda\, s^{String}.\ \lambda(k,\ z)^{Xtf \times Xtf}.\ (\texttt{Xnode } s\ k,\ k)$$

$$f_{XF} = \lambda(k_1,\ z_1)^{Xtf\ \times Xtf}.\ \lambda(k_2,\ z_2)^{Xtf\ \times Xtf}.$$

$$(\texttt{Cons } k_1\ k_2,\ \texttt{concat } z_1\ z_2)$$

$$\text{in } f_T\ (\texttt{nil, nil})\ f_{XF}$$

$$\text{in } g.1$$

**subtreesdfs**  For object $m$ of type $Xtf$, and if $m$ is an XML tree, $\texttt{subtreesdfs}$ :
$Xtf \rightarrow Xtf$ performs a depth-first traversal and returns all the subtrees
of $m$, including itself; if $m$ is an XML forest, $\texttt{subtreesdfs}$ traverses
each tree element $t_i$ of $m$, obtains subtrees for each $t_i$, and concatenates
them, such that:

$$
\begin{aligned}
\texttt{subtreesdfs } (\texttt{Xnode } s\ ts) \ &= \ \texttt{Cons } (\texttt{Xnode } s\ ts)\ (\texttt{subtreesdfs } ts) \\
\texttt{subtreesdfs nil} \ &= \ \texttt{nil} \\
\texttt{subtreesdfs } (\texttt{Cons } t\ ts) \ &= \ \texttt{concat } (\texttt{subtreesdfs } t) \\
&\qquad (\texttt{subtreesdfs } ts)
\end{aligned}
$$

where the second equation does not fit iterator Equation(4.6).  The
approach in $\texttt{children}$ is adopted by defining $g : Xtf \rightarrow Xtf \times Xtf$ such

that $g\ m = (m,\ \texttt{subtreesdfs}\ m)$. For $\texttt{Xnode}\ s\ ts$,

$$
\begin{aligned}
g\ (\texttt{Xnode}\ s\ ts) &= (\texttt{Xnode}\ s\ ts,\ \texttt{Cons}\ (\texttt{Xnode}\ s\ ts)\ (\texttt{subtreesdfs}\ ts)) \\
&= \lambda(k,\ z)^{Xtf \times Xtf}.\ (\texttt{Xnode}\ s\ k,\ \texttt{Cons}\ (\texttt{Xnode}\ s\ k)\ z) \\
&\qquad (ts,\ \texttt{subtreesdfs}\ ts) \\
&= \textsf{let}\ f_T = \lambda\ s^{String}.\ \lambda(k,\ z)^{Xtf \times Xtf}. \\
&\qquad\qquad (Xnode\ s\ k,\ \texttt{Cons}\ (\texttt{Xnode}\ s\ ts)\ z) \\
&\quad \textsf{in}\ f_T\ s\ (g\ ts)
\end{aligned}
$$

Then we define $\texttt{subtreesdfs}$ is

$$
\begin{aligned}
\texttt{subtreesdfs}\ \overset{\text{def}}{=}\ &\textsf{let}\ g = \lambda\ m^{Xtf}.\ m\ [Xtf \times Xtf] \\
&\quad \textsf{let}\ f_T = \lambda\ s^{String}.\ \lambda(k,\ z)^{Xtf \times Xtf}. \\
&\qquad\qquad (Xnode\ s\ k,\ \texttt{Cons}\ (\texttt{Xnode}\ s\ ts)\ z) \\
&\qquad f_{XF} = \lambda(k_1,\ z_1)^{Xtf\ \times Xtf}.\ \lambda(k_2,\ z_2)^{Xtf\ \times Xtf}. \\
&\qquad\qquad (\texttt{Cons}\ k_1\ k_2,\ \texttt{concat}\ z_1\ z_2) \\
&\quad \textsf{in}\ f_T\ (\texttt{nil},\ \texttt{nil})\ f_{XF} \\
&\textsf{in}\ g.1
\end{aligned}
$$

### 4.2.4 Boolean Operators

XQ defines two structural boolean operators, $\texttt{equal}$ and $\texttt{less}$, that define structural relationships of the forests. Also, operator $\texttt{empty}$ is introduced to check the emptiness of a forest.

**empty** This operator tests the emptiness of an XML forest [4] with $\texttt{empty}$ :

---

[4]When $\texttt{empty}$ is applied to an XML tree, it is assumed that $\textsf{False}$ is returned.

$Xtf \rightarrow$ `Bool`, and represented as

$$\text{empty (Xnode } s \ ts) \ = \ \text{False}$$
$$= \ (\lambda \, a^{String}. \, \lambda \, b^{Bool}. \, \text{False}) \, s \, (\text{empty } ts)$$
$$\text{empty nil} \ = \ \text{True}$$
$$\text{empty (Cons } t \ ts) \ = \ \text{False}$$
$$= \ (\lambda \, x^{Bool}. \, \lambda \, y^{Bool}. \, \text{False}) \, (\text{empty } t) \, (\text{empty } ts)$$

For object $m \, : \, Xtf$, functions $f_T \, : \, String \, \rightarrow \, Bool \, \rightarrow \, Bool$, and $f_{XF} \, : \, Bool \, \rightarrow \, Bool \, \rightarrow \, Bool$, the boolean condition `empty` is defined by Equations(4.6-4.8) as follows:

$$\text{empty} \ \stackrel{\text{def}}{=} \ \lambda \, m^{Xtf}. \, m \, [Bool]$$
$$\text{let } f_T = \lambda \, x^{String}. \, \lambda \, y^{Bool}. \, \text{False}$$
$$f_{XF} = \lambda \, x^{Bool}. \, \lambda \, y^{Bool}. \, \text{False}$$
$$\text{in } f_T \ \text{True} \ f_{XF}$$

**equal** The operator, `equal`, tests whether two objects of type $Xtf$ are structurally equal.

**Definition 18 (Structural Equality)** Two objects $ts_1$ and $ts_2$ of type $Xtf$ are structurally equal, if and only if they have the same constructor at the top level, and all the substructures are inductively equal.

Thus: $ts_1$ and $ts_2$ are *structurally equal* if and only if one of the following conditions holds:

- $ts_1$ and $ts_2$ are trees with the same root label, and their children nodes are recursively *equal*;

- $ts_1$ and $ts_2$ are `nil`;

- $ts_1$ and $ts_2$ are forests with constructor `Cons`, and their tree elements are recursively *structural equals*

The operator, `equal`, recursively tests two forests with type `equal` : $Xtf \rightarrow Xtf \rightarrow$ `Bool`. The definition of `equal` needs two functions that have been defined: `rootLabel` : $Xtf \rightarrow String$, which extracts the root label of a tree; and `empty` : $Xtf \rightarrow Bool$, which checks the emptiness of a forest. Operator `equal` can be expressed by using `rootLabel` and `empty` as follows:[5]

$$\text{equal } (\text{Xnode } s \ ts) \ m \ = \ (s == \text{rootLabel m}) \text{ and}$$
$$(\text{equal } ts \ (\text{children } m))$$
$$\text{equal nil } m \ = \ \text{empty } m$$
$$\text{equal } (\text{Cons } t \ ts) \ m \ = \ (\text{equal } t \ (\text{head } m)) \text{ and } (\text{equal } ts \ (\text{tail } m))$$

These equations are then written more abstractly as

$$\text{equal } (\text{Xnode } s \ ts) \ = \ \lambda \, m^{Xtf}. \, (s == \text{rootLabel m}) \text{ and}$$
$$(\text{equal } ts \ (\text{children } m))$$
$$= \ \text{let } f_T = \lambda \, s^{String}. \ \lambda \, z^{Xtf \rightarrow Bool}. \ \lambda \ m^{Xtf}.$$
$$(s == \text{rootLabel } m) \text{ and } (z \ (\text{children } m))$$
$$\text{in } f_T \ s \ (\text{equal } ts)$$
$$\text{equal nil } = \ \text{empty}$$

[5]Boolean operators `and` and `or` have been represented in *System F*[7].

$$\mathtt{equal}\ (\mathtt{Cons}\ t\ ts) \quad = \quad \lambda\, m^{Xtf}.\,(\mathtt{equal}\ t\ (\mathtt{head}\ m))\ \mathtt{and}$$

$$(\mathtt{equal}\ ts\ (\mathtt{tail}\ m))$$

$$= \quad \mathtt{let}\ f_{XF} = \lambda\, x^{Xtf \to Bool}.\ \lambda\ y^{Xtf \to Bool}.\ \lambda\, m^{Xtf}.$$

$$(x\ (\mathtt{head}\ m))\ \mathtt{and}\ (y\ (\mathtt{tail}\ m))$$

$$\mathtt{in}\ f_{XF}\ (\mathtt{equal}\ t)\ (\mathtt{equal}\ ts)$$

which matches the forms of Equations(4.6-4.8) with

$$f_T : String \to (Xtf \to Bool) \to Xtf \to Bool$$

$$f_{XF} : (Xtf \to Bool) \to (Xtf \to Bool) \to Xtf \to Bool$$

leading to

$$\mathtt{equal}\ \overset{\mathrm{def}}{=}\ \lambda\, n^{Xtf}.\ n\ [Xtf \to Bool]$$

$$\mathtt{let}\ f_T = \lambda\, s^{String}.\ \lambda\, z^{Xtf \to Bool}.\ \lambda\ m^{Xtf}.$$

$$(s == \mathtt{rootLabel}\ \mathtt{m})\ \mathtt{and}\ (z\ (\mathtt{children}\ m))$$

$$f_{XF} = \lambda\, x^{Xtf \to Bool}.\ \lambda\ y^{Xtf \to Bool}.\ \lambda\, m^{Xtf}.$$

$$(x\ (\mathtt{head}\ m))\ \mathtt{and}\ (y\ (\mathtt{tail}\ m))$$

$$\mathtt{in}\ f_T\ \mathtt{empty}\ f_{XF}$$

**less** The boolean operator **less** examines the structural ordering of XML trees or forests.

**Definition 19 (Structurally Less)** For two objects $m$ and $n$ of type $Xtf$, $m$ is structurally less than $n$, if and only if one of the following conditions recursively holds:

- $m$ and $n$ are trees, and one of the following conditions exclusively holds:

  - their root nodes satisfy the *less*[6] condition: `rootLabel` $m$ is a string that is literally less than `rootLabel` $n$;

  - `children` $m$ is structurally *less* than `children` $n$

- $m$ and $n$ are forests, and one of the following conditions exclusively holds,

  - $m$ is `nil`, and $n$ has at least one tree element;

  - their first element satisfies the structurally *less* ordering; that is, `head` $m$ is `less` than `head` $n$;

  - their remaining elements satisfy the structurally *less* ordering, that is, `tail` $m$ is `less` than `tail` $n$

From the definition of structurally less, for two objects of type *Xtf*, $m$ and $n$, the function `less` is expressed in *System F* by

$$\texttt{less } (\texttt{Xnode } s \ ts) \ n \ = \ (s < \texttt{rootLabel } n) \text{ or } (\texttt{less } ts \ (\texttt{children } n))$$

$$\texttt{less nil } n \ = \ \texttt{not } (\texttt{empty } n)$$

$$\texttt{less } (\texttt{Cons } t \ ts) \ n \ = \ (\texttt{less } t \ (\texttt{head } n)) \text{ or } (\texttt{less } ts \ (\texttt{tail } n))$$

---

[6]The boolean operator *less*,"$<$", for strings can also be defined in *System F*.

where `not`, `or` are boolean operators defined in *System F*. The previous functions are rewritten as

$$\texttt{less (Xnode } s \ ts) \ = \ \lambda \, n^{Xtf}. \ (s < \texttt{rootLabel n}) \text{ or } (\texttt{less } ts \ (\texttt{children } n))$$

$$= \ \texttt{let } f_T = \lambda \, s^{String}. \ \lambda \, z^{Xtf \to Bool}. \ \lambda \ n^{Xtf}.$$

$$(s < (\texttt{rootLabel n})) \text{ or } (z \ (\texttt{children } n))$$

$$\texttt{in } f_T \ s \ (\texttt{less } ts)$$

$$\texttt{less nil} \ = \ \texttt{not empty}$$

$$\texttt{less (Cons } t \ ts) \ = \ \lambda \, n^{Xtf}. \ (\texttt{less } t \ (\texttt{head } n)) \text{ or } (\texttt{less } ts \ (\texttt{tail } n))$$

$$= \ \texttt{let } f_{XF} = \lambda \, x^{Xtf \to Bool}. \ \lambda \ y^{Xtf \to Bool}. \ \lambda \, n^{Xtf}.$$

$$(x \ (\texttt{head } n)) \text{ or } (y \ (\texttt{tail } n))$$

$$\texttt{in } f_{XF} \ (\texttt{less } t) \ (\texttt{less } ts)$$

By using Equation(4.12), `less` is defined as follows:

$$\texttt{less} \ \overset{\text{def}}{=} \ \lambda \, m^{Xtf}. \ m \ [Xtf \to Bool]$$

$$\texttt{let } f_T = \lambda \, s^{String}. \ \lambda \, z^{Xtf \to Bool}. \ \lambda \ n^{Xtf}. \ (s < \texttt{rootLabel n}) \text{ or}$$

$$(z \ (\texttt{children } n))$$

$$f_{XF} = \lambda \, x^{Xtf \to Bool}. \ \lambda \ y^{Xtf \to Bool}. \ \lambda \, n^{Xtf}. \ (x \ (\texttt{head } n)) \text{ or}$$

$$(y \ (\texttt{tail } n))$$

$$\texttt{in } f_T \ (\texttt{not empty}) \ f_{XF}$$

## 4.2.5 Application Operators

Since XML trees, forests, and their basic operators are represented by *System F* polymorphic functions, the focus is now on some complex functions such

as `sort`: $Xtf \rightarrow (Xtf \rightarrow Xtf \rightarrow Bool) \rightarrow Xtf$, which sorts an XML forest using a tree structural boolean condition, and `distinct` : $Xtf \rightarrow Xtf$, which filters out all the duplicate tree elements from an XML forest.

**sort** For object $ms$ of type $Xtf$ and tree boolean condition $f$, `sort` $ms$ $f$ reorders $ms$ by using tree structural order function $f$, where `sort` recursively selects current tree element $m$ in forest $ms$, compares $m$ with the remaining elements of $ms$, and then inserts $m$ into an appropriate position.

In order to express `sort`, a subsidiary function is required,

$$\texttt{insert} : Xtf \rightarrow (Xtf \rightarrow Xtf \rightarrow Bool) \rightarrow Xtf \rightarrow Xtf$$

where `insert` $m$ $f$ $ts$ inserts tree element $m$ into an appropriate position of forest $ts$ by tree boolean condition $f$. By using `insert`, `sort` is described as:

$$
\begin{aligned}
\texttt{sort}\ f\ \texttt{nil} \ &= \ \texttt{nil} \\
\texttt{sort}\ f\ (\texttt{Cons}\ t\ ts) \ &= \ \texttt{insert}\ t\ (\texttt{sort}\ f\ ts) \\
&= \ \texttt{let}\ f_{XF} = \lambda\, k^{Xtf}.\ \lambda\, z^{Xtf}.\ \texttt{insert}\ k\ z \\
& \qquad \texttt{in}\ f_{XF}\ (\texttt{sort}\ f\ t)\ (\texttt{sort}\ f\ ts)
\end{aligned}
$$

`sort` is defined as a forest operator, which returns its second argument back when applied to a tree. These equations match the forms of the horizontal iteration Equations(4.10-4.11) with $f_{XF} : Xtf \rightarrow Xtf \rightarrow Xtf$.

`sort` is defined by Equation (4.12) as follows:

$$\texttt{sort} \quad \overset{\text{def}}{=} \quad \lambda f^{Xtf \to Xtf \to Bool}.\ \lambda m^{Xtf}.\ m\ [Xtf]$$

$$\text{let } f_T = \lambda s^{String}.\ \lambda z^{Xtf}.\ \texttt{Xnode } s\ ts$$

$$f_{XF} = \lambda k^{Xtf}.\ \lambda z^{Xtf}.\ \texttt{insert } t\ z$$

$$\text{in } f_T\ \texttt{nil}\ f_{XF}$$

The remaining task is to program the function

$$\texttt{insert}:\ Xtf \to (Xtf \to Xtf \to Bool) \to Xtf \to Xtf$$

where `insert` $m\ f\ n$ inserts tree $m$ into forest $n$ when boolean condition $f$ is satisfied. `insert` behaves like a horizontal operator on forests, and is expressed as:

$$\texttt{insert } m\ f\ \texttt{nil}\ =\ \texttt{Cons } m\ \texttt{nil}$$

$$\texttt{insert } m\ f\ (\texttt{Cons } t\ ts)\ =\ \texttt{if } f\ m\ t \texttt{ then Cons } m\ (\texttt{Cons } t\ ts)$$

$$\texttt{else Cons } t\ (\texttt{insert } m\ f\ ts)$$

where the second equation does not match (4.11), and the same approach holds for the definition of `tail`. Let

$$g\ :\ Xtf \to (Xtf \to Xtf \to Bool) \to Xtf \to (Xtf \times Xtf)$$

such that $g \; m \; f \; ts = (ts, \texttt{insert} \; m \; f \; ts)$,

$$g \; m \; f \; \texttt{nil} \;\; = \;\; (\texttt{nil}, \texttt{Cons} \; m \; \texttt{nil})$$

$$g \; m \; f \; (\texttt{Cons} \; t \; ts) \;\; = \;\; (\texttt{Cons} \; t \; ts, \; \text{if} \; f \; m \; t \; \text{then} \; \texttt{Cons} \; m \; (\texttt{Cons} \; t \; ts)$$
$$\text{else} \; \texttt{Cons} \; t \; (\texttt{insert} \; m \; f \; ts))$$
$$= \;\; \lambda(k, \; z)^{Xtf \times Xtf}.(\texttt{Cons} \; t \; k, \; \text{if} \; f \; m \; t \; \text{then} \; \texttt{Cons} \; m$$
$$(\texttt{Cons} \; t \; k) \; \text{else} \; \texttt{Cons} \; t \; z) \; (ts, \; \texttt{insert} \; m \; f \; ts))$$
$$= \;\; \text{let} \; f_{XF} = \lambda \, x^{Xtf \times Xtf}. \; \lambda(k, \; z)^{Xtf \times Xtf}. \; (\texttt{Cons} \; t \; k,$$
$$\text{if} \; f \; m \; t \; \text{then} \; \texttt{Cons} \; m \; (\texttt{Cons} \; t \; k) \; \text{else} \; \texttt{Cons} \; t \; z)$$
$$\text{in} \; f_{XF} \; (g \; m \; f \; t) \; (g \; m \; f \; ts)$$

These fit Equations (4.10-4.11). By employing (4.12), the following is defined:

$$g \;\; \stackrel{\text{def}}{=} \;\; \lambda \, m^{Xtf}. \; \lambda f^{Xtf \rightarrow Xtf \rightarrow Bool}. \; \lambda \, n^{Xtf}. \; n \, [Xtf \times Xtf]$$
$$\text{let} \; f_T = \lambda \, s^{String}. \; \lambda(k, \; z)^{Xtf \times Xtf}. \; (\texttt{Xnode} \; s \; k, \; \texttt{Xnode} \; s \; k)$$
$$f_{XF} = \lambda \, x^{Xtf \times Xtf}. \; \lambda(k, \; z)^{Xtf \times Xtf}.(\texttt{Cons} \; t \; k,$$
$$\text{if} \; f \; m \; t \; \text{then} \; \texttt{Cons} \; m \; (\texttt{Cons} \; t \; k) \; \text{else} \; \texttt{Cons} \; t \; z)$$
$$\text{in} \; f_T \; (\texttt{nil}, \texttt{Cons} \; m \; \texttt{nil}) \; f_{XF}$$

and therefore, $\texttt{insert}$ is the second element in the pair of type $[Xtf \times$

*Xtf*], and

$$\texttt{insert} \;\; \overset{\text{def}}{=} \;\; \text{let } g = \lambda\, m^{Xtf}.\; \lambda\, f^{Xtf \to Xtf \to Bool}.\; \lambda\, n^{Xtf}.\; n\; [Xtf \times Xtf]$$

$$\text{let } f_T = \lambda\, s^{String}.\, \lambda(k,\; z)^{Xtf \times \texttt{Xtf}}.(\texttt{Xnode } s\; k,\; \texttt{Xnode } s\; k)$$

$$f_{XF} = \lambda\, x^{Xtf \times Xtf}.\, \lambda(k,\; z)^{Xtf \times Xtf}.(\texttt{Cons } t\; k,$$

$$\text{if } f\; m\; t \text{ then } \texttt{Cons } m\; (\texttt{Cons } t\; k) \text{ else } \texttt{Cons } t\; z)$$

$$\text{in } f_T\; (\texttt{nil}, \texttt{Cons } m\; \texttt{nil})\; f_{XF}$$

$$\text{in } g.1$$

**distinct** For XML forest $ts$, $\texttt{distinct }ts$ eliminates all the duplicate tree elements with $Xtf \to Xtf$. Given a function, $\texttt{filter} : (Xtf \to Xtf \to Bool) \to Xtf \to Xtf \to Xtf$, where $\texttt{filter } f\; m\; ts$ eliminates all tree elements $t$ that satisfy boolean condition $f\; m\; t$, the horizontal function $\texttt{distinct}$ is then expressed by $\texttt{filter}$ so that

$$\texttt{distinct nil} \;\; = \;\; \texttt{nil}$$

$$\texttt{distinct (Cons } t\; ts) \;\; = \;\; \texttt{Cons } t\; (\texttt{distinct (filter equal } t)\; ts)$$

$$= \;\; \text{let } \texttt{shift} = \lambda\, h^{\alpha \to \beta}.\, \lambda\, g^{\beta \to \gamma}.\, \lambda\, x^{\alpha}.\; g\; h\; x \text{ in}$$

$$\texttt{Cons } t\; \texttt{shift (filter equal } t)\; \texttt{distinct } ts$$

$$= \;\; \text{let } \texttt{shift} = \lambda\, h^{\alpha \to \beta}.\, \lambda\, g^{\beta \to \gamma}.\, \lambda\, x^{\alpha}.g\; h\; x \text{ in}$$

$$\text{let } f^{Xf} = \lambda\, k^{Xtf}.\; \lambda\, z^{Xtf}.$$

$$\texttt{Cons } k\; \texttt{shift (filter equal } k)\; z$$

$$\text{in } f_{XF}\; (\texttt{distinct } t)\; (\texttt{distinct } ts)$$

When $\texttt{distinct}$ is applied to tree $t$, $\texttt{distinct}$ returns $t$. Also, a general purpose function $\texttt{shift}$ is needed to switch the order of the arguments

such that

$$\texttt{shift (filter equal } t) \texttt{ distinct } ts$$

$$= \quad (\lambda\, h^{\alpha \to \beta}.\ \lambda\, g^{\beta \to \gamma}.\ \lambda\, x^{\alpha}.\ g\ h\ x)\ \texttt{(filter equal } t) \texttt{ distinct } ts$$

$$\rightsquigarrow \quad \texttt{distinct (filter equal } t)\ ts$$

With `switch`, `distinct` is signified as

$$\texttt{distinct} \quad \overset{\text{def}}{=}\quad \lambda\, m^{Xtf}.\ m\ [Xtf]$$

$$\text{let } f_T = \lambda\, s^{String}.\, \lambda\, z^{Xtf}.\ \texttt{Xnode } s\ ts$$

$$f_{XF} = \textsf{let shift} = \lambda\, h^{\alpha \to \beta}.\ \lambda\, g^{\beta \to \gamma}.\ \lambda\, x^{\alpha}.\ g\ h\ x \textsf{ in}$$

$$\lambda\, k^{Xtf}.\ \lambda\, z^{Xtf}.\ \texttt{Cons } k\ \texttt{shift (filter equal } k)\ z$$

$$\textsf{in } f_T\ \texttt{nil}\ f_{XF}$$

The definition of `distinct` still requires operator `filter` : $(Xtf \to Xtf \to Bool) \to Xtf \to Xtf \to Xtf$, where `filter` $f$ $m$ $ts$ eliminates the tree elements that satisfy boolean condition $f$ $m$ $t$. Operator `distinct` is then described as,

$$\texttt{filter } f\ b\ \texttt{nil} \quad = \quad \texttt{nil}$$

$$\texttt{filter } f\ b\ (\texttt{Cons } t\ ts) \quad = \quad \textsf{if } f\ b\ t \textsf{ then filter } f\ b\ ts$$

$$\textsf{else Cons } t\ (\texttt{filter } f\ b\ ts)$$

$$\text{let } f_{XF} = \lambda\, k^{Xtf}.\ \lambda\, z^{Xtf}.\ \textsf{if } f\ b\ k \textsf{ then } z$$

$$\textsf{else Cons } k\ z$$

$$\textsf{in } f_{XF}\ (\texttt{filter } f\ b\ t)\ (\texttt{filter } f\ b\ ts)$$

As a horizontal forest operator, `filter` returns only the second argument, if it is a tree. The previous equations match (4.10-4.11), and are

defined by Equation (4.12) with type $f_{XF} : Xtf \rightarrow Xtf \rightarrow Xtf$ such that

$$\texttt{filter} \quad \stackrel{\text{def}}{=} \quad \lambda f^{Xtf \rightarrow Xtf \rightarrow Bool}.\ \lambda b^{Xtf}.\ \lambda m^{Xtf}.\ m\ [Xtf]$$

$$\textsf{let}\ f_T = \lambda s^{String}.\lambda z^{Xtf}.\ \texttt{Xnode}\ s\ ts$$

$$f_{XF} = \lambda k^{Xtf}.\ \lambda z^{Xtf}.\ \textsf{if}\ f\ b\ k\ \textsf{then}\ z$$

$$\textsf{else}\ \texttt{Cons}\ k\ z$$

$$\textsf{in}\ f_T\ \texttt{nil}\ f_{XF}$$

## 4.3   Translation of XQ to *System F*

After the XML tree and forest are represented in *System F*, the syntax definition of XQ expressions is rewritten into more general form, where the FLWR keyword allows the arbitrary composition of query expressions; and also, all the basic operators on XML data are presented in *System F*. Encoded in *System F*, a complex XML query expression can be expressed as a sequence of function applications without any special purpose operators. The syntax of XQ in *System F* is defined in Figure 4.2, where $\textsf{XF}_n$ expresses the basic XML operators defined as polymorphic function *abstractions*, $x$ is a variable, $\varphi$ is a boolean expression, $U$ is a type, and $e_{1_{SF}}$, $e_{2_{SF}}$, ..., $e_{k_{SF}}$ denote XQ expressions encoded in *System F*. For example,the `head` is represented in *System F* as

$$\texttt{head}\ m \quad \stackrel{\text{def}}{=} \quad m\ [Xtf]\ (\lambda s^{String}.\lambda b^{Xtf}.\texttt{Xnode}\ s\ ts)\ \texttt{default}$$

$$(\lambda t^{Xtf}.\lambda z^{Xtf}.\ t)$$

$$
\begin{aligned}
XQexp \quad &::= \quad x \\
&\mid \quad \text{let } x = e_{SF} \text{ in } e'_{SF} \\
&\mid \quad \text{for } x \text{ in } e_{1_{SF}} \text{ do } e_{2_{SF}} \\
&\mid \quad \text{where } \varphi \text{ return } e_{SF} \\
&\mid \quad \mathsf{XF_n} \; e_{1_{SF}} \; \ldots e_{k_{SF}} \\
Boolexp \quad &::= \quad \mathsf{True} \mid \mathsf{False} \\
&\mid \quad \texttt{empty } e_{SF} \\
&\mid \quad \texttt{less } e_{SF} \; e'_{SF} \\
&\mid \quad \texttt{equal } e_{SF} \; e'_{SF} \\
\mathsf{XF_n} \; e_{1_{SF}} \; e_{2_{SF}} \; \ldots \; e_{k_{SF}} \quad &::= \quad e_{SF} \; [U] \; e_T \; e_{nil} \; e_{XF}
\end{aligned}
$$

Figure 4.2: XQ in *System F*

where "[*Xtf*]" denotes the resulting type of `head`, and the following three terms, $e_T$, $e_{nil}$, and $e_{XF}$, correspond iterators inductively applied to `Xnode`, `nil` and `Cons`.

## 4.4   Semantic Analysis

Since XQ has been encoded in *System F*, the semantics of XQ query expressions are described as semantic equations, and encoded in *System F*.

### 4.4.1   Environment

To understand the semantics of XQ expressions, it is defined *semantic* equations that map XQ expressions to the meanings that these phrases denote. Clearly, XQ expressions have corresponding values which are XML forests,

similar to those of boolean expressions that have boolean values True or False.
However, the *meanings* or *denotations* of XQ expressions are much more com-
plex than the values, because the values of XQ expressions depend on the
values of the variables of the XQ expressions. More abstractly, the values de-
pend on an *environment, Env*, which maps each variable to its corresponding
value.

**Definition 20 (Environments)** An environment, *Env*, denotes the map-
ping of each variable with its value, and is described as a list of paired
variables and values with the following definitions:

$$Env \overset{\text{def}}{=} (Var \times Xtf) \; \texttt{list}$$

$$Var \times Xtf \overset{\text{def}}{=} (Var \to Xtf \to Pair) \to Pair$$

For $v : Var$, $ts : Xtf$, pair $(v, \; ts) : Var \times Xtf$ is expressed by the following
polymorphic function:

$$(v, \; ts) \overset{\text{def}}{=} \lambda f^{Var \to Xtf \to Pair}. \; f \; v \; ts$$

and the project operations are

$$(v, \; ts).0 \;=\; (v, \; ts) \, [Var] \, (\lambda x^{Var}. \; \lambda y^{Xtf}. \; x)$$

$$(v, \; ts).1 \;=\; (v, \; ts) \, [Xtf] \, (\lambda x^{Var}. \; \lambda y^{Xtf}. \; y)$$

where .0 returns the first element, variable $v$, in the pair, and .1 receives value $ts$, shown as

$$
\begin{aligned}
(v,\ ts).0 \ &\stackrel{\mathrm{def}}{=} \ (v,\ ts)\,[\mathit{Var}]\,(\lambda\,x^{\mathit{Var}}.\ \lambda\,y^{\mathit{Xtf}}.\ x) \\
&= \ (\lambda\,f^{\mathit{Var}\to\mathit{Xtf}\to\mathit{Pair}}.\ f\ v\ ts)\,[\mathit{Var}]\,(\lambda\,x^{\mathit{Var}}.\ \lambda\,y^{\mathit{Xtf}}.\ x) \\
&\leadsto \ (\lambda\,f^{\mathit{Var}\to\mathit{Xtf}\to\mathit{Var}}.\ f\ v\ ts)\,(\lambda\,x^{\mathit{Var}}.\ \lambda\,y^{\mathit{Xtf}}.\ x) \\
&\leadsto \ (\lambda\,x^{\mathit{Var}}.\ \lambda\,y^{\mathit{Xtf}}.\ x)\,(\lambda\,x^{\mathit{Var}}.\ \lambda\,y^{\mathit{Xtf}}.\ x)\ v\ ts \\
&\leadsto \ (\lambda\,y^{\mathit{Xtf}}.\ v)\ ts \\
&\leadsto \ v
\end{aligned}
$$

and similarly,

$$
(v,\ ts).1 \leadsto^{*} ts
$$

## 4.4.2 Operations of the Environment

The *environment* describes the mapping of variable $v$ to its corresponding value $ts$. In order to operate the values stored in the environment, two operations must be defined to manipulate the environment.

**put** Operation `put` $v$ $ts$ $env$ mapping value $ts$ to variable $v$ in environment $env$. If $v$ already exists and stores value $ts'$, $ts'$ is replaced by new value $ts$; otherwise, a new mapping $(v,\ ts)$ is added to $env$. The operation `put` : $\mathit{Var} \to \mathit{Xtf} \to \mathit{Env} \to \mathit{Env}$ is expressed as

$$
\begin{aligned}
\texttt{put}\ v\ ts\ \texttt{nil} \ &= \ \texttt{Cons}\ (v,\ ts)\ \texttt{nil} \\
\texttt{put}\ v\ ts\ (\texttt{Cons}\ a\ as) \ &= \ \text{if}\ v == a.0\ \text{then}\ \texttt{Cons}\ (v,\ ts)\ as \\
&\qquad\qquad\quad \text{else}\ \texttt{Cons}\ a\ (\texttt{put}\ v\ ts\ as)
\end{aligned}
$$

where *env* is the type of $(Var \times Xtf)$ `list`. To define the function `put`, a function $g$,

$$g : Var \rightarrow Xtf \rightarrow Env \rightarrow (Env \times Env)$$

with $g\ v\ ts\ m = (m,\ \texttt{put}\ v\ ts\ m)$, is introduced.

For variable $v : Var$ and value $ts : Xtf$, $g$ is expressed as

$$g\ v\ ts\ \texttt{nil}\ =\ (\texttt{nil, Cons}\ (v,\ ts)\ \texttt{nil})$$

$$g\ v\ ts\ (\texttt{Cons}\ a\ as)\ =\ (\texttt{Cons}\ a\ as,\ \texttt{if}\ v == a.0$$
$$\texttt{else Cons}\ a\ (\texttt{put}\ v\ ts\ as))$$
$$=\ \lambda(k,\ z)^{Env \times Env}.\ (\texttt{Cons}\ a\ k,\ \texttt{if}\ v == a.0$$
$$\texttt{then Cons}\ (v,\ ts)\ k\ \texttt{else Cons}\ a\ z)$$
$$a\ (as,\ \texttt{Cons}\ a\ (\texttt{put}\ v\ ts\ as))$$
$$=\ \texttt{let}\ f = \lambda\ a^{Var \times Xtf}.\ \lambda(k,\ z)^{Env \times Env}.\ (\texttt{Cons}\ a\ k,$$
$$\texttt{if}\ v == a.0\ \texttt{then Cons}\ (v,\ ts)\ k\ \texttt{else Cons}\ a\ z$$
$$\texttt{in}\ f\ a\ (g\ v\ ts\ as)$$

and matches the forms of `list` Equation (3.6), and is defined as

$$g\ \stackrel{\text{def}}{=}\ \lambda\ v^{Var}.\ \lambda\ ts^{Xtf}.\ \lambda\ m^{Env}.\ m\ [Env \times Env]$$
$$\texttt{let}\ f = \lambda\ a^{Var \times Xtf}.\ \lambda(k,\ z)^{Env \times Env}.\ (\texttt{Cons}\ a\ k,\ \texttt{if}\ v == a.0$$
$$\texttt{then Cons}\ (v,\ ts)\ k\ \texttt{else Cons}\ a\ z$$
$$\texttt{in}\ (\texttt{nil, Cons}\ (v,\ ts)\ \texttt{nil})\ f$$

and function `put` is the second element of $g$, so `put` $\stackrel{\text{def}}{=} g.1$

**get** Operation `get` $v$ $env$ retrieves the value of variable $v$ stored in the environment $env$ if $v$ exists; otherwise, `get` returns nil, where `get` has type `get` : $Var \rightarrow Env \rightarrow Xtf$, and is described as follows:

$$\text{get } v \text{ nil} \ = \ \text{nil}$$

$$\text{get } v \ (\text{Cons } a \ as) \ = \ \text{if } v == a.0 \text{ then } a.1 \text{ else get } v \ as$$

$$= \ \text{let } f = \lambda a^{Var \times Xtf}. \ \lambda z^{Xtf}. \text{ if } v == a.0 \text{ then } a.1$$

$$\text{else } z$$

$$\text{in } f \ a \ (\text{get } v \ as)$$

which fits the iterator form of list in Definition 11. Consequently, `get` is denoted as

$$\text{get} \ \overset{\text{def}}{=} \ \lambda v^{Var}. \ \lambda m^{Env}. \ m \ [Xtf]$$

$$\text{let } f = \lambda a^{Var \times Xtf}. \ \lambda z^{Xtf}. \text{ if } v == a.0 \text{ then } a.1 \text{ else } z$$

$$\text{in nil } f$$

### 4.4.3   Semantic Functions of XQ

With defining $Env$ and its operations in *System F*, an XQ query evaluation can be described as a set of computations from the initial state (environment) to the final stat(environment) by using *state-transformation* functions [7],

$$\llbracket - \rrbracket_{XQexp} E \ = \ XQexp \rightarrow Env \rightarrow Xtf$$

$$\llbracket - \rrbracket_{Boolexp} E \ = \ Boolexp \rightarrow Env \rightarrow Bool$$

where $E$ is an environment of type $Env$.

The semantics of XQ expressions are indeed a sequence of operations that

---

[7]The expression type is added to the semantic functions.

change the environments. In this thesis, the intention is to represent these
semantic functions and describe them in *System F*.

**SEM EQ: Variables**

$$[\![v]\!]_{XQexp} E \stackrel{\mathrm{def}}{=} E\ v = \texttt{get}\ v\ E$$

which means that the variable expression returns its value in environment $E$.

**SEM EQ: Let-Assignment**

$$
\begin{aligned}
[\![\mathsf{let}\ x = e\ \mathsf{in}\ e'\ ]\!]_{XQexp} E \ \stackrel{\mathrm{def}}{=}\ & [\![e']\!]_{XQexp}(E[x := ([\![e]\!]_{XQexp}E)] \\
=\ & (\lambda\ E'^{Env}.\ [\![e']\!]_{XQexp}E')\ (\lambda f^{Env \to Env}.\ f\ E) \\
& (\texttt{put}\ x\ [\![e]\!]_{XQexp}E) \\
=\ & \mathsf{let}\ v' = [\![e]\!]_{XQexp}E \\
& \qquad f = \lambda\ v'^{Xtf}.\ \texttt{put}\ x\ v' \\
& \mathsf{in}\ [\![e']\!](f\ v'\ E)
\end{aligned}
$$

The meaning of $\mathsf{let}\ x := e\,\mathsf{in}\ e'$ is that, after $e$ is evaluated in $E$ obtaining
a value $v'$, $\texttt{put}\ x\ v'$ changes environment $E$ to $E'$ by mapping $x$ to $v$,
and then expression $e'$ is evaluated in new environment $E'$.

**SEM EQ: Boolean Condition**

$$
\begin{aligned}
[\![\texttt{equal}\ e\ e'\ ]\!]_{Boolexp} E\ &\stackrel{\mathrm{def}}{=}\ \texttt{equal}\ ([\![e]\!]_{XQexp}E)\ ([\![e']\!]_{XQexp}E) \\
[\![\texttt{less}\ e\ e'\ ]\!]_{Boolexp} E\ &\stackrel{\mathrm{def}}{=}\ \texttt{less}\ ([\![e]\!]_{XQexp}E)\ ([\![e']\!]_{XQexp}E) \\
[\![\texttt{empty}\ e\ ]\!]_{Boolexp} E\ &\stackrel{\mathrm{def}}{=}\ \texttt{empty}\ ([\![e]\!]_{XQexp}E)
\end{aligned}
$$

The purpose of the boolean condition is to evaluate each *XQexp* expression in environment $E$, and then check the boolean relationship of results.

**SEM EQ: Where-Return**

$$[\![\texttt{where } \varphi \text{ return } e]\!]_{XQexp} E \quad \stackrel{\text{def}}{=} \quad \text{if } ([\![\varphi]\!]_{Boolexp} E) \text{ then } ([\![e]\!]_{XQexp} E)$$
$$\text{else nil}$$

where $\varphi$ is a boolean condition whose semantic equation is defined above. If condition $\varphi$ holds, *XQexp* expression $e$ is evaluated in environment $E$, and the result of type *Xtf* is returned.

**SEM EQ: For-Do**

$$[\![\text{for } x \in e \text{ do } e']\!]_{XQexp} E \quad \stackrel{\text{def}}{=} \quad [\![e']\!]_{XQexp}(E[x := v_1]) @ \ldots @$$
$$[\![e']\!]_{XQexp} E([x := v_k])$$
$$\text{where } [v_1, \ldots, v_k] = [\![e]\!] E$$

The semantic function of For-Do is complex. The operation for $x \in e$ do $e'$ first evaluates expression $e$ in environment $E$, returning a resulting list $[v_1, \ldots, v_k]$, and then '$x \in$" transfers environment $E$ to be a list of new environments $[E_1, \ldots, E_k]$. Finally expression $e'$ is evaluated on the list of the new environments.

In order to encode the For-Do semantic equation, a *mapping* function is employed which has

$$\texttt{mapf} : (\theta \to \theta') \to \theta \text{ list } \to \theta' \text{ list}$$

whose definition is

$$\texttt{mapf} \overset{\text{def}}{=} \lambda f^{\theta \to \theta'}. \lambda x^{\theta \text{ list}}.x \ [\theta' \text{ list}] \ \texttt{nil} \ (\lambda t^{\theta}. \lambda z^{\theta' \text{ list}}.\texttt{Cons} \ (f \ t) \ z)$$

which applies function $f$ to each element in list $x$.

For the For-Do clause, there are a set of state-transformation actions:

- the evaluation of expression $e$ in environment $E$ creates a *Xtf* list for $x$

$$newVals = [\![e]\!]E \ : \textit{Xtf} \ \text{list}$$

- $\texttt{mapPut} \ x \ newVals \ E$ creates a list of new environments $newEs$ by applying $\texttt{put} \ x \ e \ E$ to value $e$ in the $newVals$ list, where

$$
\begin{aligned}
\texttt{mapPut} \ &= \ \lambda x^{Var}.\texttt{mapf} \ (\texttt{put} \ x) \\
&= \ \lambda x^{Var}. \lambda newVals^{Xtf \ \text{list}}. \lambda E^{Env}. \ newVals \ [Env \ \text{list}] \\
&\qquad \texttt{nil} \ (\lambda e^{XQexp}. \lambda z^{Env \ \text{list}}. \texttt{Cons} \ (\texttt{put} \ x \ e \ E) \ z) \\
\texttt{mapPut} \ &: \quad \textit{Var} \to \textit{Xtf} \ \text{list} \to \textit{Env} \to \textit{Env} \ \text{list}
\end{aligned}
$$

- $\texttt{mapSem} \ e' \ E'$ evaluates $e'$ in each environment in $newEs$ list, and concatenates the results:

$$
\begin{aligned}
\texttt{mapSem} \ &= \ \lambda e'^{XQexp}. \texttt{mapf} \ ([\![e']\!]) \\
&= \ \lambda e'. \lambda newEs^{Env \ \text{list}}. \ newE \ [Xtf \ \text{list}] \ \texttt{nil} \\
&\qquad (\lambda E^{Env}. \lambda z^{\text{list} \ Xtf}. \texttt{Concat} \ ([\![e']\!]E) \ z) \\
\texttt{mapSem} \ &: \quad \textit{XQexp} \to \textit{Env} \ \text{list} \to \textit{Xtf} \ \text{list}
\end{aligned}
$$

With these functions, the semantic equation for `For-Do` is defined as

$\llbracket \text{for } x \ \in \ e \ \text{do } e' \rrbracket_{XQexp} E \overset{\text{def}}{=}$

     let $newVals = \llbracket e \rrbracket E$

         $\text{mapPut} = \lambda x^{Var}. \ \lambda newVals^{Xtf \ \text{list}}. \ \lambda E^{Env}. \ newVals \ [Env \ \text{list}]$

             $\text{nil } (\lambda e^{XQexp}. \ \lambda z^{Env \ \text{list}}. \ \text{Cons } (\text{put } x \ e \ E) \ z)$

           $\text{mapSem} = \lambda e'^{XQexp}. \ \lambda newEs^{Env \ \text{list}}. \ newE \ [Xtf \ \text{list}]$

             $\text{nil } (\lambda E^{Env}. \ \lambda z^{Xtf \ \text{list}}. \ \text{Concat } (\llbracket e' \rrbracket E) \ z)$

     in $\text{mapSem } e' \ (\text{mapPut } x \ newVals)$

**SEM EQ: $\text{XF}_\text{n}$**

$$\llbracket \text{XF}_\text{n}(e_1, \ldots, e_k) \rrbracket E \quad \overset{\text{def}}{=} \quad \text{XF}_\text{n} \ (\llbracket e_1 \rrbracket E, \ldots, \llbracket e_k \rrbracket E)$$

where $\text{XF}_\text{n}$ is an XML operator, such as `select`, `subtreesdfs`, with $k$ arguments. $[e_1, \ldots, e_k]$ is a expression list, $es$, where $e_i$ is an *XQexp* expression. The argument list, $es$, is first evaluated in environment $E$, and the resulting list $vs$ is supplied to $\text{XF}_\text{n}$ as arguments. First, `mapArg` is defined as

$\begin{aligned}
\text{mapArg} \ &= \ \text{mapf } (\lambda e^{XQexp}. \lambda E^{Env}. \llbracket e \rrbracket E) \\
&= \ (\lambda f^{\theta \to \theta'}. \lambda x^{\theta \ \text{list}}. x \ [\theta' \ \text{list}] \ \text{nil } (\lambda t^{\theta}. \lambda z^{\text{list} \ \theta'}. \text{Cons } (f \ t) \ z)) \\
&\quad \ (\lambda e^{XQexp}. \lambda E^{Env}. \llbracket e \rrbracket E) \\
&= \ \lambda argList^{XQexp \ \text{list}}. \ \lambda E^{Env}. argList \ [Xtf \ \text{list}] \\
&\qquad\qquad \text{nil } (\lambda e^{XQexp}. \lambda z^{Xtf \ \text{list}}. \text{Cons } (\llbracket e \rrbracket E) \ z) \\
\text{mapArg} \ &: \ \ XQexp \ \text{list} \to Env \to Xtf \ \text{list}
\end{aligned}$

where `mapArg` *es E* creates an *Xtf* list by applying $[\![e]\!]E$ to each *XQexp*
expression *e* in expression list *es*. Finally, the semantic meaning of $\mathtt{XF_n}$
is defined as

$$[\![\mathtt{XF_n}(e_1,\ldots,e_k)]\!]E \overset{\text{def}}{=}$$

$$\mathsf{let}\ argList = [e_1,\ldots,e_k]$$

$$\mathtt{mapArg} = \lambda\ argList^{XQexp\ \mathsf{list}}.\ \lambda\ E^{Env}.\ argList\ [Xtf\ \mathsf{list}]$$

$$\mathtt{nil}\ (\lambda\ e^{XQexp}.\ \lambda\ z^{Xtf\ \mathsf{list}}.\mathtt{Cons}\ ([\![e]\!]E)\ z)$$

$$\mathsf{in}\ \mathtt{XF_n}(\mathtt{mapArg}\ argList\ E)$$

So for, the XML data and all the basic operators defined in XQ are
represented successfully by *System F*. As a query language, XQuery has been
proved to be turing complete, which means its query evaluation termination
is undecidable. As an XQuery *fragment*, XQ has been shown to be *System
F representable*, which gives XQ a very important computability property.

**Theorem 4** Any XQ query evaluation terminates and the equivalence of
XQ queries is decidable.

This property is supported by the *strong normalization* property of *System F*
shown in Chapter 3. As an XML query language, XQ query's computation
decidability does not contradict the undecidability of the relational query
language because an XML document is ordered.

# Chapter 5

# Extension of XQ

The principal purpose of the dynamic interval method is to implement the XML query by a relation query engine, which ensures the polynomial time complexity in the query evaluation. After this is achieved, it is time to discuss the expressiveness of XQ language. Not only is its expressiveness investigate, but also its limitations are studied by examining simple queries. Although they cannot be defined in XQ language, they still can be expressed in relational SQL statements, leading us to the extension. In this chapter, first a simple query that is XQ undefinable is given, then an extension is presented to overcome this limitation, and retain XQ's *System F* encodable property to guarantee the *strong* query evaluation termination.

## 5.1   Limitation of XQ: A Case Study

As an XQuery *fragment*, XQ captures the core feature presented in XQuery/XPath expressions. XQ's simple syntax allows arbitrary compositions of basic func-

tion invocations, local variable definitions, conditional selections, and itera-tions over XML data. Although XQ is powerful, it still cannot express certain simple queries.

Figure 5.1 illustrates two XML queries, where query $Q_1$ asks to flat a multilevel list-like XML tree[1] $A$ into a two-level XML tree $B$, and query $Q_2$ conversely transforms XML tree $B$ back to tree $A$.
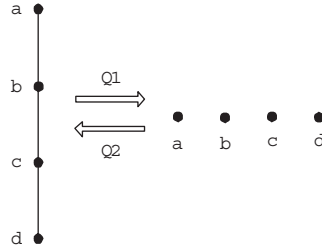


Figure 5.1: Two XML queries examples

The query $Q_1$ is represented in XQ syntax as

$$\texttt{roots (subtreesdfs } A)$$

where `subtreedfs` puts all the sub-self trees of $A$ into a forest in the docu-ment(tree) order, and the `roots` extracts the root nodes for all the sub-trees. The purpose of function composition `roots ∘ substreedfs` is to flat a com-plex multilevel hierarchical tree into a forest with each node at the same top level. In addition, the XQ expression of query $Q_1$ can be translated into an SQL statement with the XQ-to-SQL templates of `subtreedfs`, `roots`, defined in [6], and evaluated by a relational query engine. Now query $Q_2$ picks up all the nodes of tree $B$ and constructs a deep tree without *sibling* relationships. However, $Q_2$ cannot be expressed by *XQ*, even though the

---

[1]Tree and forest are used to describe the XML document in this chapter.

| s | l | r |
|---|---|---|
| a | 0 | 1 |
| b | 2 | 3 |
| c | 4 | 5 |
| d | 6 | 7 |

| s | l | r |
|---|---|---|
| a | 0 | 14 |
| b | 2 | 12 |
| c | 4 | 10 |
| d | 6 | 8 |

Figure 5.2: Table $T_B$          Figure 5.3: $T_{verticalTree}$ from Tree $B$

query still can be expressed by the SQL statement after tree $B$ is *interval encoded* into a relation.

### 5.1.1 An XQ Undefinable Query in SQL

The implementation of query $Q_2$ by the relational method includes two steps:

- encoding XML tree $B$ into an interval relation by performing a depth first traversal on the tree $B$ and marking each node's left and right end points, as depicted in Figure 5.2;

- constructing a tree by linking all the nodes of tree $B$ using the parent-child relationship with the document order, shown in Figure 5.4.

### 5.1.2 Analysis of the XQ Limitations

The reason of $Q_2$ is definable in relational SQL queries is that, when tree $B$ is encoded into relation, the node tuples can be accessed arbitrarily, and their parent-children, ancestor-descendent, and sibling relationship can be reassembled by changing the value of interval attributes $l$ and $r$. For tree $B$ $[a, b, c, d]$ (expressed simply), by using a SQL statement, query $Q_2$ can be implemented by extending the right interval value of $c$ to contain the right

```
SELECT u.s AS s,  u.l AS l,  2 * v.r − u.l AS r
FROM T_B u,
        (SELECT v.s,  v.l,  v.r
         FROM T_B v
         WHERE NOT EXISTS(
         SELECT *
         FROM T_B m
         WHERE m.l > v.r) v
```

Figure 5.4: SQL Expression for Constructing a Vertical Tree

interval of node $d$ to obtain a new tree with $c$ being the parent of $d$, and so on. With the interval relation, the *vertical* tree is successfully constructed without sibling relationship.

After the XQ *undefinable* query is expressed in SQL statements, we would discuss the limitations of XQ language, that is, why can XQ not express such a query as $Q_2$? With analyzing the syntax of XQ in Chapter 2, we find that only a *horizontal* tree constructor `Xnode` is defined, which creates a new tree using string $s$ and forest $ts$, denoted as `<s> ts </s>`. The `Xnode` expresses only the tree construction on the top level, with wrapping up a list of of XML elements by tag $s$.

**Fact 1** Only with the tree constructor `Xnode`, XQ cannot express the nesting queries whose results are much deeper than the original document, except enumerating each individual element.

In the thesis, the objective is to extend XQ by adding a new tree operator

that constructs a tree in *vertical* direction.

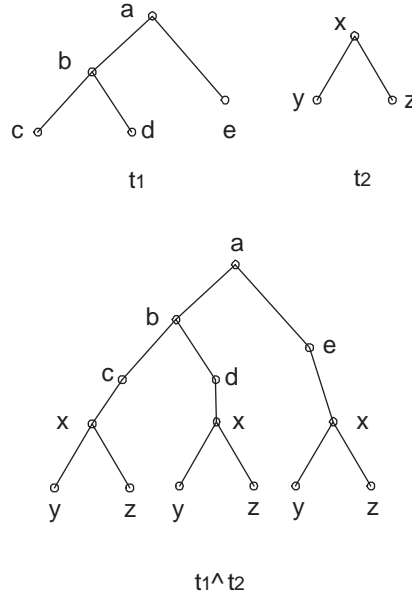## 5.2   New Tree Constructor Xtree



Figure 5.5: The operation of Xtree $t_1$ $t_2$

An example of a vertical tree construction is offered in Figure 5.5. The new tree operator is named Xtree, where Xtree $t_1$ $t_2$ appends tree $t_2$ to each leaf node of tree $t_1$. After extending XQ syntax by Xtree, Xtree will be represented in *System F* to ensure this extension retains the strong normalization property of XQ.

### 5.2.1 Representing Xtree in *System F*

As a vertical tree operation, Xtree $t_1$ $t_2$ constructs a new tree by appending $t_2$ to each leaf node of tree $t_1$ as

$$\text{Xtree (Xnode } s \text{ } ts) \text{ } t_2 \text{ } = \text{ if empty } ts \text{ then Xnode } s \text{ } t_2$$

$$\text{else Xnode } s \text{ (Xtree } ts \text{ } t_2) \qquad (5.1)$$

$$\text{Xtree nil } t_2 \text{ } = \text{ nil} \qquad (5.2)$$

$$\text{Xtree (Cons } t \text{ } ts) \text{ } t_2 \text{ } = \text{ Cons (Xtree } t \text{ } t_2) \text{ (Xtree } ts \text{ } t_2) \qquad (5.3)$$

If $t_1$ is a tree as Xnode $s$ $ts$, and if $ts$ is [ ], then $t_2$ is directly appended as the children of root $s$ as <s> $t_2$ </s>; otherwise, the iterator visits $t_1$ children's level, which is XML forest $ts$, and applies operator Xtree to each tree element of $ts$.

Equation (5.1) does not match (4.9) because of $ts$ is employed as an argument of empty. It is necessary to define a pair function $g$ to store the input argument and the corresponding result of Xtree as follows:

$$g \text{ } t_1 \text{ } t_2 \stackrel{\text{def}}{=} (t_1, \text{ Xtree } t_1 \text{ } t_2)$$

with the following details:

$$g \text{ (Xnode } s \text{ } ts) \text{ } t_2 \text{ } = \text{ (Xnode } s \text{ } ts, \text{ if empty } ts \text{ then Xnode } s \text{ } t_2$$

$$\text{else Xnode } s \text{ (Xtree } ts \text{ } t_2)) \qquad (5.4)$$

$$g \text{ nil } t_2 \text{ } = \text{ (nil, nil)}$$

$$g \text{ (Cons } t \text{ } ts) \text{ } t_2 \text{ } = \text{ (Cons } t \text{ } ts, \text{Cons (Xtree } t \text{ } t_2) \text{ (Xtree } ts \text{ } t_2))$$

In order to match the form of (4.9), with $g$ $t_1$ $t_2 = (t_1, \text{ Xtree } t_1 \text{ } t_2)$, (5.4) is

rewritten as

$$g \ (\texttt{Xnode} \ s \ ts) \ t_2 \ = \ (\lambda \, s^{string}.\lambda(k, \ z)^{Xtf \times Xtf}.(\texttt{Xnode} \ s \ k, \ \texttt{if empty} \ k$$

$$\texttt{then Xnode} \ s \ t_2 \ \texttt{else Xnode} \ s \ z)) \ s \ (ts, \ \texttt{Xtree} \ ts \ t_2)$$

$$= \ \texttt{let} \ f_T = \lambda \, s^{String}.\lambda(k, \ z)^{Xtf \times Xtf}.(\texttt{Xnode} \ s \ k,$$

$$\texttt{if empty} \ k \ \texttt{then Xnode} \ s \ t_2 \ \texttt{else Xnode} \ s \ z)$$

$$\texttt{in} \ f_T \ s \ (g \ ts)$$

which matches the form of Equation (4.9). Function $g$ is defined as

$$g \ \overset{\text{def}}{=} \ \lambda \, t_1^{Xtf}. \ t_1 \ [Xtf \times Xtf]$$

$$\texttt{let} \ f_T = \lambda \, t_2^{Xtf}.\lambda \, s^{String}. \ \lambda(k, \ z)^{Xtf \times Xtf}. \ (\texttt{Xnode} \ s \ k, \ \texttt{if empty} \ k$$

$$\texttt{then Xnode} \ s \ t_2 \ \texttt{else Xnode} \ s \ z)$$

$$f_{XF} = \lambda \, t_2^{Xtf}. \ \lambda(k, \ z)^{Xtf \times Xtf}.(\texttt{Cons} \ t \ k, \ \texttt{Cons} \ (\texttt{Xtree} \ k \ t_2)$$

$$(\texttt{Xtree} \ z \ t_2))$$

$$\texttt{in} \ f_T \ (\lambda \, t_2^{Xtf}.(\texttt{nil, nil})) \ f_{XF}$$

and operator $\texttt{Xtree}$ is defined as follows:

$$\texttt{Xtree} \ \overset{\text{def}}{=} \ g.1$$

$$= \ (\lambda \, t_1^{Xtf}. \ t_1 \ [Xtf \times Xtf]$$

$$\texttt{let} \ f_T = \lambda \, t_2^{Xtf}.\lambda \, s^{String}.\lambda(k, \ z)^{Xtf \times Xtf}. \ (\texttt{Xnode} \ s \ k,$$

$$\texttt{if empty} \ k \ \texttt{then Xnode} \ s \ t_2 \ \texttt{else Xnode} \ s \ z)$$

$$f_{XF} = \lambda \, t_2^{Xtf}.\lambda(k, \ z)^{Xtf \times Xtf}.(\texttt{Cons} \ k \ z, \ \texttt{Cons} \ (\texttt{Xtree} \ k \ t_2)$$

$$(\texttt{Xtree} \ z \ t_2))$$

$$\texttt{in} \ f_T \ (\lambda \, t_2^{Xtf}.(\texttt{nil, nil})) \ f_{XF}).1$$

```
1        CREATE VIEW Xtree(T₁, T₂) AS
2          (SELECT u.s AS s, u.l * w_T₂ AS l, u.r * w_T₂ AS r
3           FROM T₁ u)
4         UNION ALL
5          (SELECT m.s AS s, m.l + u.l * w_T₂ AS l, m.r + u.l * w_T₂ AS r
6           FROM
7             (SELECT u.s, u.l, u.r
8              FROM T₁ u
9              WHERE NOT EXISTS
10                (SELECT *
11                 FROM T₁ v
12                 WHERE u.l < v.l AND u.r > v.r),
13             (SELECT m.s, m.l, m.r
14              FROM T₂ m
15              WHERE EXISTS
16                (SELECT *
17                 FROM T₂ p
18                 WHERE p.l < m.l AND p.r > m.r)
```

Figure 5.6: The relational template (View) for `Xtree`

## 5.2.2  SQL Template for `Xtree`

After the new operator `Xtree` has been successfully represented in *System F*, indicating that the extended XQ still exhibits the strong normalization property, this thesis attempts to show that a corresponding SQL template exists for `Xtree`, ensuring the polynomial time complexity of XQ query evaluation. The `Xtree` template is signified in Figure 5.6.

First, the interval values of each node of tree $t_1$, encoded in table $T_1$, are extended by the width of tree $t_2$ (Line 2, Line 3) to confirm that each node has enough interval space to contain tree $t_2$ as its children, and then each leaf nodes of tree $t_1$ has been picked up (Line 7 to Line 12) and attached to

the children nodes of tree $t_2$ (Line 13 to Line 18) by a cross-product (Line 6 to Line 18), where the intervals of children nodes of $t_2$ have been adjust to fit the interval of each leaf node (Line 5). Finally, the tree nodes of $t_1$ with the extended interval and appended $t_2$ nodes, are returned as the result of view `Xtree`.

## 5.3   New Vertical **Vfor** Clause

XQ syntax defines FLWR-like clauses, where the for clause for $x \in e$ do $e'$ has the following semantic meaning:

$$[\![\text{for } x \ \in \ e \text{ do } e']\!]E \ \overset{\text{def}}{=} \ [\![e']\!](E[x := v_1]) \ @\ldots@ \ [\![e']\!]E([x := v_k])$$

$$\text{where } [v_1, \ \ldots, \ v_k] = [\![e]\!]E$$

The expression for $x \ \in e$ do $e'$ first evaluates expression $e$ in *Environment E*, and obtains a list of result $[v_1, \ \ldots, \ v_k]$, and then the assignment operation "$x \in$" translates *Environment E* to a list of new *Environments* $[E_1, \ \ldots, \ E_k]$. Finally, the expression $e'$ is evaluated in the list of new environments, and the results are concatenated. The concatenating operation @ has been defined as an *infix* operator in [6], and expressed as a *prefix* operator `concat` in this thesis by

$$\text{concat } ts_1 \ ts_2 \equiv ts_1 \ @ \ ts_2$$

where $ts_1$ and $ts_2$ are XML forests.

With the concatenating operator @, the semantics of the for clause is easy to express. The idea of evaluating a for XQ expression with concatenating a

sequence of individual result is used and extended to represent the semantics of vertical construction. A new vertical Vfor clause is proposed as

$$\textsf{Vfor } x \in e \textsf{ do } e'$$

and its semantic equation is expressed, with the infix form $\wedge$ of `Xtree`, as

$$[\![\textsf{Vfor } x \in e \textsf{ do } e']\!]E \overset{\text{def}}{=} [\![e']\!](E[x := v_1]) \wedge \cdots \wedge [\![e']\!]E([x := v_k])$$

$$\text{where } [v_1, \ldots, v_k] = [\![e]\!]E \tag{5.5}$$

where the expression $e'$ is first evaluated in *Environment E*, returning a list of result $[v_1, \ldots, v_k]$, and then the assignment operation "$x \in$" translates *Environment E* into a sequence of new *Environments* $[E_1, \ldots, E_k]$. Finally, the expression $e'$ is evaluated in these new environments, and the sequence of results are combined by the vertical tree operation "$\wedge$", where

$$\texttt{Xtree } t_1 \ t_2 \equiv t_1 \wedge t_2$$

with $ts_1$ and $ts_2$ being XML trees.

### 5.3.1 Expressing the Semantics of **Vfor** in *System F*

With the extension of `Xtree` operator and Vfor clause, XQ has become more powerful to express the queries that are undefinable in the original XQ syntax. The focus is still on XQ's computational properties in the attempt to represent the vertical "for" clause Vfor in *System F*, which ensures that the extended XQ language still holds the strong normalization property.

The semantic equation (5.5) of Vfor is similar to that of For-Do. Also, it is necessary to use the general list *mapping* function `mapf`, defined in

Section 4.4.3, which applies the function $f$ to each element in the list $ts$ with type

$$\texttt{mapf} : (\theta \rightarrow \theta') \rightarrow \theta \text{ list } \rightarrow \theta' \text{ list}$$

First, the expression $\mathsf{Vfor}\ x\ \in e\ \mathsf{do}\ e'$ evaluates expression $e$ in environment $E$ and obtains a list of $Xtf$ results, $[v_1,\ \ldots,\ v_k]$, defined as an $Xtf$ list, $newVals$, such that

$$newVals = [\![e]\!]E\ :\ Xtf \text{ list}$$

Next, the assignment operation "$x \in$" transfers the given environment $E$ into a list of new environments $[E_1,\ \ldots,\ E_k]$ by inserting each $x$ value into a different environment. This transformation is implemented by a *mapping* function $\texttt{mapPut}$, where $\texttt{mapPut}\ x\ newVals\ E$ creates a list of new environments, $newEs$. This is achieved by utilizing $\texttt{put}\ x\ e_j\ E$, defined in Section 4.4.3, to insert each value, $e_j,\ j = 1 \ldots k$, in the $newVals$ list into environment $E$:

$$
\begin{aligned}
\texttt{mapPut}\ &=\ \lambda\, x^{Var}.\ \texttt{mapf}\ (\texttt{put}\ x)\\
&=\ \lambda\, x^{Var}.\ \lambda\, newVals^{Xtf \text{ list}}.\ \lambda\, E^{Env}.\ newVals\ [Env \text{ list}]\ \texttt{nil}\\
&\qquad (\lambda\, e^{XQexp}.\ \lambda\, z^{Env \text{ list}}.\ \texttt{Cons}\ (\texttt{put}\ x\ e\ E)\ z)\\
\texttt{mapPut}\ &:\ \ Var \rightarrow Xtf \text{ list} \rightarrow Env \rightarrow Env \text{ list}
\end{aligned}
$$

After expression $e'$ is evaluated in the list of environments, $newEs$, the results are combined by the vertical tree operator $\texttt{Xtree}$, and expressed as

$$[\![e']\!]E_1\ \wedge \cdots \wedge\ [\![e']\!]E_k$$

The function, $\texttt{semVfor}$, expresses the evaluation of $e'$ in the environment list,

*newEs*:

$$\texttt{semVfor} \; [\![e']\!] \; \texttt{nil} \;\; = \;\; \texttt{nil}$$

$$\texttt{semVfor} \; [\![e']\!] \; (\texttt{Cons} \; E \; Es) \;\; = \;\; \texttt{Xtree} \; ([\![e']\!]E) \; (\texttt{semVfor} \; [\![e']\!] \; Es)$$

and `semVfor` is defined as

$$\texttt{semVfor} \;\; = \;\; \lambda([\![e']\!])^{Env \to Xtf}. \, \lambda \, newEs^{Env \; \mathsf{list}}. \, newEs \, [Xtf]$$

$$\texttt{nil} \; (\lambda \, E^{Env}. \; \lambda \, z^{Xtf}. \, \texttt{Xtree} \; ([\![e']\!]E) \; z)$$

$$\texttt{semVfor} \;\; : \;\; (XQexp \to Env \to Xtf) \to Env \; \mathsf{list} \to Xtf$$

Finally, the semantic equation for `Vfor` can be defined as

$$[\![\mathsf{Vfor} \; x \; \in \; e \; \mathsf{do} \; e']\!]_{XQexp} E \; \overset{\mathrm{def}}{=}$$

$$\mathsf{let} \; newVals = [\![e]\!]E$$

$$\texttt{mapPut} = \lambda \, x^{Var}. \; \lambda \, newVals^{Xtf \; \mathsf{list}}. \; \lambda \, E^{Env}. \, newVals \, [Env \; \mathsf{list}]$$

$$\texttt{nil} \; (\lambda \, e^{XQexp}. \; \lambda \, z^{Env \; \mathsf{list}}. \, \texttt{Cons} \; (\texttt{put} \; x \; e \; E) \; z)$$

$$\texttt{semVfor} = \lambda([\![e']\!])^{Env \to Xtf}. \, \lambda \, newEs^{Env \; \mathsf{list}}. \, newEs \, [Xtf]$$

$$\texttt{nil} \; (\lambda \, E^{Env}. \; \lambda \, z^{Xtf}. \, \texttt{Xtree} \; ([\![e']\!]E) \; z)$$

$$\mathsf{in} \; \texttt{semVfor} \; [\![e']\!] \; (\texttt{mapPut} \; x \; newVals)$$

## 5.3.2 Expressing the XQ Undefinable Query

With the syntax of the vertical **Vfor** clause and corresponding semantic meaning, the extended XQ can easily express query $Q_2$ in Figure 5.1, which transforms forest $B$ into a *linked-list* deep tree that has no sibling relationships:

$$\mathsf{Vfor} \; \$1 = B \; \mathsf{do} \; \$1$$

where variable $1 is bound to each tree element of $B$, which are trees $a$, $b$, $c$, and $d$ (here the trees are simply expressed by their root labels), as seen in Figure 5.7(a); the Vfor clause utilizes operator Xtree to append all trees bound by 1 in a vertical direction to a deeply linked tree, shown in Figure 5.7(b).

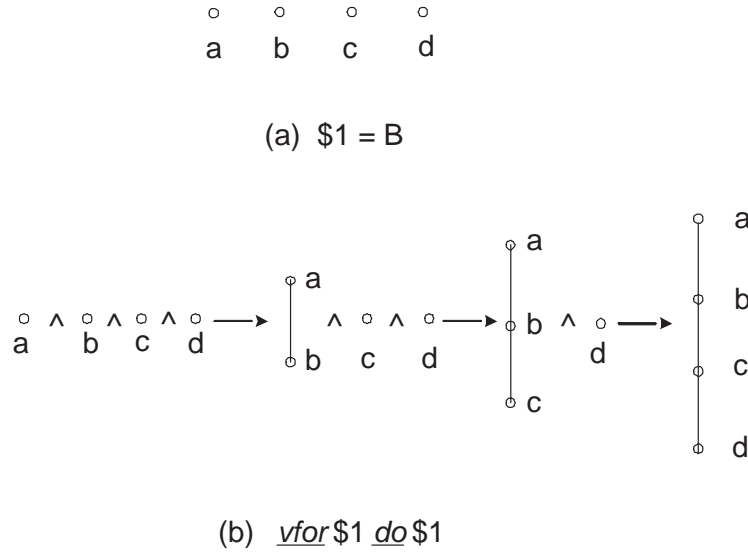The relational implementation of the Vfor $x \in e$ do $e'$ is similar to the imple-



(a)  $1 = B



(b)  *vfor* $1 *do* $1

Figure 5.7: The illustration of query $Q_2$ with the Vfor clause

mentation of for $x \in e$ do $e'$, except in the last step, the for-do horizontally concatenates, with operator @, of all the results evaluated in the sequence of environments created by $x \in e$; the Vfor applies the vertical tree opera-tor, Xtree, to append the results evaluated in the sequence of environments vertically.

With the new tree operator Xtree and the new Vfor clause, XQ has been successfully extended to be more expressive. Also, the extended XQ has been

shown to be *System F* definable and XQ-to-SQL translatable, which ensure this extension keeps the polynomial time complexity of query evaluation and *strong* termination property.

# Chapter 6

# Conclusion and Future Work

## 6.1 Conclusion

In this thesis *System F* is shown to be able to express queries in the core of XQ, an XQuery fragment proposed in [6]. As a second order typed lambda calculus, *System F* holds a proven expressive power [7] and a strong normalization property. Therefore, the XML data is represented in *System F* as inductive data types, XML tree and XML forest, with the corresponding iterators. Also, all the basic XML operators of XQ are encoded into *System F*. The semantics of XQ are expressed in *System F* by encoding the semantic environment using an *Environment* data with a set of operations. The successful encoding of XQ in *System F* guarantees the termination of XQ query evaluation and the decidability of query equivalence.

After XQ is represented in *System F*, an extension of XQ is investigated. Due to the limitations of XQ, a new tree operator, XTree, and a vertical Vfor clause are proposed to express some of the *undefinable* XQ queries. It

is demonstrated that this extension still retains the XQ-to-SQL translation property to achieve polynomial evaluation time complexity by utilizing the relation query engine, and also XQ's *System F* encodable property to ensure the termination of XQ query evaluation.

## 6.2   Future Work

The suggestions for future work include the followings:

- to further extend XQ towards its completeness with keeping its XQ-to-SQL translatable and *System F* encodable properties, along with proof of the completeness of the extended XQ;

- an XQ-to-SF parser to translate XQ queries into *System F* syntax terms; also, an XQ *System F* processor should be implemented by a functional language to perform XQ query normalization and detect the query equivalence;

- to evaluate *System F* normalized XQ queries with a relational evaluation engine, an SF-to-SQL translator should be implemented to translate the *System F* syntax fragments to SQL statements.

# Bibliography

[1] The XML Query Use Cases, Available from http://www.w3.org/TR/2005/WD-xquery-use-cases-20050404/.

[2] H.P. Barendregt. *The Lambda Calculus, Its Syntax and Semantics*, volume 103 of *Studies in Logics and the Foundations of Mathmatics*. North-Holland, 1981.

[3] H.P. Barendregt. *The Lambda Calculus with Types*. Handbook of Logic in Computer Science. Oxford University Press, 1990.

[4] M. Fernández B.Choi and J. Siméon. The XQuery Formal Semantics: A Foundation for Implementation and Optimization. Technical report, 2002.

[5] M. Fernández A. Malhotra K. Rose M. Rys D. Draper, P. Fankhauser and J. Siméon. XQuery 1.0 and XPath 2.0 Formal Semantics, Available from http://www.w3.org/TR/query-semantics/, 2002.

[6] D. DeHaan, D. Toman, M. Consens, and M. Ozsu. A Comprehensive XQuery to SQL Translation Using Dynamic Interval Coding, 2003.

[7] Y. Lafont J.-Y.Girard and P.Taylor. *Proof Theory and Types.* Cambridge Tracts in Theorectical Computer Science 7. Cambridge University Press, 1989.

[8] H. V. Jagadish, Laks V. S. Lakshmanan, Divesh Srivastava, and Keith Thompson. Tax: A tree algebra for xml. In *DBPL '01: Revised Papers from the 8th International Workshop on Database Programming Languages*, pages 149–164, London, UK, 2002. Springer-Verlag.

[9] Howard Katz. *XQuery from the Experts.* Addison-Wesley, 2004.

[10] J. C. Reynolds. Polymorphic Lambda-Calculus: Introduction. In *Logical Foundations of Functional Programming*, pages 77–86. Addison-Wesley, Reading, MA, 1990.

[11] J.C. Reynolds. *Theories of Programming Languages.* Cambridge University Press, 1998.