

A semi-monad for semi-structured data

Mary Fernandez¹, Jérôme Siméon², and Philip Wadler³

¹ AT&T Labs, mff@research.att.com

² Bell Labs, Lucent Technologies simeon@research.bell-labs.com

³ Avaya Labs, wadler@avaya.com

Abstract. This document proposes an algebra for XML Query. The algebra has been submitted to the W3C XML Query Working Group. A novel feature of the algebra is the use of regular-expression types, similar in power to DTDs or XML Schemas, and closely related to Hasoya and Pierce's work on Xduce. The iteration construct is based on the notion of a monad, and involves novel typing rules not encountered elsewhere.

1 Introduction

This document proposes an algebra for XML Query.

This work builds on long standing traditions in the database community. In particular, we have been inspired by systems such as SQL, OQL, and nested relational algebra (NRA). We have also been inspired by systems such as Quilt, UnQL, XDuce, XML-QL, XPath, XQL, XSLT, and YATL. We give citations for all these systems below.

In the database world, it is common to translate a query language into an algebra; this happens in SQL, OQL, and NRA, among others. The purpose of the algebra is twofold. First, the algebra is used to give a semantics for the query language, so the operations of the algebra should be well-defined. Second, the algebra is used to support query optimization, so the algebra should possess a rich set of laws. Our algebra is powerful enough to capture the semantics of many XML query languages, and the laws we give include analogues of most of the laws of relational algebra.

In the database world, it is common for a query language to exploit schemas or types; this happens in SQL, OQL, and NRA, among others. The purpose of types is twofold. Types can be used to detect certain kinds of errors at compile time and to support query optimization. DTDs and XML Schemas can be thought of as providing something like types for XML. Our algebra uses a simple type system that captures the essence of XML Schema. The type system is close to that used in XDuce, which in turn is based on the well-known idea of tree automata. Our type system can detect common type errors and support optimization. A novel aspect of the type system (not found in Xduce) is the description of projection in terms of iteration, and the typing rules for iteration that make this viable. The algebra is based on earlier work on the use of monads to query semi-structured data, and iteration construct satisfies the three monad laws.

To better familiarize readers with the algebra, we have implemented a type checker and an interpreter for the algebra in OCaml [8]. A demonstration version of the system is available via the authors's home pages.

This paper describes the key features of the algebra. For simplicity, we restrict our attention to only three scalar types (strings, integers, and booleans), but we believe the system will smoothly extend to cover the continuum of scalar types found in XML Schema. Other important features that we do not tackle include attributes, namespaces, element identity, collation, and key constraints, among others. Again, we believe they can be added within the framework given here.

Two earlier versions of this paper have been distributed. The first [18] used a more ad-hoc approach to typing. The second [19] used a different notation for case analysis, and has less discussion of the relation to earlier work on monads.

The paper is organized as follows. A tutorial introduction is presented in Section 2. Section 3 explains key aspects of how the algebra treats projection and iteration, and the relation to monads. The expressions of the algebra are summarized in Section 4. The type system is reviewed in Section 5. Some laws of the algebra are presented in Section 6. Finally, the static typing rules for the algebra are described in Section 7. Section 8 discusses open issues and problems.

Cited literature includes: alternative algebras for XML [3, 11], comprehensions [32, 6], monads [26, 27, 32–34], NRA and equivalents [14, 7, 22, 23, 29, 30], OQL [2, 1, 13, 9], Quilt [10], UnQL [5], SQL [16], tree Automata [15, 31], type systems [24, 25], XDuce [20, 21], XML Query [38, 39], XML Schema [40, 41], XML-QL [17], XPath [37, 35], XQL [28], XSLT [42], and YaTL [12].

2 The Algebra by Example

This section introduces the main features of the algebra, using familiar examples based on accessing a database of books.

2.1 Data and Types

Consider the following sample data:

```
<bib>
  <book>
    <title>Data on the Web</title>
    <year>1999</year>
    <author>Abiteboul</author>
    <author>Buneman</author>
    <author>Suciu</author>
  </book>
  <book>
    <title>XML Query</title>
    <year>2001</year>
    <author>Fernandez</author>
```

```

    <author>Suciu</author>
  </book>
</bib>

```

Here is a fragment of a XML Schema for such data.

```

<xsd:group name="Bib">
  <xsd:element name="bib">
    <xsd:complexType>
      <xsd:group ref="Book"
        minOccurs="0" maxOccurs="unbounded"/>
    </xsd:complexType>
  </xsd:element>
</xsd:group>

<xsd:group name="Book">
  <xsd:element name="book">
    <xsd:complexType>
      <xsd:element name="title" type="xsd:string"/>
      <xsd:element name="year" type="xsd:integer"/>
      <xsd:element name="author" type="xsd:integer"
        minOccurs="1" maxOccurs="unbounded"/>
    </xsd:complexType>
  </xsd:element>
</xsd:group>

```

This data and schema is represented in our algebra as follows:

```

type Bib =
  bib [ Book* ]
type Book =
  book [
    title [ String ],
    year [ Integer ],
    author [ String ]+
  ]
let bib0 : Bib =
  bib [
    book [
      title [ "Data on the Web" ],
      year [ 1999 ],
      author [ "Abiteboul" ],
      author [ "Buneman" ],
      author [ "Suciu" ]
    ],
    book [
      title [ "XML Query" ],

```

```

        year   [ 2001 ],
        author [ "Fernandez" ],
        author [ "Suciu" ]
    ]
]

```

The expression above defines two types, `Bib` and `Book`, and defines one global variable, `bib0`.

The `Bib` type consists of a `bib` element containing zero or more entries of type `Book`. The `Book` type consists of a `book` element containing a `title` element (which contains a string), a `year` element (which contains an integer), and one or more `author` elements (which contain strings).

The `Bib` type corresponds to a single `bib` element, which contains a *forest* of zero or more `Book` elements. We use the term *forest* to refer to an ordered sequence of (zero or more) elements. Every element can be viewed as a forest of length one.

The `Book` type corresponds to a single `book` element, which contains one `title` element, followed by one `year` element, followed by one or more `author` elements. A `title` or `author` element contains a string value and a `year` element contains an integer.

The variable `bib0` is bound to a literal XML value, which is the data model representation of the earlier XML document. The `bib` element contains two `book` elements.

The algebra is a strongly typed language, therefore the value of `bib0` must be an instance of its declared type, or the expression is ill-typed. Here the value of `bib0` is an instance of the `Bib` type, because it contains one `bib` element, which contains two `book` elements, each of which contain a string-valued `title`, an integer-valued `year`, and one or more string-valued `author` elements.

For convenience, we define a second global variable `book0`, also bound to a literal value, which is equivalent to the first book in `bib0`.

```

let book0 : Book =
  book [
    title [ "Data on the Web" ],
    year   [ 1999 ],
    author [ "Abiteboul" ],
    author [ "Buneman" ],
    author [ "Suciu" ]
  ]

```

The notation used for types is based on that used in Xduce [20], which in turn is based on tree automata [15, 31]. Linear automata come in two varieties, deterministic or non-deterministic, and these are well known to have equivalent expressive power. In contrast, tree automata come in four varieties, they may be deterministic or non-deterministic, and they may be top-down or bottom-up. Of these four varieties, three have equivalent expressive power, the odd man out being top-down deterministic tree automata, which are strictly less powerful.

While the Xduce notation expresses the full power of non-deterministic or bottom-up tree automata, both XML Schema and the type system presented here are restricted to the power of top-down and deterministic tree automata. This restriction is expressed by two requirements in XML Schema, which for compatibility are also adopted by XML Query: any sibling elements with the same name must have the same content, and all regular expressions must be one-unambiguous [4]. Adopting these restrictions makes it easy to decide when one type is smaller than another (inclusion of languages, later written as $<:$). For deterministic and top-down automata this can be determined in time quadratic in the number of states, whereas for non-deterministic or bottom-up automata it may require time exponential in the number of states.

2.2 Projection

The simplest operation is projection. The algebra uses a notation similar in appearance and meaning to path navigation in XPath.

The following expression returns all `author` elements contained in `book0`:

```

book0/author
==> author [ "Abiteboul" ],
      author [ "Buneman" ],
      author [ "Suciu" ]
:   author [ String ]+

```

The above example and the ones that follow have three parts. First is an expression in the algebra. Second, following the `==>`, is the value of this expression. Third, following the `:`, is the type of the expression, which is (of course) also a legal type for the value.

The following expression returns all `author` elements contained in `book` elements contained in `bib0`:

```

bib0/book/author
==> author [ "Abiteboul" ],
      author [ "Buneman" ],
      author [ "Suciu" ],
      author [ "Fernandez" ],
      author [ "Suciu" ]
:   author [ String ]*

```

Note that in the result, the document order of `author` elements is preserved and that duplicate elements are also preserved.

It may be unclear why the type of `bib0/book/author` contains *zero* or more authors, even though the type of a `book` element contains *one* or more authors. Let's look at the derivation of the result type by looking at the type of each sub-expression:

```

bib0           : Bib
bib0/book      : Book*
bib0/book/author : author [ String ]*

```

Recall that `Bib`, the type of `bib0`, may contain *zero* or more `Book` elements, therefore the expression `bib0/book` might contain zero `book` elements, in which case, `bib0/book/author` would contain no authors.

This illustrates an important feature of the type system: the type of an expression depends only on the type of its sub-expressions. It also illustrates the difference between an expression's run-time value and its compile-time type. Since the type of `bib0` is `Bib`, the best type for `bib0/book/author` is one listing zero or more authors, even though for the given value of `bib0` the expression will always contain exactly five authors.

One may access scalar data (strings, integers, or booleans) using the keyword `data()`. For instance, if we wish to select all author names in a book, rather than all author elements, we could write the following.

```
book0/author/data()
==> "Abiteboul",
    "Buneman",
    "Suciu"
:   String+
```

Similarly, the following returns the year the book was published.

```
book0/year/data()
==> 1999
:   Integer
```

This notation is similar to the use of `text()` in XPath. We chose the keyword `data()` because, as the second example shows, not all data items are strings.

2.3 Iteration

Another common operation is to iterate over elements in a document so that their content can be transformed into new content. Here is an example of how to process each book to list the authors before the title, and remove the year.

```
for b in bib0/book do
  book [ b/author, b/title ]
==> book [
  author [ "Abiteboul" ],
  author [ "Buneman" ],
  author [ "Suciu" ],
  title [ "Data on the Web" ]
],
book [
  author [ "Fernandez" ],
  author [ "Suciu" ],
  title [ "XML Query" ]
]
```

```

:   book [
      author[ String ]+,
      title[ String ]
    ]*

```

The `for` expression iterates over all `book` elements in `bib0` and binds the variable `b` to each such element. For each element bound to `b`, the inner expression constructs a new `book` element containing the book's authors followed by its title. The transformed elements appear in the same order as they occur in `bib0`.

In the result type, a `book` element is guaranteed to contain one or more authors followed by one title. Let's look at the derivation of the result type to see why:

```

bib0/book      : Book*
b              : Book
b/author       : author [ String ]+
b/title        : title  [ String ]

```

The type system can determine that `b` is always `Book`, therefore the type of `b/author` is `author[String]+` and the type of `b/title` is `title[String]`.

In general, the value of a `for` loop is a forest. If the body of the loop itself yields a forest, then all of the forests are concatenated together. For instance, the expression:

```

for b in bib0/book do
  b/author

```

is exactly equivalent to the expression `bib0/book/author`.

Here we have explained the typing of `for` loops by example. In fact, the typing rules are rather subtle, and one of the more interesting aspects of the algebra, and will be explained further below.

2.4 Selection

Projection and `for` loops can serve as the basis for many interesting queries. The next three sections show how they provide the power for selection, quantification, join, and regrouping.

To select values that satisfy some predicate, we use the `where` expression. For example, the following expression selects all `book` elements in `bib0` that were published before 2000.

```

for b in bib0/book do
  where b/year/data() <= 2000 do
    b
==> book [
  title [ "Data on the Web" ],
  year  [ 1999 ],
  author [ "Abiteboul" ],

```

```

        author [ "Buneman" ],
        author [ "Suciu" ]
    ]
:   Book*

```

An expression of the form

```
where  $e_1$  do  $e_2$ 
```

is just syntactic sugar for

```
if  $e_1$  then  $e_2$  else ()
```

where e_1 and e_2 are expressions. Here $()$ is an expression that stands for the empty sequence, a forest that contains no elements. We also write $()$ for the type of the empty sequence.

According to this rule, the expression above translates to

```
for b <- bib0/book in
  if b/year/data() < 2000 then b else ()
```

and this has the same value and the same type as the preceding expression.

2.5 Quantification

The following expression selects all book elements in `bib0` that have *some* author named "Buneman".

```

    for b in bib0/book do
      for a in b/author do
        where a/data() = "Buneman" do
          b
  ==> book [
    title [ "Data on the Web" ],
    year  [ 1999 ],
    author [ "Abiteboul" ],
    author [ "Buneman" ],
    author [ "Suciu" ]
  ]
:   Book*

```

In contrast, we can use the `empty` operator to find all books that have *no* author whose name is Buneman:

```

for b in bib0/book do
  where empty(for a in b/author do
    where a/data() = "Buneman" do
      a) do
    b

```

```

==> book [
  title [ "XML Query" ],
  year  [ 2001 ],
  author [ "Fernandez" ],
  author [ "Suciu" ]
]
: Book*

```

The `empty` expression checks that its argument is the empty sequence `()`.

We can also use the `empty` operator to find all books where all the authors are Buneman, by checking that there are no authors that are not Buneman:

```

for b in bib0/book do
  where empty(for a in b/author do
    where a/data() <> "Buneman" do
      a) do
    b
==> ()
: Book*

```

There are no such books, so the result is the empty sequence. Appropriate use of `empty` (possibly combined with `not`) can express universally or existentially quantified expressions.

Here is a good place to introduce the `let` expression, which binds a local variable to a value. Introducing local variables may improve readability. For example, the following expression is exactly equivalent to the previous one.

```

for b in bib0/book do
  let nonbunemans = (for a in b/author do
    where a/data() <> "Buneman" do
      a) do
    where empty(nonbunemans) do
      b

```

Local variables can also be used to avoid repetition when the same subexpression appears more than once in a query.

Later we will introduce `match` expressions, and we will see how to define `empty` using `match` in Section 6.

2.6 Join

Another common operation is to *join* values from one or more documents. To illustrate joins, we give a second data source that defines book reviews:

```

type Reviews =
  reviews [
    book [
      title [ String ],

```

```

        review [ String ]
    ]*
]
let review0 : Reviews =
  reviews [
    book [
      title [ "XML Query" ],
      review [ "A darn fine book." ]
    ],
    book [
      title [ "Data on the Web" ],
      review [ "This is great!" ]
    ]
  ]
]

```

The `Reviews` type contains one `reviews` element, which contains zero or more `book` elements; each `book` contains a `title` and a `review`.

We can use nested `for` loops to join the two sources `review0` and `bib0` on `title` values. The result combines the `title`, `authors`, and `reviews` for each `book`.

```

    for b in bib0/book do
      for r in review0/book do
        where b/title/data() = r/title/data() do
          book [ b/title, b/author, r/review ]
==>
book [
  title [ "Data on the Web" ],
  author [ "Abiteboul" ],
  author [ "Buneman" ],
  author [ "Suciu" ]
  review [ "A darn fine book." ]
],
book [
  title [ "XML Query" ],
  author [ "Fernandez" ],
  author [ "Suciu" ]
  review [ "This is great!" ]
]

: book [
  title [ String ],
  author [ String ]+
  review [ String ]
]*

```

Note that the outer-most `for` expression determines the order of the result. Readers familiar with optimization of relational join queries know that relational

joins commute, i.e., they can be evaluated in any order. This is not true for the XML algebra: changing the order of the first two `for` expressions would produce different output. In Section 8, we discuss extending the algebra to support unordered forests, which would permit commutable joins.

2.7 Restructuring

Often it is useful to regroup elements in an XML document. For example, each book element in `bib0` groups one title with multiple authors. This expression regroups each author with the titles of his/her publications.

```

for a in distinct(bib0/book/author) do
  biblio [
    a,
    for b in bib0/book do
      for a2 in b/author do
        where a/data() = a2/data() do
          b/title
      ]
  ]
==> biblio [
  author [ "Abiteboul" ],
  title [ "Data on the Web" ]
],
biblio [
  author [ "Buneman" ],
  title [ "Data on the Web" ]
],
biblio [
  author [ "Suciu" ],
  title [ "Data on the Web" ],
  title [ "XML Query" ]
],
biblio [
  author [ "Fernandez" ],
  title [ "XML Query" ]
]
: biblio [
  author [ String ],
  title [ String ]*
]*

```

Readers may recognize this expression as a self-join of books on authors. The expression `distinct(bib0/book/author)` produces a forest of author elements with no duplicates. The outer `for` expression binds `a` to each author element, and the inner `for` expression selects the title of each book that has some author equal to `a`.

Here `distinct` is an example of a built-in function. It takes a forest of elements and removes duplicates.

The type of the result expression may seem surprising: each `biblio` element may contain *zero* or more `title` elements, even though in `bib0`, every `author` co-occurs with a `title`. Recognizing such a constraint is outside the scope of the type system, so the resulting type is not as precise as we might like.

2.8 Aggregation

The algebra has five built-in aggregation functions: `avg`, `count`, `max`, `min` and `sum`. This expression selects books that have more than two authors:

```
for b in bib0/book do
  where count(b/author) > 2 do
    b
==> book [
  title [ "Data on the Web" ],
  year  [ 1999 ],
  author [ "Abiteboul" ],
  author [ "Buneman" ],
  author [ "Suciu" ]
]
: Book*
```

All the aggregation functions take a forest with repetition type and return an integer value; `count` returns the number of elements in the forest.

2.9 Functions

Functions can make queries more modular and concise. Recall that we used the following query to find all books that do not have “Buneman” as an author.

```
for b in bib0/book do
  where empty(for a in b/author do
    where a/data() = "Buneman" do
      a) do
    b
==> book [
  title [ "XML Query" ],
  year  [ 2001 ],
  author [ "Fernandez" ],
  author [ "Suciu" ]
]
: Book*
```

A different way to formulate this query is to first define a function that takes a string `s` and a book `b` as arguments, and returns true if book `b` does not have an author with name `s`.

```

fun notauthor (s : String; b : Book) : Boolean =
  empty(for a in b/author do
    where a/data() = s do
      a)

```

The query can then be re-expressed as follows.

```

for b in bib0/book do
  where notauthor("Buneman"; b) do
    b
==> book [
  title [ "XML Query" ],
  year  [ 2001 ],
  author [ "Fernandez" ],
  author [ "Suciu" ]
]
: Book*

```

We use semicolon rather than comma to separate function arguments, since comma is used to concatenate forests.

Note that a function declaration includes the types of all its arguments and the type of its result. This is necessary for the type system to guarantee that applications of functions are type correct.

In general, any number of functions may be declared at the top-level. The order of function declarations does not matter, and each function may refer to any other function. Among other things, this allows functions to be recursive (or mutually recursive), which supports structural recursion, the subject of the next section.

2.10 Structural Recursion

XML documents can be recursive in structure, for example, it is possible to define a **part** element that directly or indirectly contains other **part** elements. In the algebra, we use recursive types to define documents with a recursive structure, and we use recursive functions to process such documents. (We can also use mutual recursion for more complex recursive structures.)

For instance, here is a recursive type defining a part hierarchy.

```

type Part =
  Basic | Composite
type Basic =
  basic [
    cost [ Integer ]
  ]
type Composite =
  composite [
    assembly_cost [ Integer ],

```

```

    subparts [ Part+ ]
  ]

```

And here is some sample data.

```

let part0 : Part =
  composite [
    assembly_cost [ 12 ],
    subparts [
      composite [
        assembly_cost [ 22 ],
        subparts [
          basic [ cost [ 33 ] ]
        ]
      ],
      basic [ cost [ 7 ] ]
    ]
  ]

```

Here vertical bar (|) is used to indicate a choice between types: each part is either basic (no subparts), and has a cost, or is composite, and includes an assembly cost and subparts.

We might want to translate to a second form, where every part has a total cost and a list of subparts (for a basic part, the list of subparts is empty).

```

type Part2 =
  part [
    total_cost [ Integer ],
    subparts [ Part2* ]
  ]

```

Here is a recursive function that performs the desired transformation. It uses a new construct, the *match* expression.

```

fun convert(p : Part) : Part2 =
  match p
  case b : Basic do
    part[
      total_cost[ b/cost/data() ],
      subparts[]
    ]
  case c : Composite do
    let s = (for q in children(c/subparts) do convert(q)) in
    part[
      total_cost[
        c/assembly_cost/data() + sum(s/total_cost/data())
      ],
      subparts[ s ]
    ]

```

```

    ]
  else error()

```

Each branch of the match expression is labeled with a type, **Basic** or **Composite**, and with a corresponding variable, **b** or **c**. The evaluator checks the type of the value of **p** at *run-time*, and evaluates the corresponding branch. If the first branch is taken then **b** is bound to the value of **p**, and the branch returns a new part with total cost the same as the cost of **b**, and with no subparts. If the second branch is taken then **c** is bound to the value of **p**. The function is recursively applied to each of the subparts of **c**, giving a list of new subparts **s**. The branch returns a new part with total cost computed by adding the assembly cost of **c** to the sum of the total cost of each subpart in **s**, and with subparts **s**.

One might wonder why **b** and **c** are required, since they have the same value as **p**. The reason why is that **p**, **b**, and **c** have different types.

```

p : Part
b : Basic
c : Composite

```

The types of **b** and **c** are more precise than the type of **p**, because which branch is taken depends upon the type of value in **p**.

Applying the query to the given data gives the following result.

```

convert(part0)
==> part [
  total_cost [ 74 ],
  subparts [
    part [
      total_cost [ 55 ],
      subparts [
        part [
          total_cost [ 33 ],
          subparts []
        ]
      ]
    ],
    part [
      total_cost [ 7 ],
      subparts []
    ]
  ]
]
: Part2

```

Of course, a match expression may be used in any query, not just in a recursive one.

2.11 Processing any well-formed document

Recursive types allow us to define a type that matches any well-formed XML document. This type is called `UrTree`:

```
type UrTree = UrScalar | ~[UrType]
type UrType = UrTree*
```

Here `UrScalar` is a built-in scalar type. It stands for the most general scalar type, and all other scalar types (like `Integer` or `String`) are subtypes of it. The tilde (`~`) is used to indicate a wild-card type. In general, `~[t]` indicates the type of elements that may have any element name, but must have children of type `t`. So an `UrTree` is either an `UrScalar` or a wildcard element with zero or more children, each of which is itself an `UrTree`. In other words, any single element or scalar has type `UrTree`.

Types analogous to `UrType` and `UrScalar` appear in XML Schema. The use of tilde is a significant extension to XML Schema, because XML Schema has no type corresponding to `~[t]`, where `t` is some type other than `UrType`. It is not clear that this extension is necessary, since the more restrictive expressiveness of XML Schema wildcards may be adequate.

In particular, our earlier data also has type `UrTree`.

```
book0 : UrTree
==> book [
  title [ "Data on the Web" ],
  year  [ 1999 ],
  author [ "Abiteboul" ],
  author [ "Buneman" ],
  author [ "Suciu" ]
]
: UrTree
```

A specific type can be indicated for any expression in the query language, by writing a colon and the type after the expression.

As an example, we define a recursive function that converts any XML data into HTML. We first give a simplified definition of HTML.

```
type HTML =
  ( UrScalar
  | b [ HTML ]
  | ul [ (li [ HTML ])* ]
  )*
```

An HTML body consists of a sequence of zero or more items, each of which is either: a scalar; or a `b` element (boldface) with HTML content; or a `ul` element (unordered list), where the children are `li` elements (list item), each of which has HTML content.

Now, here is the function that performs the conversion.

```

fun html_of_xml( t : UrTree ) : HTML =
  match t
  case s : UrScalar do
    s
  case e : ~[UrType] do
    b [ name(e) ],
    ul [ for c in children(e) do li [ html_of_xml(c) ] ]
  else error()

```

The `case` expression checks whether the value of `t` is data or an element, and evaluates the corresponding branch. If the first branch is taken, then `s` is bound to the value of `t`, which must be a scalar, and the branch returns the scalar. If the second branch is taken, then `e` is bound to the value of `t`, which must be an element. The branch returns the name of the element in boldface, followed by a list containing one item for each child of the element. The function is recursively applied to get the content of each list item.

Applying the query to the book element above gives the following result.

```

      html_of_xml(book0)
==> b [ "book" ],
      ul [
        li [ b [ "title" ], ul [ li [ "Data on the Web" ] ] ],
        li [ b [ "year" ], ul [ li [ 1999 ] ] ],
        li [ b [ "author" ], ul [ li [ "Abiteboul" ] ] ],
        li [ b [ "author" ], ul [ li [ "Buneman" ] ] ],
        li [ b [ "author" ], ul [ li [ "Suciu" ] ] ]
      ]
: HTML

```

2.12 Top-level Queries

A query consists of a sequence of top-level expressions, or *query items*, where each query item is either a type declaration, a function declaration, a global variable declaration, or a query expression. The order of query items is immaterial; all type, function, and global variable declarations may be mutually recursive.

A query can be evaluated by the query interpreter. Each query expression is evaluated in the environment specified by all of the declarations. (Typically, all of the declarations will precede all of the query expressions, but this is not required.) We have already seen examples of type, function, and global variable declarations. An example of a query expression is:

```

query html_of_xml(book0)

```

To transform any expression into a top-level query, we simply precede the expression by the `query` keyword.

3 Projection and iteration

This section describes key aspects of projection and iteration.

3.1 Relating projection to iteration

The previous examples use the `/` operator liberally, but in fact we use `/` as a convenient abbreviation for expressions built from lower-level operators: `for` expressions, the `children` function, and `match` expressions.

For example, the expression:

```
book0/author
```

is equivalent to the expression:

```
for c in children(book0) do
  match c
    case a : author[UrType] do a
    else ()
```

Here the `children` function returns a forest consisting of the children of the element `book0`, namely, a title element, a year element, and three author elements (the order is preserved). The `for` expression binds the variable `v` successively to each of these elements. Then the `match` expression selects a branch based on the value of `v`. If it is an `author` element then the first branch is evaluated, otherwise the second branch. If the first branch is evaluated, the variable `a` is bound to the same value as `x`, then the branch returns the value of `a`. The type of `a` is `author[String]`, which is the the intersection of the type of `c` and the type `author[UrType]`. If the second branch is evaluated, then the branch returns `()`, the empty sequence.

To compose several expressions using `/`, we again use `for` expressions. For example, the expression:

```
bib0/book/author
```

is equivalent to the expression:

```
for c in children(bib0) do
  match c
    case b : book[UrType] do
      for d in children(b) do
        match d
          case a : author[UrType] do a
          else ()
        else ()
    else ()
```

The `for` expression iterates over all `book` elements in `bib0` and binds the variable `b` to each such element. For each element bound to `b`, the inner expression returns all the `author` elements in `b`, and the resulting forests are concatenated together in order.

In general, an expression of the form `e / a` is converted to the form

```

for v1 in e do
  for v2 in children(v1) do
    match v2
      case v3 : a[UrType] do v3
      else ()

```

where e is an expression, a is an element name, and v_1 , v_2 , and v_3 are fresh variables (ones that do not appear in the expression being converted).

According to this rule, the expression `bib0/book` translates to

```

for v1 in bib0 do
  for v2 in children(v1) do
    match v2
      case v3 : book[UrType] do v3
      else ()

```

In Section 3.3 we discuss laws which allow us to simplify this to the previous expression

```

for v2 in children(bib0) do
  match v2
    case v3 : book[UrType] do v3
    else ()

```

Similarly, the expression `bib0/book/author` translates to

```

for v4 in (for v2 in children(bib0) do
  match v2
    case v3 : book[UrType] do v3
    else ()) do
  for v5 in children(v4) do
    match v5
      case v6 : author[UrType] do v6
      else ()

```

Again, the laws will allow us to simplify this to the previous expression

```

for v2 in children(bib0) do
  match v2
    case v3 : book[UrType] do
      for v5 in children(v3) do
        match v5
          case v6 : author[UrType] do v6
          else ()
      else ()

```

These examples illustrate an important feature of the algebra: high-level operators may be defined in terms of low-level operators, and the low-level operators may be subject to algebraic laws that can be used to further simplify the expression.

3.2 Typing iteration

The typing of `for` loops is rather subtle. We give an intuitive explanation here, and cover the detailed typing rules in Section 7.

A *unit* type is either an element type $a[t]$, a wildcard type $\sim[t]$, or a scalar type s . A `for` loop

```
for v in e1 do e2
```

is typed as follows. First, one finds the type of expression e_1 . Next, for each unit type in this type one assumes the variable v has the unit type and one types the body e_2 . Note that this means we may type the body of e_2 several times, once for each unit type in the type of e_1 . Finally, the types of the body e_2 are combined, according to how the types were combined in e_1 . That is, if the type of e_1 is formed with sequencing, then sequencing is used to combine the types of e_2 , and similarly for choice or repetition.

For example, consider the following expression, which selects all `author` elements from a book.

```
for c in children(book0) do
  match c
    case a : author do a
    else ()
```

The type of `children(book0)` is

```
title[String], year[Integer], author[String]+
```

This is composed of three unit types, and so the body is typed three times.

```
assuming c has type title[String] the body has type ()
"                year[Integer]    "                ()
"                author[String]   "                author[String]
```

The three result types are then combined in the same way the original unit types were, using sequencing and iteration. This yields

```
(), (), author[String]+
```

as the type of the iteration, and simplifying yields

```
author[String]+
```

as the final type.

As a second example, consider the following expression, which selects all `title` and `author` elements from a book, and renames them.

```
for c in children(book0) do
  match c
    case t : title[String] do titl [ t/data() ]
    case y : year[Integer] do ()
    case a : author[String] do auth [ a/data() ]
    else error()
```

Again, the type of `children(book0)` is

```
title[String], year[Integer], author[String]+
```

This is composed of three unit types, and so the body is typed three times.

```
assuming c has type title[String] the body has type titl[String]
"                year[Integer]    "                ()
"                author[String]   "                auth[String]
```

The three result types are then combined in the same way the original unit types were, using sequencing and iteration. This yields

```
titl[String], (), auth[String]+
```

as the type of the iteration, and simplifying yields

```
titl[String], auth[String]+
```

as the final type. Note that the title occurs just once and the author occurs one or more times, as one would expect.

As a third example, consider the following expression, which selects all basic parts from a sequence of parts.

```
for p in children(part0/subparts) do
  match p
  case b : Basic      do b
  case c : Composite do ()
  else error()
```

The type of `children(part0/subparts)` is

```
(Basic | Composite)+
```

This is composed of two unit types, and so the body is typed two times.

```
assuming p has type Basic the body has type Basic
"                Composite "                ()
```

The two result types are then combined in the same way the original unit types were, using sequencing and iteration. This yields

```
(Basic | ())+
```

as the type of the iteration, and simplifying yields

```
Basic*
```

as the final type. Note that although the original type involves repetition one or more times, the final result is a repetition zero or more times. This is what one would expect, since if all the parts are composite the final result will be an empty sequence.

In this way, we see that `for` loops can be combined with `match` expressions to select and rename elements from a sequence, and that the result is given a sensible type.

In order for this approach to typing to be sensible, it is necessary that the unit types can be uniquely identified. However, the type system given here satisfies the following law.

$$a[t_1 \mid t_2] = a[t_1] \mid a[t_2]$$

This has one unit type on the left, but two distinct unit types on the right, and so might cause trouble. Fortunately, our type system inherits an additional restriction from XML Schema: we insist that the regular expressions can be recognized by a top-down deterministic automaton. In that case, the regular expression must have the form on the left, the form on the right is outlawed because it requires a non-deterministic recognizer. With this additional restriction, there is no problem.

The method of translating projection to iteration described in the previous section combined with the typing rules given here yield optimal types for projections, in the following sense. Say that variable x has type t , and the projection x / a has type t' . The type assignment is *sound* if for every value of type t , the value of x / a has type t' . The type assignment is *complete* if for every value y of type t' there is a value x of type t such that $x / a = y$. In symbols, we can see that these conditions are complementary.

$$\begin{array}{ll} \text{sound:} & \forall x \in t. \exists y \in t'. x / a = y \\ \text{complete:} & \forall y \in t'. \exists x \in t. x / a = y \end{array}$$

Any sensible type system must be sound, but it is rare for a type system to be complete. But, remarkably, the type assignment given by the above approach is both sound and complete.

3.3 Monad laws

Investigating aspects of homological algebra in the 1950s, category theorists uncovered the concept of a *monad*, which among other things generalizes set, bag, and list types. Investigating programming languages based on lists in the 1970s, functional programmers adapted from set theory the notion of a *comprehension*, which expresses iteration over set, bag, and list types. In the early 1990s, a precise connection between monads and comprehension notation was uncovered by Wadler [32–34], who was inspired by Moggi’s work applying monads to describe features of programming languages [26, 27]. As the decade progressed this was applied by researchers at the University of Pennsylvania to database languages [6], particularly nested relational algebra (NRA) [7, 23, 22] and the Kleisli system [36].

The iteration construct of the algebra corresponds to the structure of a monad. The correspondence is close, but not exact. Each monad is based on a unary type constructor, such as $Set(t)$ or $List(t)$, representing a homogenous

set or list where all elements are of type t . In contrast, here we have more complex and heterogenous types, such as a forest consisting of a title, a year, and a sequence of one or more authors. Also, one important component of a monad is the unit operator, which converts an element to a set or list. If x has type t , then $\{x\}$ is a unit set of type $Set(t)$ or $[x]$ is a unit list of type $List(t)$. In contrast, here we simply write, say, `author["Buneman"]`, which stands for both a tree and for the unit forest containing that tree.

One can define comprehensions in terms of iteration:

$$\begin{aligned} [e_0 \mid x_1 <- e_1] &= \text{for } x_1 \text{ in } e_1 \text{ do } e_0 \\ [e_0 \mid x_1 <- e_1, x_2 <- e_2] &= \text{for } x_1 \text{ in } e_1 \text{ do for } x_2 \text{ in } e_2 \text{ do } e_0 \\ &\dots \end{aligned}$$

Conversely, one can define iterations in terms of comprehension:

$$\text{for } x \text{ in } e_1 \text{ do } e_2 = [y \mid x <- e_1, y <- e_2]$$

Here y is a fresh variable name.

Monads satisfy three laws, and three corresponding laws are satisfied by the iteration notation given here.

First, iteration over a unit forest can be replaced by substitution. This is called the *left unit law*.

$$\text{for } v \text{ in } e_1 \text{ do } e_2 = e_2\{v := e_1\}$$

provided that e_1 is a unit type (e.g., is an element or a scalar constant). We write $e_1\{v := e_2\}$ to denote the result of taking expression e_1 and replacing occurrences of the variable v by the expression e_2 . For example

$$\text{for } v \text{ in } \text{author}["\text{Buneman}"] \text{ do } \text{auth}[v/\text{data}()] = \text{auth}["\text{Buneman}"]$$

Second, an iteration that returns the iteration variable is equivalent to the identity. This is called the *right unit law*.

$$\text{for } v \text{ in } e \text{ do } v = e$$

For example

$$\text{for } v \text{ in } \text{book0} \text{ do } v = \text{book0}$$

An important feature of the type system described here is that the left side of the above equation always has the same type as the right side. (This was not true for an earlier version of the type system [18].)

Third, there are two ways of writing an iteration over an iteration, both of which are equivalent. This is called the *associative law*.

$$\begin{aligned} \text{for } v_2 \text{ in } (\text{for } v_1 \text{ in } e_1 \text{ do } e_2) \text{ do } e_3 \\ = \text{for } v_1 \text{ in } e_1 \text{ do } (\text{for } v_2 \text{ in } e_2 \text{ do } e_3) \end{aligned}$$

For example, a projection over a forest includes an implicit iteration, so $e / a = \text{for } v \text{ in } e \text{ do } v / a$. Say we define a forest of bibliographies, `bib1 = bib0`,

bib0. Then bib1/book/author is equivalent to the first expression below, which in turn is equivalent to the second.

```
for b in (for a in bib1 do a/book) do b/author
= for a in bib1 do (for b in a/book do b/author
```

With nested relational algebra, the monad laws play a key role in optimizing queries. For instance, they are exploited extensively in the Kleisli system for biomedical data, developed by Limsoon Wong and others at the University of Pennsylvania and Kent Ridge Digital Labs, and now sold commercially [36]. Similarly, the monad laws can also be exploited for optimization in this context.

For example, if `b` is a book, the following finds all authors of the book that are not Buneman:

```
for a in b do
  where a/data() != Buneman do
    a
```

If `l` is a list of authors, the following renames all `author` elements to `auth` elements:

```
for a' in l do
  auth[ a'/data() ]
```

Combining these, we select all authors that are not Buneman, and rename the elements:

```
for a' in (for a in b do
  where a/data() != Buneman do
  a) do
  auth[ a'/data() ]
```

Applying the associative law for a monad, we get:

```
for a in b do
  for a' in (where a/data() != Buneman do a) do
    auth[ a'/data() ]
```

Expanding the `where` clause to a conditional, we get:

```
for a in b do
  for a' in (if a/data() != Buneman then a else ()) do
    auth[ a'/data() ]
```

Applying a standard for loops over conditionals gives:

```
for a in b do
  if a/data() != Buneman then
    for a' in a do
      auth[ a'/data() ]
  else ()
```

Applying the left unit law for a monad, we get:

```
for a in b do
  if a/data() != Buneman then
    auth[ a/data() ]
  else ()
```

And replacing the conditional by a `where` clause, we get:

```
for a in b do
  where a/data() != Buneman do
    auth[ a/data() ]
```

Thus, simple manipulations, including the monad laws, fuse the two loops.

Section 3.1 ended with two examples of simplification. Returning to these, we can now see that the simplifications are achieved by application of the left unit and associative monad laws.

4 Expressions

Figure 1 contains the grammar for the algebra, i.e., the convenient concrete syntax in which a user may write a query. A few of these expressions can be rewritten as other expressions in a smaller *core* algebra; such derived expressions are labeled with “*”. We define the algebra’s typing rules on the smaller core algebra. In Section 6, we give the laws that relate a user expression with its equivalent expression in the core algebra. Typing rules for the core algebra are defined in Section 7.

We have seen examples of most of the expressions, so we will only point out a few details here. We include only two operators, `+` and `=`, and one aggregate function `sum` in the formal syntax, adding others is straightforward.

A query consists of a sequence of *query items*, where each query item is either a type declaration, a function declaration, a global variable declaration, or a query expression. The order of query items is immaterial; all type, function, and global variable declarations may be mutually recursive. Each query expression is evaluated in the environment specified by all of the declarations. (Typically, all of the declarations will precede all of the query expressions, but this is not required.)

We define a subset of expressions that correspond to *data values*. An expression is a data value if it consists only of scalar constant, element, sequence, and empty sequence expressions.

5 Types

Figure 2 contains the grammar for the algebra’s type system. We have already seen many examples of types. Here, we point out some details.

tag	a	
function	f	
variable	v	
integer	$c_{\text{int}} ::= \dots \mid -1 \mid 0 \mid 1 \mid \dots$	
string	$c_{\text{str}} ::= "" \mid "a" \mid "b" \mid \dots \mid "aa" \mid \dots$	
boolean	$c_{\text{bool}} ::= \text{false} \mid \text{true}$	
constant	$c ::= c_{\text{int}} \mid c_{\text{str}} \mid c_{\text{bool}}$	
expression	$e ::= c$	scalar constant
	v	variable
	$a[e]$	element
	$\sim e[e]$	computed element
	e, e	sequence
	$()$	empty sequence
	if e then e else e	conditional
	let $v = e$ do e	local binding
	for v in e do e	iteration
	match e_0	match
	case $v_1 : t_1$ do e_1	
	...	
	case $v_n : t_n$ do e_n	
	else e_{n+1}	
	$f(e; \dots; e)$	function application
	error $()$	error
	$e + e$	plus
	$e = e$	equal
	sum (e)	aggregation
	children (e)	children
	name (e)	element name
	e / a	element projection *
	$e / \text{data}()$	scalar projection *
	where e then e	conditional *
	empty (e)	empty test *
query item	$q ::= \text{type } x = t$	type declaration
	fun $f(v:t; \dots; v:t) : t = e$	function declaration
	let $v : t = e$	global declaration
	query e	query expression
data	$d ::= c$	scalar constant
	$a[d]$	element
	d, d	sequence
	$()$	empty sequence

Fig. 1. Expressions

Our algebra uses a simple type system that captures the essence of XML Schema [40]. The type system is close to that used in XDuCE [20].

In the type system of Figure 2, a scalar type may be a `UrScalar`, `Boolean`, `Integer`, or `String`. In XML Schema, a scalar type is defined by one of fourteen

element name	a	
type name	x	
scalar type	$s ::=$	Integer
		String
		Boolean
		UrScalar
type	$t ::=$	x type name
		s scalar type
		$a[t]$ element
		$\sim[t]$ wildcard
		t, t sequence
		$t t$ choice
		t^* repetition
		$()$ empty sequence
		\emptyset empty choice
unit type	$u ::=$	$a[t]$ element
		$\sim[t]$ wildcard
		s scalar type

Fig. 2. Types

primitive datatypes and a list of facets. A type hierarchy is induced between scalar types by containment of facets. The algebra's type system can be generalized to support these types without much increase in its complexity. We added **UrScalar**, because XML Schema does not support a most general scalar type.

A type is either: a type variable; a scalar type; an element type with element name a and content type t ; a *wildcard* type with an unknown element name and content type t ; a sequence of two types, a choice of two types; a repetition type; the empty sequence type; or the empty choice type.

The algebra's external type system, that is, the type definitions associated with input and output documents, is XML Schema. The internal types are in some ways more expressive than XML Schema, for example, XML Schema has no type corresponding to **Integer*** (which is required as the type of the argument to an aggregation operator like **sum** or **min** or **max**), or corresponding to $\sim[t]$ where t is some type other than **UrTree***. In general, mapping XML Schema types into internal types will not lose information, however, mapping internal types into XML Schema may lose information.

5.1 Relating values to types

Recall that *data* is the subset of expressions that consists only of scalar constant, element, sequence, and empty sequence expressions. We write $\vdash d : t$ if data d has type t . The following type rules define this relation.

$$\frac{}{\vdash c_{\text{int}} : \text{Integer}}$$

$$\frac{}{\vdash c_{\text{str}} : \text{String}}$$

$$\frac{}{\vdash c_{\text{bool}} : \text{Boolean}}$$

$$\frac{}{\vdash c : \text{UrScalar}}$$

$$\frac{\vdash d : t}{\vdash a[d] : a[t]}$$

$$\frac{\vdash d : t}{\vdash a[d] : \sim[t]}$$

$$\frac{\vdash d_1 : t_1 \quad \vdash d_2 : t_2}{\vdash d_1, d_2 : t_1, t_2}$$

$$\frac{}{\vdash () : ()}$$

$$\frac{\vdash d : t_1}{\vdash d : t_1 \mid t_2}$$

$$\frac{\vdash d : t_2}{\vdash d : t_1 \mid t_2}$$

$$\frac{\vdash d_1 : t \quad \vdash d_2 : t^*}{\vdash (d_1, d_2) : t^*}$$

$$\frac{}{\vdash () : t^*}$$

We write $t_1 <: t_2$ if for every data d such that $\vdash d : t_1$ it is also the case that $\vdash d : t_2$, that is t_1 is a subtype of t_2 . It is easy to see that $<:$ is a partial order, that is it is reflexive, $t <: t$, and it is transitive, if $t_1 <: t_2$ and $t_2 <: t_3$ then $t_1 <: t_3$. Here are some of the inequalities that hold.

\emptyset	$<: t$
t	$<: \text{UrType}$
t_1	$<: t_1 \mid t_2$
t_2	$<: t_1 \mid t_2$
Integer	$<: \text{UrScalar}$
String	$<: \text{UrScalar}$
Boolean	$<: \text{UrScalar}$
$a[t]$	$<: \sim[t]$

Further, if $t <: t'$ then

$$\begin{aligned} a[t] &<: a[t'] \\ t* &<: t'* \end{aligned}$$

And if $t_1 <: t'_1$ and $t_2 <: t'_2$ then

$$\begin{aligned} t_1, t_2 &<: t'_1, t'_2 \\ t_1 \mid t_2 &<: t'_1 \mid t'_2 \end{aligned}$$

We write $t_1 = t_2$ if $t_1 <: t_2$ and $t_2 <: t_1$. Here are some of the equations that hold.

$$\begin{aligned} \text{UrScalar} &= \text{Integer} \mid \text{String} \mid \text{Boolean} \\ (t_1, t_2), t_3 &= t_1, (t_2, t_3) \\ t, () &= t \\ (), t &= t \\ t_1 \mid t_2 &= t_2 \mid t_1 \\ (t_1 \mid t_2) \mid t_3 &= t_1 \mid (t_2 \mid t_3) \\ t \mid \emptyset &= t \\ \emptyset \mid t &= t \\ t_1, (t_2 \mid t_3) &= (t_1, t_2) \mid (t_1, t_3) \\ (t_1 \mid t_2), t_3 &= (t_1, t_3) \mid (t_2, t_3) \\ t, \emptyset &= \emptyset \\ \emptyset, t &= \emptyset \\ t* &= () \mid t, t* \end{aligned}$$

We also have that $t_1 <: t_2$ if and only if $t_1 \mid t_2 = t_2$.

We define $t?$ and $t+$ as abbreviations, by the following equivalences.

$$\begin{aligned} t? &= () \mid t \\ t+ &= t, t* \end{aligned}$$

We define the *intersection* $t_1 \wedge t_2$ of two types t_1 and t_2 to be the largest type t that is smaller than both t_1 and t_2 . That is, $t = t_1 \wedge t_2$ if $t <: t_1$ and $t <: t_2$ and if for any t' such that $t' <: t_1$ and $t' <: t_2$ we have $t' <: t$.

6 Equivalences and Optimization

6.1 Equivalences

Here are the laws that define derived expressions (those labeled with * in Figure 1) in terms of other expressions.

$$\begin{aligned} e / a & \\ &= \text{for } v_1 \text{ in } e \text{ do} & (1) \\ &\quad \text{for } v_2 \text{ in children}(v_1) \text{ do} \\ &\quad \text{match } v_2 \\ &\quad \quad \text{case } v_3 : a[\text{UrType}] \text{ do } v_3 \\ &\quad \quad \text{else } () \end{aligned}$$

```

e / data()
= for v1 in e do
  for v2 in children(v1) do
    match v2
      case v3 : UrScalar do v3
      else ()

```

(2)

```

where e1 then e2
= if e1 then e2 else ()

```

(3)

```

empty(e)
= match e case v : () do true else false

```

(4)

Law 1 rewrites the element projection expression e / a , as described previously. Law 2 rewrites the scalar projection expression $e / \text{data}()$, similarly. Law 3 rewrites a **where** expression as a conditional, as described previously. Law 4 rewrites an **empty** test using a **match** expression.

6.2 Optimizations

In a relational query engine, algebraic simplifications are often applied by a query optimizer before a physical execution plan is generated; algebraic simplification can often reduce the size of the intermediate results computed by a query interpreter. The purpose of our laws is similar – they eliminate unnecessary **for** or **match** expressions, or they enable other optimizations by reordering or distributing computations. The set of laws given here is suggestive, rather than complete.

Here are some simplification laws.

```

for v in () do e = ()

```

(1)

```

for v in (e1, e2) do e3
= (for v in e1 do e3), (for v in e2 do e3)

```

(2)

```

for v in e1 do e2
= e2{v := e1}, if e : u

```

(3)

```

match a[e0] case v : a[t] do e1 else e2
= e1{v := a[e0]}, if e0 : t' and t' <: t

```

(4)

```

match a'[e0] case v : a[t] do e1 else e2
= e2, if a ≠ a'

```

(5)

```

for v in e do v = e

```

(6)

Laws 1, 2, and 3 simplify iterations. Law 1 rewrites an iteration over the empty sequence as the empty sequence. Law 2 distributes iteration through

sequence: iterating over the sequence e_1, e_2 is equivalent to the sequence of two iterations, one over e_1 and one over e_2 . Law 3 is the left unit law for a monad. If e_1 is a unit type, then e_1 can be substituted for occurrences of v in e_2 . Laws 4 and 5 eliminate trivial `match` expressions. Law 6 is the right unit law for a monad.

The remaining laws commute expressions. Each law actually abbreviates a number of other laws, since the *context variable* E stands for a number of different expressions. The notation $E[e]$ stands for one of the six expressions given with expression e replacing the hole $[]$ that appears in each of the alternatives.

$$\begin{aligned}
 E ::= & \text{if } [] \text{ then } e_1 \text{ else } e_2 \\
 & | \text{let } v = [] \text{ do } e \\
 & | \text{for } v \text{ in } [] \text{ do } e \\
 & | \text{match } [] \\
 & \quad \text{case } v_1 : t_1 \text{ do } e_1 \\
 & \quad \dots \\
 & \quad \text{case } v_n : t_n \text{ do } e_n \\
 & \quad \text{else } e_{n+1}
 \end{aligned}$$

Here are the laws for commuting expressions.

$$\begin{aligned}
 E[\text{if } e_1 \text{ then } e_2 \text{ else } e_3] \\
 = \text{if } e_1 \text{ then } E[e_2] \text{ else } E[e_3]
 \end{aligned} \tag{7}$$

$$\begin{aligned}
 E[\text{let } v = e_1 \text{ do } e_2] \\
 = \text{let } v = e_1 \text{ do } E[e_2]
 \end{aligned} \tag{8}$$

$$\begin{aligned}
 E[\text{for } v \text{ in } e_1 \text{ do } e_2] \\
 = \text{for } v \text{ in } e_1 \text{ do } E[e_2]
 \end{aligned} \tag{9}$$

$$\begin{aligned}
 E[\text{match } e_0 \\
 \quad \text{case } v_1 : t_1 \text{ do } e_1 \\
 \quad \dots \\
 \quad \text{case } v_n : t_n \text{ do } e_n \\
 \quad \text{else } e_{n+1}] \\
 = \text{match } e_0 \\
 \quad \text{case } v_1 : t_1 \text{ do } E[e_1] \\
 \quad \dots \\
 \quad \text{case } v_n : t_n \text{ do } E[e_n] \\
 \quad \text{else } E[e_{n+1}]
 \end{aligned} \tag{10}$$

Each law has the same form. Law 7 commutes conditionals, Law 8 commutes local bindings, Law 9 commutes iterations, and Law 10 commutes match expressions. For instance, one of the expansions of Law 9 is the following, when E is taken to be `for v in $[]$ do e` .

$$\begin{aligned}
 & \text{for } v_2 \text{ in } (\text{for } v_1 \text{ in } e_1 \text{ do } e_2) \text{ do } e_3 \\
 & = \text{for } v_1 \text{ in } e_1 \text{ do } (\text{for } v_2 \text{ in } e_2 \text{ do } e_3)
 \end{aligned}$$

This will be recognized as the associative law for a monad.

7 Type Rules

We explain our type system in the form commonly used in the programming languages community. For a textbook introduction to type systems, see, for example, Mitchell [24, 25].

7.1 Environments

The type rules make use of an environment that specifies the types of variables and functions. The type environment is denoted by Γ , and is composed of a comma-separated list of variable types, $v : t$ or function types, $f : (t_1; \dots; t_n) \rightarrow t$. We retrieve type information from the environment by writing $(v : t) \in \Gamma$ to look up a variable, or by writing $(f : (t_1; \dots; t_n) \rightarrow t) \in \Gamma$ to look up a function.

7.2 Type rules

We write $\Gamma \vdash e : t$ if in environment Γ the expression e has type t . Below are all the rules except those for `for` and `match` expressions, which are discussed in the following subsections.

$$\frac{}{\Gamma \vdash c_{\text{int}} : \text{Integer}}$$

$$\frac{}{\Gamma \vdash c_{\text{str}} : \text{String}}$$

$$\frac{}{\Gamma \vdash c_{\text{bool}} : \text{Boolean}}$$

$$\frac{(v : t) \in \Gamma}{\Gamma \vdash v : t}$$

$$\frac{\Gamma \vdash e : t}{\Gamma \vdash a[e] : a[t]}$$

$$\frac{\Gamma \vdash e_1 : \text{String} \quad \Gamma \vdash e_2 : t}{\Gamma \vdash \sim_{e_1}[e_2] : \sim[t]}$$

$$\frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2}{\Gamma \vdash e_1, e_2 : t_1, t_2}$$

$$\frac{}{\Gamma \vdash () : ()}$$

$$\frac{\Gamma \vdash e_1 : \text{Boolean} \quad \Gamma \vdash e_2 : t_2 \quad \Gamma \vdash e_3 : t_3}{\Gamma \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 : (t_2 | t_3)}$$

$$\frac{\Gamma \vdash e_1 : t_1 \quad \Gamma, v : t_1 \vdash e_2 : t_2}{\Gamma \vdash \text{let } v = e_1 \text{ do } e_2 : t_2}$$

$$\frac{\begin{array}{c} (f : (t_1; \dots; t_n) \rightarrow t) \in \Gamma \\ \Gamma \vdash e_1 : t'_1 \quad t'_1 <: t_1 \\ \dots \\ \Gamma \vdash e_n : t'_n \quad t'_n <: t_n \end{array}}{\Gamma \vdash f(e_1; \dots; e_n) : t}$$

$$\frac{}{\Gamma \vdash \text{error}() : \emptyset}$$

$$\frac{\Gamma \vdash e_1 : \text{Integer} \quad \Gamma \vdash e_2 : \text{Integer}}{\Gamma \vdash e_1 + e_2 : \text{Integer}}$$

$$\frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2}{\Gamma \vdash e_1 = e_2 : \text{Boolean}}$$

$$\frac{\Gamma \vdash e : \text{Integer}^*}{\Gamma \vdash \text{sum } e : \text{Integer}}$$

7.3 Typing for expressions

The type rule for **for** expressions uses the following auxiliary judgement. We write $\Gamma \vdash \text{for } v : t \text{ do } e : t'$ if in environment Γ when the bound variable of an iteration v has type t then the body e of the iteration has type t' . Recall that u is a unit type.

$$\frac{\Gamma, v : u \vdash e : t'}{\Gamma \vdash \text{for } v : u \text{ do } e : t'}$$

$$\frac{}{\Gamma \vdash \text{for } v : () \text{ do } e : ()}$$

$$\frac{\Gamma \vdash \text{for } v : t_1 \text{ do } e : t'_1 \quad \Gamma \vdash \text{for } v : t_2 \text{ do } e : t'_2}{\Gamma \vdash \text{for } v : t_1, t_2 \text{ do } e : t'_1, t'_2}$$

$$\frac{}{\Gamma \vdash \text{for } v : \emptyset \text{ do } e : \emptyset}$$

$$\frac{\Gamma \vdash \text{for } v : t_1 \text{ do } e : t'_1 \quad \Gamma \vdash \text{for } v : t_2 \text{ do } e : t'_2}{\Gamma \vdash \text{for } v : t_1 | t_2 \text{ do } e : t'_1 | t'_2}$$

$$\frac{\Gamma \vdash \text{for } v : t \text{ do } e : t'}{\Gamma \vdash \text{for } v : t^* \text{ do } e : t'^*}$$

Given the above, the type rule for `for` expressions is immediate.

$$\frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash \text{for } v : t_1 \text{ do } e_2 : t_2}{\Gamma \vdash \text{for } v \text{ in } e_1 \text{ do } e_2 : t_2}$$

Note this rule is unusual in that it may type check the body of the iteration many times, once for each unit type that is iterated over.

7.4 Typing match expressions

Due to the rule for iteration, it is possible that the body of an iteration is checked many times. Thus, when a match expression is checked, it is possible that quite a lot is known about the type of the expression being matched, and one can determine that only some of the clauses of the match apply. The definition of `match` uses the auxiliary judgments to check whether a given clause is applicable.

We write $\Gamma \vdash \text{case } v : t \text{ do } e : t'$ if in environment Γ when the bound variable of the case v has type t then the body e of the case has type t' . Note the type of the body is irrelevant if $t = \emptyset$.

$$\frac{t \neq \emptyset \quad \Gamma, v : t \vdash e : t'}{\Gamma \vdash \text{case } v : t \text{ do } e : t'}$$

$$\frac{}{\Gamma \vdash \text{case } v : \emptyset \text{ do } e : \emptyset}$$

We write $\Gamma \vdash t <: t' \text{ else } e : t''$ if in environment Γ when $t <: t'$ does not hold then the body e of the `else` clause has type t'' . Note that the type of the body is irrelevant if $t <: t'$.

$$\frac{t <: t'}{\Gamma \vdash t <: t' \text{ else } e : \emptyset}$$

$$\frac{t \not<: t' \quad \Gamma \vdash e : t''}{\Gamma \vdash t <: t' \text{ else } e : t''}$$

Given the above, it is straightforward to construct the typing rule for a `match` expression. Recall that we write $t \wedge t'$ for the intersection of two types.

$$\frac{\begin{array}{c} \Gamma \vdash e_0 : t_0 \\ \Gamma \vdash \text{case } v_1 : t_0 \wedge t_1 \text{ do } e_1 : t'_1 \\ \dots \\ \Gamma \vdash \text{case } v_n : t_0 \wedge t_n \text{ do } e_n : t'_n \\ \Gamma \vdash t_0 <: t_1 \mid \dots \mid t_n \text{ else } e_{n+1} : t'_{n+1} \end{array}}{\Gamma \vdash \left(\begin{array}{l} \text{match } e_0 \\ \text{case } v_1 : t_1 \text{ do } e_1 \\ \dots \\ \text{case } v_n : t_n \text{ do } e_n \\ \text{else } e_{n+1} \end{array} \right) : t'_1 \mid \dots \mid t'_{n+1}}$$

7.5 Top-level expressions

We write $\Gamma \vdash q$ if in environment Γ the query item q is well-typed.

$$\frac{}{\Gamma \vdash \mathbf{type} \ x = t}$$

$$\frac{\Gamma, v_1 : t_1, \dots, v_n : t_n \vdash e : t' \quad t' <: t}{\Gamma \vdash f(v_1:t_1; \dots; v_n:t_n) : t = e}$$

$$\frac{\Gamma \vdash e : t' \quad t' <: t}{\Gamma \vdash \mathbf{let} \ v : t = e}$$

$$\frac{\Gamma \vdash e : t}{\Gamma \vdash \mathbf{query} \ e}$$

We extract the relevant component of a type environment from a query item q with the function $environment(q)$.

$$\begin{aligned} environment(\mathbf{type} \ x = t) &= () \\ environment(\mathbf{fun} \ f(v_1:t_1; \dots; v_n:t_n) : t) &= f : (t_1; \dots; t_n) \rightarrow t \\ environment(\mathbf{let} \ v : t = e) &= v : t \end{aligned}$$

We write $\vdash q_1 \dots q_n$ if the sequence of query items $q_1 \dots q_n$ is well typed.

$$\frac{\Gamma = environment(q_1), \dots, environment(q_n) \quad \Gamma \vdash q_1 \quad \dots \quad \Gamma \vdash q_n}{\vdash q_1 \dots q_n}$$

8 Discussion

The algebra has several important characteristics: its operators are strongly typed, and they obey laws of equivalence and optimization.

There are many issues to resolve in the completion of the algebra. We enumerate some of these here.

Data Model. Currently, all forests in the data model are ordered. It may be useful to have unordered forests. The **distinct** operator, for example, produces an inherently unordered forest. Unordered forests can benefit from many optimizations for the relational algebra, such as commutable joins.

The data model and algebra do not define a global order on documents. Querying global order is often required in document-oriented queries.

Currently, the algebra does not support reference values, which are defined in the XML Query Data Model. The algebra's type system should be extended to support reference types and the data model operators **ref** and **deref** should be supported.

Type System. As discussed, the algebra's internal type system is closely related to the type system of XDuce. A potentially significant problem is that the algebra's types may lose information when converted into XML Schema types, for example, when a result is serialized into an XML document and XML Schema.

The type system is currently first order: it does not support function types nor higher-order functions. Higher-order functions are useful for specifying, for example, sorting and grouping operators, which take other functions as arguments.

The type system is currently monomorphic: it does not permit the definition of a function over generalized types. Polymorphic functions are useful for factoring equivalent functions, each of which operate on a fixed type. The lack of polymorphism is one of the principal weaknesses of the type system.

Recursion. Currently, the algebra does not guarantee termination of recursive expressions. In order to ensure termination, we might require that a recursive function take one argument that is a singleton element, and any recursive invocation should be on a descendant of that element; since any element has a finite number of descendants, this avoids infinite regress. (Ideally, we should have a simple syntactic rule that enforces this restriction, but we have not yet devised such a rule.)

References

1. Francois Bancilhon, Paris Kanellakis, Claude Delobel. *Building an Object-Oriented Database System*. Morgan Kaufmann, 1990.
2. Catriel Beeri and Yoram Kornatzky. Algebraic Optimization of Object-Oriented Query Languages. *Theoretical Computer Science* 116(1&2):59–94, August 1993.
3. Catriel Beeri and Yariv Tzaban, SAL: An Algebra for Semistructured Data and XML, *International Workshop on the Web and Databases (WebDB'99)*, Philadelphia, Pennsylvania, June 1999.
4. Anne Brüggemann-Klein, Derick Wood, One-Unambiguous Regular Languages, *Information and Computation*, 140(2): 229–253 and 142(2): 182–206, 1998.
5. P. Buneman, M. Fernandez, D. Suciu. UnQL: A query language and algebra for semistructured data based on structural recursion. *VLDB Journal*, to appear.
6. Peter Buneman, Leonid Libkin, Dan Suciu, Val Tannen, and Limsoon Wong. Comprehension Syntax. *SIGMOD Record*, 23:87–96, 1994.
7. Peter Buneman, Shamim Naqvi, Val Tannen, Limsoon Wong. Principles of programming with complex object and collection types. *Theoretical Computer Science* 149(1):3–48, 1995.
8. The Caml Language. <http://caml.inria.fr/>.
9. R. G. Cattell. *The Object Database Standard: ODMG 2.0*. Morgan Kaufmann, 1997.
10. Don Chamberlin, Jonathan Robie, and Daniela Florescu. Quilt: An XML Query Language for Heterogeneous Data Sources. *International Workshop on the Web and Databases (WebDB'2000)*, Dallas, Texas, May 2000.

11. Vassilis Christophides and Sophie Cluet and Jérôme Siméon. On Wrapping Query Languages and Efficient XML Integration. *Proceedings of ACM SIGMOD Conference on Management of Data*, Dallas, Texas, May 2000.
12. S. Cluet, S. Jacqmin and J. Siméon The New YATL: Design and Specifications. *Technical Report*, INRIA, 1999.
13. S. Cluet and G. Moerkotte. Nested queries in object bases. *Workshop on Database Programming Languages*, pages 226–242, New York, August 1993.
14. L. S. Colby. A recursive algebra for nested relations. *Information Systems* 15(5):567–582, 1990.
15. H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi, *Tree Automata Techniques and Applications*, 1997.
16. Hugh Darwen (Contributor) and Chris Date. *Guide to the SQL Standard: A User's Guide to the Standard Database Language SQL* Addison-Wesley, 1997.
17. A. Deutsch, M. Fernandez, D. Florescu, A. Levy, and D. Suciu. A query language for XML. In *International World Wide Web Conference*, 1999. <http://www.research.att.com/~mff/files/final.html>
18. Mary Fernandez, Jerome Simeon, Philip Wadler. An Algebra for XML Query, Draft manuscript, June 2000.
19. Mary Fernandez, Jerome Simeon, Philip Wadler. An Algebra for XML Query, *Foundations of Software Technology and Theoretical Computer Science (FSTTCS 2000)*, New Delhi, December 2000.
20. Haruo Hosoya, Benjamin Pierce, XDuce : A Typed XML Processing Language (Preliminary Report). In *WebDB Workshop 2000*.
21. Haruo Hosoya, Benjamin Pierce, Jerome Vouillon, Regular Expression Types for XML. In *International Conference on Functional Programming (ICFP)*, September 2000.
22. Leonid Libkin and Limsoon Wong. Query languages for bags and aggregate functions. *Journal of Computer and Systems Sciences*, 55(2):241–272, October 1997.
23. Leonid Libkin, Rona Machlin, and Limsoon Wong. A query language for multi-dimensional arrays: Design, implementation, and optimization techniques. *SIGMOD* 1996.
24. John C. Mitchell, *Foundations for Programming Languages*. MIT Press, 1998.
25. John C. Mitchell, Type systems for programming languages, In J. van Leeuwen, editor, *Handbook of Theoretical Computer Science*, Volume B, MIT Press, 1998.
26. E. Moggi, Computational lambda-calculus and monads. In *Symposium on Logic in Computer Science*, Asilomar, California, IEEE, June 1989.
27. E. Moggi, Notions of computation and monads. *Information and Computation*, 93(1), 1991.
28. J. Robie, editor. XQL '99 Proposal, 1999. <http://metalab.unc.edu/xql/xql-proposal.html>.
29. H.-J. Schek and M. H. Scholl. The relational model with relational-valued attributes. *Information Systems* 11(2):137–147, 1986.
30. S. J. Thomas and P. C. Fischer. Nested Relational Structures. In *Advances in Computing Research: The Theory of Databases*, JAI Press, London, 1986.
31. W. Thomas, Chapter 4.8, Tree Automata, In J. van Leeuwen, editor, *Handbook of Theoretical Computer Science*, Volume B, MIT Press, 1998.
32. Philip Wadler. Comprehending monads. *Mathematical Structures in Computer Science*, 2:461-493, 1992.

33. P. Wadler, Monads for functional programming. In M. Broy, editor, *Program Design Calculi*, NATO ASI Series, Springer Verlag, 1993. Also in J. Jeuring and E. Meijer, editors, *Advanced Functional Programming*, LNCS 925, Springer Verlag, 1995.
34. P. Wadler, How to declare an imperative. *ACM Computing Surveys*, 29(3):240–263, September 1997.
35. Philip Wadler. A formal semantics of patterns in XSLT. Markup Technologies, Philadelphia, December 1999.
36. Limsoon Wong. An introduction to the Kleisli query system and a commentary on the influence of functional programming on its implementation. *Journal of Functional Programming*, to appear.
37. World-Wide Web Consortium, XML Path Language (XPath): Version 1.0. November, 1999. <http://www.w3.org/TR/xpath.html>
38. World-Wide Web Consortium, XML Query: Requirements, Working Draft. August 2000. <http://www.w3.org/TR/xmlquery-req>
39. World-Wide Web Consortium, XML Query: Data Model, Working Draft. May 2000. <http://www.w3.org/TR/query-datamodel/>
40. World-Wide Web Consortium, XML Schema Part 1: Structures, Working Draft. April 2000. <http://www.w3.org/TR/xmlschema-1>
41. World-Wide Web Consortium, XML Schema Part 2: Datatypes, Working Draft, April 2000. <http://www.w3.org/TR/xmlschema-2>.
42. World-Wide Web Consortium, XSL Transformations (XSLT), Version 1.0. W3C Recommendation, November 1999. <http://www.w3.org/TR/xslt>.